

INAUGURAL-DISSERTATION ZUR

ERLANGUNG DES AKADEMISCHEN GRADES

doctor rerum naturalium (Dr. rer. nat.)

IN THEORETISCHER PHYSIK

der Mathematisch-Naturwissenschaftlichen Fakultät der Universität zu Köln

Network inference and response prediction in biological systems

vorgelegt von Niklas Bonacker aus Köln

Gutachter: Professor Dr. Johannes Berg Privatdozent Dr. Rochus Klesse

Erklärung zur Dissertation

gemäß der Promotionsordnung vom 02. Februar 2006 mit den Änderungsordnungen vom 10. Mai 2012, 16. Januar 2013 und 21. Februar 2014.

Ich versichere, dass ich die von mir vorgelegte Dissertation selbständig angefertigt, die benutzten Quellen und Hilfsmittel vollständig angegeben und die Stellen der Arbeit – einschließlich Tabellen, Karten und Abbildungen – , die anderen Werken im Wortlaut oder dem Sinn nach entnommen sind, in jedem Einzelfall als Entlehnung kenntlich gemacht habe; dass diese Dissertation noch keiner anderen Fakultät oder Universität zur Prüfung vorgelegen hat; dass sie – abgesehen von unten angegebenen Teilpublikationen – noch nicht veröffentlicht worden ist, sowie, dass ich eine solche Veröffentlichung vor Abschluss des Promotionsverfahrens nicht vornehmen werde. Die Bestimmungen der Promotionsordnung sind mir bekannt. Die von mir vorgelegte Dissertation ist von Professor Dr. Johannes Berg betreut worden.

Niklas Bonacker Köln, den 15. November 2021

Kurzzusammenfassung

Anomale biologische Informationsverarbeitung spielt eine entscheidende Rolle bei der Entstehung und Ausbreitung von Krebs. Die mathematische Modellierung und Inferenz von Signalwegen ist eine zentrale Herausforderung auf dem Gebiet der Krebsforschung [1].

Wir schlagen ein stochastisches Modell der Genregulation vor und lösen unser Modell im Rahmen einer Gaußschen Theorie. Wir entwickeln eine Maximum-Likelihood-Methode, die auf der Gaußschen Theorie basiert, um regulatorische Interaktionen aus zeitunabhängigen Genexpressionsdaten zu inferieren.

Innerhalb eines simulierten Datensatzes vergleichen wir unseren Ansatz mit Methoden der kleinsten Quadrate, die Standardmethoden für die Inferenz von genregulatorischen Netzwerken sind. Unsere Inferenz liefert eine genauere Rekonstruktion des Netzwerks, wenn ein erheblicher stochastischer Beitrag zur Systemdynamik vorliegt. Auf der Grundlage von Störexperimenten in der SK-MEL-133 Melanom-Zelllinie stellen wir fest, dass unsere Maximum-Likelihood-Methode zu genaueren Vorhersagen der Genexpression führt als Methoden der kleinsten Quadrate.

Die hohe Variabilität im Ansprechen der Patienten beeinträchtigt erheblich den klinischen Erfolg der Krebsimmuntherapie. Das Verständnis der Determinanten, die die Immunantwort und unerwünschte Nebenwirkungen steuern, ist eine zentrale wissenschaftliche Frage für den Fortschritt von Krebsimmuntherapien [2].

Mehrere Determinanten der Immunantwort werden in der Literatur kontrovers diskutiert [3, 4, 5]. Wir suchen nach patientenspezifischen Informationen über das Ansprechen auf eine Krebsimmuntherapie auf der Grundlage des Tumorgenoms. Innerhalb des Tumorgenoms ist der Beitrag von Frameshift-Mutationen für den Erfolg der Immuntherapie nicht gut untersucht [6]. Von Frameshifts erzeugte Peptide unterscheiden sich stark von Selbstpeptiden und sind ein potenzielles Ziel für das Immunsystem. In unserer Analyse konzentrieren wir uns daher auf Frameshifts.

Wir finden Hinweise darauf, dass Frameshifts mit dem Ansprechen auf eine Immuntherapie zusammenhängen. Dennoch ergibt unsere statistische Analyse, dass Frameshifts nicht signifikant mit dem Erfolg der Immuntherapie verbunden sind. Wir finden einige Hinweise, dass ein versteckter Faktor, z. B. die Mutationsrate, sowohl die Zahl der unbekannten immunogenen Mutationen als auch die Zahl der Frameshifts erhöht. Es gibt jedoch keine Hinweise darauf, dass die Frameshifts kausal für den Erfolg einer Immuntherapie sind.

Abstract

Anomalous biological information processing plays a crucial role in the formation and spread of cancer. Mathematical modelling and inference of signalling pathways is a central challenge in the field of cancer research [1].

We propose a stochastic model of gene regulation and solve our model within a Gaussian theory. We develop a maximum-likelihood estimate based on the Gaussian theory to infer regulatory interactions from steady state gene expression data.

Within a simulated dataset, we compare our method to least squares fits, which are standard methods for gene regulatory network inference [7, 8]. Our estimate provides a more accurate network reconstruction in the regime of a sizeable stochastic contribution to the system dynamics. Based on perturbation experiments in the SK-MEL-133 melanoma cell line, we find that our maximum likelihood method leads to more accurate response predictions than least squares methods.

High variability in patient response encumbers the clinical use of cancer immunotherapy. The understanding of determinants that drive immune response, resistance, and adverse side effects is a central scientific issue to move the field of cancer immunotherapy forward [2].

Several hypothetical response determinants are controversially discussed in the literature [3, 4, 5]. We search for patient-specific information about cancer immunotherapy response based on the tumour genome. Within the tumour genome, the contribution of frameshift mutations to immunotherapy response is less well studied [6]. Frameshift-derived peptides are very different from self-peptides and are a potential prime target for the immune system. Within our analysis, we focus, therefore, on frameshift-derived peptides.

We find slight evidence that frameshift mutations are related to immunotherapy response. Nonetheless, our statistical analysis revealed that frameshift-derived peptides are not significantly associated with immunotherapy response. We find some evidence that a hidden factor, e.g. the mutation rate, increases both the number of unknown immunogenic mutations and the number of frameshift mutations. Still, there is no evidence that the frameshift mutations are causal for immunotherapy response.

Table of contents

1.	Intro	oductio	n	1			
	1.1.	Gene r	regulatory network inference	1			
	1.2.	Cancer	r immunotherapy response prediction	3			
2.	Fou	Foundations					
	2.1.	Theore	etical foundations	5			
		2.1.1.	Information theory and statistical mechanics	6			
		2.1.2.	Stochastic dynamics out of equilibrium	8			
		2.1.3.	Inverse problems and statistical inference	12			
	2.2.	Biologi	ical foundations	14			
		2.2.1.	Central dogma of molecular biology	15			
		2.2.2.	Aspects of genetic mutation within cancer	17			
		2.2.3.	Immune recognition of genetic mutations	19			
3.	Gen	Gene regulatory network inference					
	3.1.	Motiva	tion to gene regulatory network inference	26			
	3.2.	Review	v of experimental protocols and mathematical models	28			
		3.2.1.	Experimental protocols	28			
		3.2.2.	Mathematical models	31			
	3.3.	Stocha	stic model of gene regulation	35			
		3.3.1.	Stochastic model of gene regulation	35			
		3.3.2.	Steady state of our stochastic model of gene regulation	37			
		3.3.3.	Connection to the asymmetric Ising model	39			
	3.4.	Forwar	rd problem	40			
		3.4.1.	Mean field theory	41			
		3.4.2.	Gaussian theory	51			
		3.4.3.	Comparison of mean field theory and Gaussian theory	54			
	3.5.	Inverse	e problem	57			
		3.5.1.	Least squares methods	58			
		3.5.2.	Maximum likelihood method	59			
		3.5.3.	Comparison of least squares and maximum likelihood				
			method	60			
	3.6.	Inferen	ace and response prediction in a melanoma cell line	63			
		3.6.1.	Perturbation experiment	64			
		3.6.2.	Response prediction	64			

		3.6.3. Network reconstruction	65			
4.	. Cancer immunotherapy response prediction 7					
	4.1. Cancer immunotherapy by immune checkpoint blockade					
	4.2.	Hypothetical frameshift-based response determinants	74			
	4.3.	. Mathematical foundations of statistical classification and sur-				
		vival analysis	75			
		4.3.1. Statistical classification	76			
		4.3.2. Survival analysis	79			
	4.4.	Statistical classification	81			
		4.4.1. Hypothetical response determinants	81			
		4.4.2. Response prediction	82			
	4.5.	Survival analysis	84			
		4.5.1. Hypothetical survival determinants	84			
		4.5.2. Survival prediction	88			
	4.6.	Conclusion about the immunogenic potential of frameshift				
		mutation	90			
5.	5 Conclusion					
-	5.1 Likelihood-based gene regulatory network inference					
	5.2.	Frameshift-based cancer immunotherapy response prediction				
	Bibl	iography	95			
	8p3					
	Appendicies					
Ι.	Gene regulatory network inference					
	I.1.	Parallel Glauber dynamics	103			
	I.2.	Mean field theory	103			
		I.2.1. Mean gene expression in second order MFT	104			
		I.2.2. Covariance of gene expression in second order MFT .	107			
	I.3.	Network inference and response prediction	110			
II.	I. Cancer immunotherapy response prediction					
	II.1.	Frameshift-derived peptide sequences	115			
	II.2.	Response classification and survival analysis	115			

1. Introduction

The job of the first eight pages is not to have the reader want to throw the book at the wall, during the first eight pages.

David Foster Wallace

Systems biology is an interdisciplinary field of research that studies complex interactions within biological systems. This thesis combines computational methods from broadly varying research fields such as information theory, computational biology, and statistical mechanics. We unify these methods to learn about complex biological systems based on an interplay between information, probability, and logic.

We address two biological problems within this thesis. In section 1.1, we introduce the inference of gene regulatory networks. The second problem, the response prediction to cancer immunotherapy, is posed in section 1.2.

1.1. Gene regulatory network inference

The regulation of gene expression is fundamental in the complexity and diversity of life. Gene regulation controls everything from the response of unicellular organisms to environmental changes up to cell differentiation and self-organisation in multicellular organisms. The multitude of genes forms a complex gene regulatory network (GRN). The GRN encodes response patterns, cell differentiation, and self-organisation. Obtaining reliable information about gene regulation is essential to understanding the complexity and diversity of biological organisms.

Anomalous gene regulation plays a crucial role in the formation and spread of cancer. Therefore, knowledge about GRNs in increasing levels of complexity, from small signalling cascades to large gene regulatory networks consisting of thousands of genes, is important in cancer biology and designing new targeted therapies against cancer. Mathematical modelling of gene regulation, inference of regulatory interactions, and prediction of gene expression is a promising scientific issue in the field of cancer research [1]. We use a stochastic model of gene expression dynamics based on an interaction network to infer regulatory relationships from gene expression data.

Gene expression is, due to the small number of molecules involved and randomness in transcription and translation, an intrinsic stochastic process [9]. We investigate whether one can learn about regulatory interactions from correlations between fluctuations in gene expression. To this end, we propose a system of stochastic differential equations to model gene expression and employ stochastic calculus to characterise the steady state distribution of our system.

Based on the steady state characterisation, we solve the forward problem, calculating mean gene expression and their covariance given a set of model parameters. For this purpose, we extend approximation methods developed in statistical physics and artificial neural network research.

Gene regulatory network inference (GRNI), the reconstruction of regulatory interactions given gene expression measurements, is an inverse problem. For the solution of this inverse problem, we employ statistical inference and construct a maximum posterior estimate based on our forward problem solution.

We address the question of whether our maximum posterior estimate yields a more accurate network reconstruction than standard methods based on a least squares fit [7]. We reconstruct networks based on simulated data and predict gene expression within a cell line experiment to compare the approaches.

The reconstruction of a GRN based on gene expression data is a computationally challenging problem. Without assumptions on the network structure, the computation of a maximum posterior estimate for the structure in an undirected regulatory model is an NP-hard problem [10].

The number of possible configurations for a model with n nodes and d directed interactions is d^{n^2} . Accordingly, the combinatorial explosion of configurations is a computational challenge [11]. Due to the computational complexity, we focus on small regulatory networks to provide proof of principle.

On account of experimental limitations, quantifying gene expression in single-cell high-throughput experiments usually incurs cell destruction. Therefore we focus on the inference without time-series data, which is typical in high-throughput experiments in the context of GRNI. Because the inverse problem of inference from a single-time point without time-series data is ambiguous, perturbation experiments are used in the field of systems biology to infer causality in GRNs. The idea of a perturbation experiment is to quantify steady state gene expression without any perturbation and then apply a known combination of drugs to the system and measure the gene expression after settling into the perturbed steady state. On account of systematic perturbations, more information is available to determine regulatory mechanisms. The approach of perturbation experiments has been successfully employed to infer GRNs [8].

1.2. Cancer immunotherapy response prediction

Immune checkpoints play a crucial role in the regulation of the immune system. Genetic mutations occur that affect these immune checkpoints and disable the immune system's ability to recognise and destroy tumour cells during the development of many cancers. Checkpoint blockade immunotherapy (CBI), which inhibits these immune checkpoints, enables an immune response again.

For a minority of cancer patients, checkpoint inhibition has an outstanding clinical benefit [12]. Nonetheless, adverse effects and high variability in patient response limits clinical success. The understanding of determinants to CBI response is a key scientific issue in the field of immuno-oncology [2].

The obstacle in predicting response to CBI is the complex behaviour of the immune system in cancer [13, 14]. Due to a large number of hypothetical decisive response determinants, careful analysis is required.

The focus of CBI response prediction based on genetic alterations has been mainly on point mutations. The contribution of frameshift mutations is less well studied. Swanton et al. find that the number of frameshift mutations is significantly associated with CBI response across three melanoma cohorts [6].

We focus on frameshift-derived peptide sequences that are entirely different from self-peptides. These frameshift-derived peptides are a hypothetical rich source of immunogenic targets. We investigate whether one can predict cancer immunotherapy response based on frameshift mutations within the tumour genome. To answer this question, the information content of hypothetical frameshift-related response determinants is analysed. We use tools from computational biology to process patient-specific mutation data. To investigate information about CBI response within frameshift mutations in the tumour genome, we implement statistical inference. We aim to identify patients that likely benefit from CBI and support clinical decision-making.

2. Foundations

In any field, the establishment is seldom in pursuit of the truth, because it is composed of those who sincerely believe that they are already in possession of it.

Edwin Thompson Jaynes

This chapter provides an overview of the foundations on which our network inference and response prediction is based.

We introduce, in section 2.1, computational methods employed within this thesis. We provide references to pioneering papers in which those methods were developed.

In section 2.2, we outline the biological foundations. We refer to excellent and extensive textbooks about biological information processing and the immune system's role in cancer.

2.1. Theoretical foundations

In this section, we outline the theoretical foundations of biological information processing and statistical inference, the mathematical framework to draw conclusions in the presence of uncertainty.

The laws of statistical mechanics characterise biological information processing. An information theoretical view of statistical mechanics, a subjective inference of system properties based on data, is outlined in subsection 2.1.1.

For our stochastic model of gene regulation, we employ a system of stochastic differential equations. We define stochastic differential equations, equilibrium, and detailed balance in subsection 2.1.2. In equilibrium statistical mechanics, interactions are generally assumed to be symmetric. There are no symmetric interactions in various inhomogeneous and irregular biological models, such as our stochastic gene expression model. Biological systems generally run far from equilibrium. Permanent consumption and dissipation of energy lead to a non-equilibrium activity that is at the heart of biological organisms [15].

In subsection 2.1.3 we introduce statistical inference, the update of beliefs to account for new data, as an unbiased approach to inverse problems in the field of statistical mechanics.

2.1.1. Information theory and statistical mechanics

Statistical mechanics deals with systems consisting of a multitude of interacting components. On account of the large system size, the general idea of statistical mechanics is to gain information about macroscopic observables based on microscopic laws describing the interacting components. Thus one leaves out the ambition to gain information about the system on a microscopic level.

Within the last decades, the interplay between statistical mechanics and information science has become more important because applications within information science moved towards large, interacting systems [16].

We retrace a general idea of an information-theoretical approach to statistical mechanics developed by Edwin Thompson Jaynes in the 1957 paper [17]. Jaynes, who had a penchant for logic, constructed probability distributions over system states based on a maximum entropy estimate. He extended the methods of Bayesian inference and employed information theory to interpret statistical mechanics. Edwin Thompson Jaynes' book "Probability Theory: The Logic of Science" [18] was published posthumously in 2003. In this exceptional book, he compiled findings on Bayesian probability and statistical inference.

The theoretical basis for this subjective construction of statistical mechanics is the concept of information entropy,

$$S = -k \sum_{\alpha} p_{\alpha} \ln\left(p_{\alpha}\right) \,, \tag{2.1}$$

as a measure for the uncertainty of an observer about the probabilities of system configurations, p_i . This concept was first discussed by Claude Shannon in his landmark paper "A Mathematical Theory of Communication" [19]. The choice of a basis for the logarithm and the value of constant k in the definition (2.1) merely specify the unit in which uncertainty is quantified. Independent of the unit, information entropy measures the average quantity of information required to represent an event from a probability distribution.

One employs information entropy as a starting point to estimate macroscopic observables within statistical inference. The probability distribution that maximises information entropy subject to constraints, which remain to be specified, is the unbiased estimate of the system state. One requires a normalisation,

$$\sum_{\alpha} p_{\alpha} = 1 \,, \tag{2.2}$$

for a well defined probability distribution as a constrain. Despite normalisation one allows for N additional constraints regarding expectation values,

$$\forall \mathbf{n} \in \{1, \dots \mathbf{N}\} : \sum_{\alpha} p_{\alpha} f_{\mathbf{n}}(\mathbf{x}_{\alpha}) = \langle f_{\mathbf{n}}(\mathbf{x}) \rangle , \qquad (2.3)$$

such as system energy or other constants of motion based on experimental measurements. Thus, one employs Lagrange multiplier, $\lambda_0, \ldots, \lambda_N$, to maximise 2.1 subject to 2.2 and 2.3. A maximum entropy estimate of the probabilities,

$$p_{\alpha} = \exp\left(-\lambda_{0} - \sum_{n} \lambda_{n} f_{n}\left(\mathbf{x}_{\alpha}\right)\right), \qquad (2.4)$$

is obtained. Within the estimated distribution the Lagrange multiplier are determined by the required normalisation and expectation values,

$$\begin{split} \lambda_0 &= \ln\left(Z(\lambda_1, \dots, \lambda_N)\right) \text{ and} \\ \langle f_n(\mathbf{x}) \rangle &= -\frac{\partial}{\partial \lambda_n} \ln\left(Z(\lambda_1, \dots, \lambda_N)\right) \\ \text{with } Z(\lambda_1, \dots, \lambda_N) &= \sum_{\alpha} \exp\left(-\sum_n \lambda_n f_n\left(\mathbf{x}_{\alpha}\right)\right) \,. \end{split} \tag{2.5}$$

In equation 2.5 we defined the partition function, $Z(\lambda_1, ..., \lambda_N)$, which encodes how the p_{α} are partitioned among the microstates within our maximum entropy estimate. Finally, the entropy of the inferred distribution is given by

$$S = \lambda_{0} + \sum_{n} \lambda_{n} \left\langle f_{n} \left(\mathbf{x} \right) \right\rangle \,. \tag{2.6}$$

Estimates about macroscopic observables can be expressed via the partition function, $Z(\lambda_1, ..., \lambda_N)$, and partial derivertives with respect to λ_n .

Within the standard statistical mechanics approach, the probability measure (2.4) is known as a Boltzmann distribution. The Boltzmann distribution,

$$p_i = \frac{1}{Z} e^{-\beta \varepsilon_i} \,, \tag{2.7}$$

gives the probability that a system is in state i depending on energy level, ε_i , and a temperature dependent constant, β . Within an information-theoretical approach, this is the maximum entropy estimate under an expected energy constraint and the corresponding Lagrange multiplier β .

The maximum entropy principle in an information-theoretical approach to statistical mechanics is not a physical law like in an objective interpretation but rather a first principle to avoid evidence-less assumptions about the system to make unbiased predictions about macroscopic observables.

Standard statistical mechanic arguments based on microscopic laws of motion lead to identical macroscopic predictions of time-independent observables. There is no additional information in the laws of motion for the statistical inference of these macroscopic observables than the measurement constraints. We encounter a similar property in our network inference approach, which is based on time-independent observables. We employ relations between these observables without additional information on the system dynamics within our unbiased estimate to gain information about the gene regulatory network.

2.1.2. Stochastic dynamics out of equilibrium

A stochastic process is a mathematical object used to model irregularly fluctuating dynamical systems. The theory of stochastic processes is considered as "one of the most important mathematical developments of the twentieth century" [20] with applications that range from statistical mechanics to biology and finance [21].

We introduce stochastic processes in the first paragraph of this subsection. Within the second paragraph, we briefly discuss steady states, equilibrium and detailed balance in the context of stochastic processes and statistical mechanics.

Stochastic dynamics

The French physicist Paul Langevin modelled stochastic dynamics of molecular systems and invented an analytical approach to stochastic processes in his landmark paper "Sur la théorie du mouvement brownien" [22] in the year 1908. The original Langevin equation describes Brownian motion, the random movement of a particle in a liquid due to collisions with liquid molecules. Langevin applied Newtonian dynamics and modelled the effect of interactions between the particle and the fluid molecules with an irregularly random fluctuating contribution. In reward to his pioneering work the equation

$$\frac{\mathrm{d}}{\mathrm{dt}}x(t) = a(x,t) + b(x,t)\xi(t), \qquad (2.8)$$

which is defined by deterministic functions, a(x,t) and b(x,t), and an irregularly random fluctuating function, $\xi(t)$, is named Langevin equation. For $\xi(t)$ an expectation value of zero and no correlation between $\xi(t)$ and $\xi(t')$ for $t \neq t'$ are required,

$$\langle \xi(t) \rangle = 0 \tag{2.9}$$

$$\langle \xi(t)\xi(t')\rangle = \delta(t-t'). \qquad (2.10)$$

Accordingly $\xi(t)$ has the interesting property of diverging variance. To construct a solution of the Langevin equation the integral over $\xi(t)$ is defined,

$$u(t) = \int_0^t dt' \xi(t') \,. \tag{2.11}$$

The following definition of an Ito stochastic integral is a stochastic generalisation of the well-known Riemann-Stieltjes integral. On account of the vanishing correlation of $\xi(t)$ the function u(t) has the property of a Markov process, such that u(t) - u(t') is independent of u(t'') for all t'' < t' < t. The Markov process u(t) is described by zero drift and a time independent global diffusion equal to one,

$$\lim_{\Delta t \to 0} \frac{1}{\Delta t} \left\langle \int_{t}^{t+\Delta t} dt' \xi(t') \right\rangle \stackrel{(2.9)}{=} 0$$

$$\lim_{\Delta t \to 0} \frac{1}{\Delta t} \left\langle \int_{t}^{t+\Delta t} dt' \int_{t}^{t+\Delta t} dt'' \xi(t') \xi(t'') \right\rangle \stackrel{(2.10)}{=} 1.$$
(2.12)

Such a stochastic process is named in the literature Wiener process in honour of the American mathematician and philosopher Norbert Wiener. The Wiener process is not differentiable, and the Langevin equation (2.8) does not exist in a strict mathematical sense. Therefore one studies the corresponding integral equation,

$$x(t) - x(t_0) = \int_{t_0}^t \mathrm{d}t' a(x,t') + \int_{t_0}^t \mathrm{d}t' b(x,t')\xi(t') \,, \tag{2.13}$$

and defines the Ito stochastic integral as a limit of $N \to \infty$ partial sums over intermittent points, t_i , such that

$$\int_{t_0}^t \mathrm{d}t' G(t')\xi(t') = \lim_{N \to \infty} \left(\sum_{i=1}^N G(t_{i-1}) \left(u(t_i) - u(t_{i-1}) \right) \right) \,. \tag{2.14}$$

9

The Japanese mathematician Kiyosi Itô (伊藤 清) invented the method of stochastic integration and stochastic differential equations, nowadays known as Itô calculus [23]. The basic concept of Itô calculus is the Itô stochastic integral, which extends the rigorous methods of calculus to the study of stochastic processes. Nowadays, stochastic calculus is applied in various fields, from molecular dynamics to finance, as the mathematical description of a stochastic process. Gardiner's "Handbook of Stochastic Methods" [24] is an excellent textbook to understand stochastic differential equations and to gain confidence in their application. The following methods, mean value formula and Ito's formula, are comprehensively covered in Gardiner's textbook.

To calculate expectation values within a stochastic process, like mean gene expression in our stochastic model, one uses the mean value formula,

$$\left\langle \int_{t_0}^t \mathrm{d}t'\xi(t')G(t') \right\rangle \stackrel{(2.14)}{=} 0, \qquad (2.15)$$

for non-anticipating G(t). A function G(t) is non-anticipating if G(t) does not depend on the future values of the stochastic process. Therefore it is a reasonable assumption for a function with a physical interpretation to be non-anticipating.

We make use of Ito's formula,

$$df(x) = a(x,t)f'(x) dt + \frac{1}{2}b(x,t)^2 f''(x) dt + b(x,t)f'(x)\xi(t)dt, \quad (2.16)$$

which can be derived in an heuristic approach by Taylor expansion of the function f(x). A rigorous proof of Ito's formula is based on the limit of a sequence of random variables. With Ito's formula one can calculate derivatives of functions depending on a single-variable stochastic process described be the Langevin equation (2.8). Our stochastic model of gene regulation is a multi-variable stochastic process. For such a many variable system described by

$$dx_{i} = a_{i}(\mathbf{x}, t) dt + \sum_{j} b_{ij}(\mathbf{x}, t)\xi_{j}(t)dt \qquad (2.17)$$

one can employ a many-variable variant of Ito's formula,

$$df(\mathbf{x}) = \sum_{i} a_{i}(\mathbf{x}, t) \frac{\partial}{\partial x_{i}} (f(\mathbf{x})) dt + \frac{1}{2} \sum_{ij} (\mathbf{b}\mathbf{b}^{\mathrm{T}})_{ij}(\mathbf{x}, t) \frac{\partial^{2}}{\partial x_{i} \partial x_{j}} (f(\mathbf{x})) dt + \sum_{ij} b_{ij}(\mathbf{x}, t) \frac{\partial}{\partial x_{i}} (f(\mathbf{x})) \xi_{j}(t) dt(t) , \qquad (2.18)$$

to calculate derivatives of a function dependent on a multivariate stochastic process, $f(\mathbf{x})$. Such like stochastic processes and quantitative methods have proven to be useful in various applications in natural and social science [21].

Steady states, equilibrium, and detailed balance

A system is in a steady state if the variables which characterise the system are constant in time. Steady states of our stochastic model of gene regulation are the theoretical foundation of our inference approach. A system is in a stead state if and only if the steady state relation,

$$\forall \mathbf{x} : \frac{\partial}{\partial t} p\left(\mathbf{x}\right) = 0, \qquad (2.19)$$

holds, where $p(\mathbf{x})$ is the probability that the system is in the state \mathbf{x} .

In general, steady states are divided into two categories (equilibrium steady states and non-equilibrium steady states). Equilibrium steady states form a subset characterised by further properties, making the equilibrium steady state distribution more accessible. Analytical solutions for non-equilibrium systems exist mostly in one dimension. Recently, a lot of research has been done on non-equilibrium mechanics because many biological and engaging small systems are usually out of equilibrium [15].

An equilibrium steady state is characterised by detailed balance. A system fulfils the detailed balance condition if and only if there exists a unique solution such that the relation

$$\forall \alpha, \beta : p(\mathbf{x}_{\alpha} | \mathbf{x}_{\beta}) p(\mathbf{x}_{\beta}) = p(\mathbf{x}_{\beta} | \mathbf{x}_{\alpha}) p(\mathbf{x}_{\alpha})$$
(2.20)

holds. This condition demands that the net probability flow between each pair of configurations, \mathbf{x}_{α} and \mathbf{x}_{β} , is equal to zero. Detailed balance corresponds to a time-reversal symmetry, where every transition process is balanced out by its reversed process.

Detailed balance is linked to the presence of an energy function associated with a Boltzmann distribution.

In the case of discrete system with configurations, $\{\mathbf{x}_{\alpha}, \dots, \mathbf{x}_{\omega}\}$, we construct an energy function based on the detailed balance relation 2.20,

$$\ln\left(p(\mathbf{x}_{\beta})\right) = \ln\left(\frac{p(\mathbf{x}_{\beta}|\mathbf{x}_{\alpha})}{p(\mathbf{x}_{\alpha}|\mathbf{x}_{\beta})}\right) + \ln\left(p(\mathbf{x}_{\alpha})\right) .$$
(2.21)

For this purpose we fix the energy for an arbitrary state, $\varepsilon_{\alpha} = \varepsilon_0$, and determine the energies,

$$\varepsilon_{\beta} = \ln\left(\frac{p(\mathbf{x}_{\beta}|\mathbf{x}_{\alpha})}{p(\mathbf{x}_{\alpha}|\mathbf{x}_{\beta})}\right) + \varepsilon_{a}$$
(2.22)

with non-vanishing transition-probabilities $p(\mathbf{x}_{\alpha}|\mathbf{x}_{\beta}) \neq 0$. This assignment is iteratively repeated for the entire configuration space. Kolmogorov's criterion assures that for an irreducible process the detailed balance condition is fulfilled if and only if

$$p(\mathbf{x}_{\alpha}|\mathbf{x}_{\beta})p(\mathbf{x}_{\beta}|\mathbf{x}_{\gamma})\cdots p(\mathbf{x}_{\psi}|\mathbf{x}_{\omega})p(\mathbf{x}_{\omega}|\mathbf{x}_{\alpha})$$

= $p(\mathbf{x}_{\alpha}|\mathbf{x}_{\omega})p(\mathbf{x}_{\omega}|\mathbf{x}_{\psi})\cdots p(\mathbf{x}_{\gamma}|\mathbf{x}_{\beta})p(\mathbf{x}_{\beta}|\mathbf{x}_{\alpha})$ (2.23)

holds for all sequences $(\alpha, \beta, \gamma, \dots, \psi, \omega)$ [25]. Thus, Kolmogorov's criterion guarantees a consistent energy distribution for each irreducible subset of discrete system configuration.

For a continuous system that obeys the detailed balance condition for all times one can define an energy function,

$$\varepsilon(\mathbf{x}) = -\beta^{-1} \log(p^{\mathbf{s}}(\mathbf{x})), \qquad (2.24)$$

based on the corresponding steady states distribution, $p^{s}(\mathbf{x})$. Gardiner establishes in Stochastic Methods section 5.3.4. [24] necessary and sufficient conditions for a system described by the multivariate Langevin equation (2.17) to have a stationary solution that satisfies the detailed balance condition. Under the assumption of linearity and constant coefficients (2.17) can be written as

$$\mathrm{d}x_{\mathrm{i}} = \sum_{\mathrm{j}} a_{\mathrm{ij}} x_{\mathrm{j}} \,\mathrm{d}t + \sum_{\mathrm{j}} b_{\mathrm{ij}} \xi_{\mathrm{j}}(t) dt \,. \tag{2.25}$$

Within this system, Gardiner shows in section 5.3.6. that the detailed balance condition implies that $p^{s}(\mathbf{x})$ is Gaussian distributed,

$$p^{\mathbf{s}}(\mathbf{x}) = Z^{-1} \exp\left(-\frac{1}{2}\mathbf{x}^{\mathrm{T}} \sigma^{-1} \mathbf{x}\right), \qquad (2.26)$$

with $\mathbf{a} \sigma + \sigma \mathbf{a}^{\mathrm{T}} = -\mathbf{b}$.

2.1.3. Inverse problems and statistical inference

Statistical inference is an unbiased update of beliefs to account for new data, \mathcal{D} . The inference is accomplished by employing Bayes theorem

$$\mathcal{P}(\theta|\mathcal{D},\mathcal{M}) = \frac{P(\mathcal{D}|\theta,\mathcal{M})P(\theta|\mathcal{M})}{P(\mathcal{D}|\mathcal{M})}$$
(2.27)

to calculate the posterior probability [26, 18]. The posterior, $\mathcal{P}(\theta|\mathcal{D}, \mathcal{M})$, characterises the probability of a set of model parameters θ , treated as a random variable, under the condition of evidence from new data and a model hypothesis, \mathcal{M} . The prior, $P(\theta|\mathcal{M})$, is a probability distribution encoding beliefs about the model parameter before some evidence from data is taken into account. The likelihood of a model parameter set, $P(\mathcal{D}|\theta, \mathcal{M})$, is the conditional probability of measuring \mathcal{D} given \mathcal{M} and the θ . One obtains the evidence,

$$P(\mathcal{D}|\mathcal{M}) = \int \mathrm{d}\theta P(\mathcal{D}|\theta, \mathcal{M}) P(\theta|\mathcal{M}), \qquad (2.28)$$

by integrating out the model parameter. The evidence is a θ -independent normalising constant, which can be neglected within a maximum posterior estimate.

Within the language of statistical mechanics, the forward problem is the calculation of an observable, \mathcal{D} , based on a model, \mathcal{M} , equipped with a set of model parameters, θ . Thus, the quantification of the likelihood, $P(\mathcal{D}|\theta, \mathcal{M})$, is connected to the forward problem in a statistical mechanics problem.

The inverse problem, the estimation of model parameters based on macroscopic observables, is linked to statistical inference, which provides us with an unbiased approach to inverse problems.

Conjugate prior

For the construction of a suitable prior, we use conjugate prior distributions within this thesis. A prior distribution, $P(\theta|\mathcal{M})$, is a conjugate prior for the likelihood if and only if the posterior, $\mathcal{P}(\theta|\mathcal{DM})$, is in the same family of probability distribution as the prior itself. A conjugate prior shows transparently how our beliefs are updated to account for new data because the update modifies only some parameters in the prior distribution. In case of repetitive updates, a conjugate prior is algebraically facile to handle.

We employ a Gaussian distribution within the response prediction to cancer immunotherapy as a likelihood function. In the standard form,

$$P(x|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(x-\mu)\right),$$
 (2.29)

the distribution is characterised by two parameters, mean, μ , and variance, σ^2 . Assuming $\mu = 0$ is fixed, then the conjugate prior for σ^2 is an inverse Gamma distribution,

$$P(\sigma^2|\alpha,\beta) = \frac{\beta^{\alpha}}{\Gamma(\alpha)} (\sigma^2)^{-\alpha-1} \exp\left(-\frac{\beta}{\sigma^2}\right).$$
 (2.30)

Thus, the posterior is once more an inverse Gamma distributions and in the same family as the prior,

$$P(\sigma^2|x,\alpha,\beta) \propto (\sigma^2)^{-(\alpha+\frac{1}{2})-1} \exp\left(-\frac{\beta+\frac{1}{2}x}{\sigma^2}\right).$$
 (2.31)

Credible region

A credible region is a region of the model parameter space within which a parameter configuration falls with a particular probability. Within our statistical inference, we employ the concept of credible regions to compare a set of model parameters with a posterior probability distribution.

The φ credible region is a parameter space region, A, that is not necessarily connected, such that A contains the fraction φ of the posterior probability,

$$\int_{A} P(\theta) \,\mathrm{d}\theta = \varphi \,. \tag{2.32}$$

The φ credible region with the smallest volume, which contains the probability φ according to constrain (2.32),

$$A^{\star} = \arg\min_{A} \int_{A} \mathrm{d}\theta \,, \qquad (2.33)$$

is the highest posterior density region. Within our analysis we refer to highest posterior density regions as a $n\sigma$ credible region. This region is defined via

$$\varphi = \operatorname{erf}\left(\frac{n}{\sqrt{2}}\right),$$
 (2.34)

such that 1σ corresponds to $\varphi \approx 0.68$ and 2σ correspond to $\varphi \approx 0.95$.

2.2. Biological foundations

Within this section, we introduce the biological principles of our network inference and response prediction.

In subsection 2.2.1 we outline the basic biochemical principles of gene regulation. Based on these, we build our stochastic model of gene regulation, which is the starting point for our network inference.

Genetic mutations in cancer, discussed in subsection 2.2.2, and the immune recognition of these mutations, outlined in the closing subsection 2.2.3, are the biological foundation of our response prediction.

2.2.1. Central dogma of molecular biology

Biological information, the determination of nucleic acid sequences in the genome and amino acid sequences in protein, plays a central role in the complexity and diversity of life. The genetic code enables biological systems to synthesise a multitude of structurally and functionally diverse proteins based on sequential information within the genome.

The central dogma of molecular biology describes the flow of biological information within the cell. Constituents of this process are deoxyribonucleic acid (DNA) and ribonucleic acid (RNA) molecules, both composed of a nucleic acid sequence, as well as proteins, consisting of an amino acid sequence. The general flow of biological information consists of three sub-processes according to the central dogma of molecular biology. These processes are depicted in figure 2.1 and outlined in the following paragraphs. In the last paragraph, we briefly discuss some mechanisms of error recognition within biological information processing.

Replication

The biochemical process of DNA replication is the building of two identical DNA replicas from one template. DNA replication is the fundamental process in cell division, providing each cell with a complete set of biological information, and is thus part of information transfer in all living organisms. DNA consists of two complementary nucleic acid strands, which split into single strands during replication. Each DNA single-strand carries the complete genetic information and serves as a building plan for a DNA replica. The cell-division cycle regulates the replication of DNA.

Transcription

Transcription is the synthesis of a single-stranded RNA molecule based on the nucleotide sequence of a double-stranded DNA molecule. The central player of transcription is the protein complex RNA polymerase, which produces a complementary RNA strand. This stand is called primary transcript.

One refers to RNA segments that encode sequential information for protein synthesis as messenger RNA (mRNA), while all other parts are called non-coding RNA (ncRNA). The ncRNA include several functional subtypes, which account for the majority of human RNA and have a widespread role in cells [27]. The biochemical process of RNA splicing transforms newlytranscribed precursor mRNA molecules into mature mRNA molecules. Noncoding mRNA regions, referred to as introns, are removed during RNA splicing, and the coding regions, called exons, are joined together. The regulation of transcription, a fundamental building block in the complexity and diversity of life, is accomplished by transcription factors (TFs). These are proteins that control the transcription rate by binding to DNA. Thus the interplay of TFs provides the right amount of transcripts according to cell type, homeostasis, and external factors.

Translation

Translation is the protein synthesis according to sequential information encoded in an mRNA molecule. An initiation complex binds to the mRNA molecule, and a ribosome assembles around the targeted mRNA to start the synthesis. A ribosome is a macromolecular machine, found within all living cells, that performs protein synthesis. In this process, the sequential information is translated into an amino acid sequence, which is the basis of a functional protein. The mechanisms of transcription regulation are pretty diverse. They range from specific chemical catalysts to a shutdown of initiation by modification of an initiation complex. However, the most crucial target of gene regulation within the human organism is not the translation, but the regulation of transcription by TFs [28].

Despite these fundamental processes of sequential information flow, posttranslational modification, the chemical modification of a protein that occurs after translation, is an essential component in cell signalling and information processing. An essential modification in the context of signalling is phosphorylation, the reversible attachment of a phosphoryl group to a protein. Within this thesis, we refer to the mitogen-activated protein kinase (MAPK) pathway as a central signalling pathway in the development and spread of cancer. The MAPK pathway consists of proteins that communicate a growth signal via phosphorylation from a cell-surface receptor to the cell nucleus. Thus abnormal information processing within the MAPK pathway is a key mechanism for the formation of many cancer types.

Error recognition

Based on the stop codon, a nucleotide triplet that terminates the protein synthesis, two mechanisms prevent erroneous mRNA from translation.

One mechanism is nonsense-mediated mRNA decay (NMD), which recognises a premature stop codon in mRNA and the prevention of truncated protein expression. During initial translation, the ribosome removes exon-exon junction complexes, which are protein complexes connecting exons. Exonexon junction complexes located after a stop codon remain bound to the

Figure 2.1.: General flow of biological sequence information according to the central dogma of molecular biology. Created with BioRender [30].



mRNA molecule because the ribosome complex does pass them. In case an exon-exon junction complex is still bound to the mRNA after initial translation, NMD takes place. With this mechanism, premature stop codons about 50 nucleotides before an exon-exon junction, and therefore not in the last exon, are detected [29].

The other mechanism is non-stop decay (NSD). This control mechanism prevents mRNA molecules that do not contain a stop codon from translation. Ribosomes are detached from the mRNA molecule after they pass a stop codon. The mRNA decay is initiated by ribosomes that reach the end of an mRNA molecule. Thus, NSD safeguards against point mutations within the stop codon and frameshift mutations resulting in mRNA molecules without a proper stop codon.

2.2.2. Aspects of genetic mutation within cancer

All types of cancer are a consequence of genetic alterations that have occurred purely by chance in the genome of cancer cells. These genetic alternations comprise genetic mutation and chromosomal aberration. Chromosomal aberrations are either numeric, a change in chromosome number, or structural, an abnormal spacial chromosome configuration. Errors during cell division are typically the origin of these chromosomal aberrations, some of which are central in the formation of cancer. A change in sequential information is called a genetic mutation. These genetic mutations are gained during DNA reproduction or as a consequence of external factors. Specific sequential changes enable the development and spread of cancer.

Genetic mutations within a multicellular organism are classified according to their origin in germline and somatic mutations. A germline mutation is a genetic mutation within a germ cell. These cells divide to produce all of the organism's cells, causing germline mutations to be present in every cell in the organism and to be inherited to offspring. In contrast, somatic mutations occur in cells other than germline cells and are acquired by internal or external factors during the organism's lifespan. While somatic mutations are not inherited to offspring, somatic mutations will be present in all cells which are derived from the mutated cell by cell division. The development of cancer is the result of an accumulation of somatic mutations in the course of life.

In cancer research, a driver mutation is defined as a somatic mutation within a signal transduction pathway that provides a growth advantage to the tumour cell while maintaining a beneficial microenvironment, thereby promoting cancer cell proliferation. In contrast, passenger mutations do not promote cancer cell proliferation. Genomic instability and high mutation rates cause cancerous cells to acquire numerous genetic mutations. Still, most mutations are classified as passengers because they do not provide a growth advantage to the tumour cell.

In addition to the property of tumour growth promotion, cancer cell fraction (CCF) is an essential characteristic of a genetic alteration, especially concerning an immune response. The CCF is the fraction of cancer cells within which a specific genetic alteration is present. Thus one refers to a genetic alteration as clonal if the corresponding CCF is close to one, otherwise as non-clonal. We will employ CCF in our immunotherapy response prediction.

Mutations can alter the sequential information within a gene in numerous ways. Therefore multiple types of genetic alterations exist. Within this thesis, we focus on frameshift mutations, which cause peptide sequences highly distinct from self-peptides. Point mutations, which are mentioned within our response prediction, are represented in figure 2.2. We refer to missense mutations, which cause a change in one amino acid, nonsense mutations, resulting in a premature stop codon, and silent mutations, which do not change the amino acid sequence.

Figure 2.3 shows a graphical representation of how frameshift mutations

change amino acid sequences of proteins. Protein biosynthesis is based on a non-injective mapping from a nucleotide triplet, referred to as codon, to an amino acid. Thus, the insertion or deletion of a number of nucleotides, which is not a multiple of three, results in a frameshift-derived peptide sequence that is entirely different from the original sequence and therefore predestined for immune recognition.

2.2.3. Immune recognition of genetic mutations

The immune system is an interacting defence system of biological organisms, consisting of a multifold of biochemical structures. These constituents collaborate to respond to diverse pathogens, such as microorganisms, viruses, toxins, and genetic mutations. The complexity and beauty of the whole human immune system are presented in the textbook [14]. This subsection focuses on the immune reaction caused by genetic mutations in the development and spread of cancer.

The immune system consists of two distinct subsystems, which respond to pathogens on different levels of complexity and time scales. The innate immune system provides a pre-programmed response on short time scales to common pathogens, whereas the sophisticated adaptive immune response is highly specific to each pathogen. Based on information about encountered pathogens, the adaptive immune system can learn and provide a reinforced response when exposed to the same pathogen again. Within this subsection, we focus on the adaptive subsystem, which is decisive for an anti-tumour response.

A critical part of pathogen recognition by the adaptive immune system is the binding of antibodies to antigens. An antibody molecule is a large protein that binds to a specific pathogen molecule, called an antigen. Antibody binding initiates a rigid response against the encountered pathogen by destructive parts of the immune system.

Antigen-presenting cells (APCs) form a heterogeneous group of cells that mediate immune response by presenting antigen. These cells express major histocompatibility complex (MHC) molecules on their cell surface to display peptides consisting out of 8 to 11 amino acids. Almost all cells present peptides out of their proteome, whereas professional APCs present foreign antigens out of the cellular environment to orchestrate the immune response.

Pathogen recognition is performed by lymphocytes such as B cells, which develop in the bone marrow and T cells, which migrate to the thymus gland to mature. B cells provide antibodies, which bind to a specific antigen, and additionally present antigens on their cell surface such that they are classified as professional APCs. To detect somatic mutations and the differentiation



Figure 2.2.: Point mutations. Created with BioRender [30].



Figure 2.3.: Frameshift mutations. Created with BioRender [30].

between self and non-self protein, the function of T cells is crucial. Thus T cells provide a safeguard against the development and spread of cancer. The T-cell receptor (TCR) is a protein complex on the T cell surface that is subject to recognising non-self peptides, which are bound to MHC molecules. The great diversity of TCRs, enabling the recognition of a multitude of hypothetical non-self peptides, results from genetic DNA segment recombination in individual T cells. For self non-self discrimination, T cells go through positive selection within the thymus. Such that almost exclusively non-self peptides trigger TCR-mediated cell destruction. The activation of a T cell requires an extracellular stimulatory signal that professional APCs mediate. Figure 2.4 depicts an overview of an anti-tumour immune response.

Immune checkpoints are activating and inhibitory regulatory mechanisms of the immune system, which are crucial for self-tolerance. Nevertheless, some cancers hide from detection by the immune system by using inhibitory checkpoints such as anti-programmed cell death protein 1 (PD-1), antiprogrammed cell death ligand 1 (PD-L1), and anti-cytotoxic T-lymphocyteassociated protein 4 (CTLA-4).

PD-1 is a cell-surface protein on T and B cells. PD-L1 binding to PD-1 triggers inhibitory signals downstream of the TCR, which block TCR-mediated cell destruction. This binding hinders autoimmune diseases but may also prevent the immune system from a forceful anti-cancer response.

CTLA-4 is expressed by activated T cells, preventing an overreaction of the immune system. Thus CTLA-4 is classified as an inhibitory checkpoint. As well as PD-1 and PD-L1, CTLA-4 is a hypothetical target for checkpoint blockade cancer immunotherapy.

Figure 2.4.: T-cell mediated immune response to tumour cells. Created with BioRender [30].



3. Gene regulatory network inference

If people do not believe that mathematics is simple, it is only because they do not realise how complicated life is.

John von Neumann

The problem of gene regulatory network inference (GRNI) is the primary objective of this chapter. We address the question of whether one can learn something out of correlations between gene expression measurements for the reconstruction of regulatory interactions. Therefore we develop a maximum likelihood method (MLM), and we compare the performance of our MLM to least squares fits, which are standard methods for GRNI [7, 8].

In section 3.1, we highlight research areas where knowledge of GRN is essential. We motivate the reconstruction of small signalling cascades to large gene regulatory networks consisting of thousands of genes.

The revealing of unknown regulatory relationships is based on experimental protocols to quantify gene expression, a mathematical model of the gene regulatory network (GRN), and the inference of model parameters given gene expression data. In section 3.2 the state-of-art of experimental protocols and mathematical approaches to GRNI is reviewed.

Gene expression is an intrinsic stochastic process. We try to learn something from correlations between fluctuations in gene expression about the regulatory interactions. Therefore, we propose a stochastic model of gene regulation. We model gene regulation dynamics by a system of stochastic differential equations. Our stochastic model of gene regulation, which we define in section 3.3, is based on a deterministic model developed by Nelander [7].

We divide the complex challenge of GRNI up into two problems. The first problem is the forward problem, the calculation of mean and covariance of gene expression given model parameters. The second problem is the inverse problem, the inference of model parameters given gene expression data. An accurate solution to the forward problem is crucial for our inference approach, as our MLM depends to a large extent on the forward problem of our stochastic model. Therefore, we focus in the first step on a precise solution to the forward problem.

To solve the forward problem, we employ a Gaussian theory (GT) developed by Mézard and Sakallariou in the context of the asymmetric Ising model [31] in section 3.4. We compute mean gene expression levels and their correlations in the steady state. The GT results are compared to the results of mean field theory (MFT) and a numeric solution of our stochastic model. We show that the GT outperforms the mean field approach in the regime of strong gene interaction. We expected this behaviour because our MFT is based on an expansion in the interactions around a factorising model. In contrast, the GT is independent of the coupling strength.

We focus on inverse problems where time-series data is not available, which is typical for high-throughput experiments in the context of GRNI. To solve the inverse problem, we apply a maximum likelihood method (MLM) in section 3.5. Our MLM is based on the GT for the solution of the forward problem. We compare the results of our MLM to standard inference approaches based on least squares fits of the first moments of the steady state distribution [7, 8]. We demonstrate that our MLM outperforms least squares methods in the regime of a significant contribution of stochastic noise to the system dynamics.

Finally, we perform an inference and response prediction of a signalling network in a melanoma cell line based on experimental data in section 3.6. Our findings give evidence that our MLM results in a more precise response prediction than least squares methods.

3.1. Motivation to gene regulatory network inference

The regulation of gene expression encodes the complexity and diversity of life. Unravelling regulatory motives is, therefore, a fundamental step in understanding basic mechanisms of life and the design of targeted therapies against various diseases. GRNI based on high-throughput datasets is an important and unsolved problem. In this section, we describe applications that require GRNI to some extent.

Cellular differentiation is the metamorphosis in which a cell develops from one cell type into another type of cell. The differentiation includes changes in physical shape, metabolic activity, and response to extracellular signals. The metamorphosis of cells is based on altered gene expression. Therefore the understanding of gene regulation is fundamental in understanding the process of cellular differentiation.

Morphogenesis is the developmental process that causes pattern formation and spatial organisation within organisms. The morphogenesis of organisms is based on programmed and self-organised information processing, which is encoded within the GRN [32].

In the field of evolutionary developmental biology, the developmental processes of organisms are compared to gain information about ancestral relationships and the evolution of organisms. Knowledge about the architecture of gene regulatory networks is one building block in evolutionary developmental biology [33].

Another branch of biological research where knowledge about GRNs is essential is the broad field of molecular medicine. In molecular medicine, biochemical structures and mechanisms are studied, and fundamental molecular or genetic causes of diseases are identified. One can use information about gene regulation to develop therapies and personalised medicine to relieve symptoms or to cure diseases [34].

Cancer is worldwide one of the most common causes of death [35]. Preventing premature death from cancer is therefore crucial for global health and the extension of life expectancy.

Knowledge about signalling pathways in both healthy and cancerous cells is essential in the field of cancer research. A signalling pathway consists of a chain of proteins that transfer information by biochemical activation. Signalling pathways regulate basic cell functions and coordinate these functions within their extracellular environment. The ability to perceive information and respond to the environment is the basis of a functional organism and tissue homeostasis. Anomalous biochemical information flow and errors in cellular information processing can lead to structurally and functionally modification of the cell. These modifications may cause diseases such as autoimmunity and cancer [36, 37]. Knowledge about tumour specific anomalous regulation is therefore essential in the design of new targeted drugs against cancer.

Targeted therapies, which use drugs that inhibit specific signalling pathways, are a promising alternative to conventional chemotherapy [38, 39]. A hypothetical target to inhibit is a signalling pathway, which enables or promotes tumour growth. A targeted drug blocks signals that lead to cell growth, cell division or increased lifespan of tumour cells. The systematic design of new targeted drugs is based on GRNI and the prediction of gene expression under the influence of specific drugs or drug combinations.

Based on GRNI, we predict the response in gene expression to targeted drugs in the SK-MEL-133 melanoma cell line. The SK-MEL-133 cell line has functional mutations within the MAPK pathway, which we introduced as a phosphorylation-based signalling cascade in section 2.2.1. The MAPK pathway consists of a set of proteins that communicate a signal from a cell surface receptor to the cell nucleus. The signalling cascade of the MAPK pathway regulates various cellular processes such as response to environmental conditions, proliferation, differentiation and cell death. A mutation within the MAPK pathway can lead to abnormal information flow. Abnormal information flow plays a key role in the development of melanoma [40] and is a necessary step in the formation and spread of many types of cancer. Chemical compounds that target the MAPK pathway are being investigated as promising and contributed to immense progress in the therapeutic treatment of melanoma [41]. Thus knowledge about gene regulation in the SK-MEL-133 melanoma cell line is a starting point for designing new therapies.

3.2. Review of experimental protocols and mathematical models

The challenge of GRNI is to gain information about regulatory relationships based on measurements of gene expression. The revealing of unknown signalling pathways is based on experimental protocols to quantify gene expression, a mathematical model of the GRN, and the inference of model parameters given gene expression data.

Before we outline mathematical approaches to GRNI based on gene expression data in subsection 3.2.2, we give a brief overview of state-of-the-art experimental protocols to quantify gene expression in the first subsection 3.2.1.

3.2.1. Experimental protocols

The transcriptome is the set of RNA, which is expressed at a specific moment in an individual cell or a cell group. Transcriptomics technologies are experimental protocols and corresponding data analysis methods to study the whole or cell-specific transcriptome.

The focus of this section is on protocols quantifying the expression of mRNA. We neglect ncRNAs, which carry out diverse functions in protein synthesis but are not central for the regulation of gene expression.
The historical evolution of transcriptomics technologies is characterised by breakthrough technologies, opening new horizons and making old techniques obsolete. In 1991, experiments identified 609 mRNA sequences of the human brain [42]. Currently, complete functional genomics data and transcriptomes of hundreds of organisms under a multitude of different conditions and diseases are publicly accessible [43, 44].

There are two current experimental technologies for the simultaneous measurement of RNA transcription levels. DNA microarrays quantify the expression of a predefined set of genes [45, 46], whereas DNA sequencing is principally capable of quantifying the expression level of arbitrary sequences [47, 48].

The measurement of protein expression in a single cell is technically challenging. Some protocols have made it possible to quantify the expression of numerous proteins in a single cell [49]. Ongoing technological advances increase the coverage and sensitivity of these approaches and make it possible to measure in parallel mRNA and protein concentrations in a single cell [50]. Combining mRNA and protein expression data could be the starting point for the inference of multi-level gene regulation. However, the focus of this study is to draw conclusions based on single-level expression data.

DNA microarray

In this paragraph, we give a short introduction about DNA microarray technology, and a fully detailed description can be found in the literature [45, 46]. The idea of DNA microarray technology is the binding of fluorescently labelled single-stranded sequences with complementary probes on the microarray and the measurement of DNA concentration by the fluorescence pattern of the microarray. The heart of the technology is a microarray chip. The chip consists of a few up to as many as thousands of DNA probes. Each probe consists of thousands of identical DNA sequences attached to the solid microarray surface. One refers to the spatial location of a specific probe on the microarray spot or feature.

The fluorescently labelled DNA sequences, of which one measures the concentration, are called targets. The targets are applied as a solution onto the microarray chip and bind to their complementary DNA sequences, immobilised on a specific spot. Binding, the hybridisation between two complementary strands of nucleic acids, is based on non-covalent interactions between the molecules. One can estimate the concentration of each DNA sequence from the fluorescence pattern after washing the unbound sequences.

For the measurement of mRNA concentration, one employs reverse tran-



Figure 3.1.: Sequencing with an DNA microarray. Created with BioRender [30].

scription in a preparation step. Reverse transcription is performed by the reverse transcriptase enzyme, which synthesises complementary DNA (cDNA) from an RNA template. In a follow-up step, one quantifies cDNA concentration by microarray technology and deduces the mRNA concentration.

Approaches based on hybridisation are relatively inexpensive and scalable up to high throughput experiments. On the downside, DNA microarray has intrinsic limitations. They are based on a predefined set of probes. Cross-hybridisation, the formation of double-stranded nucleic acid strains between two molecules with similar but not identical complementary sequences, causes experimental noise. Furthermore, the comparison of experimental data from different setups is challenging and requires elaborate normalisation methods.

DNA sequencing

One can find an extensive review of DNA sequencing technology in the literature [47, 48, 51], and we give an overview of the involved technology in this paragraph. Central for the quantification of mRNA expression is the nucleic acid sequence determination of the cDNA. Therefore mRNA is reverse transcribed into cDNA fragments with adaptors at one end, in the case of single-end sequencing, or both ends, in the case of paired-end sequencing. After the reverse transcription of the mRNA follows the next-generation sequencing (NGS) of the cDNA.

The term NGS is a collective name used to describe different highthroughput sequencing technologies. The sequencing process can be divided into preparing a nucleic acid source, binding of labelled constituents, and sequence determination. There are two main approaches to sequencing. The first approach is sequencing by synthesis (SBS), and the second is sequencing by ligation (SBL). In figure 3.2 the main steps of SBS and SBL are depicted. SBS is based on the synthesis of fluorescently labelled nucleotides. In each sequencing cycle, one nucleotide is added.

SBL uses DNA ligase, an enzyme that facilitates the joining of DNA strands and is sensitive to base-pairing mismatches. DNA ligase preferentially joins a complementary probe out of a pool of short labelled DNA strands. Based on the fluorescence pattern, one can reveal the unknown DNA sequence in both sequencing approaches.

mRNA sequencing is capable of quantifying single-cell gene expression on account of the small amount of needed mRNA. Moreover, mRNA sequencing covers an extensive range over which expression levels can be accurately quantified [52]. Due to experimental limitations, the quantification of gene expression in single-cell experiments incurs cell destruction. Therefore the focus in this chapter is on GRNI without time-series data.

3.2.2. Mathematical models

The complex system of gene regulation can be modelled as a network. The regulatory network is composed of nodes, representing genes, and edges, representing regulatory relationships between the genes. The construction of a mathematical model is the theoretical basis for GRNI. The biological interpretation of an inferred regulatory relationship heavily depends on the underlying mathematical model.

In this subsection, we introduce four major approaches (an informationtheoretic model, a Gaussian model, a Bayesian network, and a differential equation model), of which we will combine two models. In section 3.3, we propose a stochastic model of gene regulation based on a system of differential equations. Our MLM for the network reconstruction, which we introduce in subsection 3.5.2, is a Gaussian model of gene expression based on the differential equation model. Figure 3.2.: Next-generation sequencing of a nucleic acid source by synthesis of flourescently labelled nucleotides and ligation of flourescently labelled oligonucleotides. Created with BioRender [30].



Information-theoretic model

Information-theoretic approaches are based on a measure of statistical dependency. One assigns the statistical dependency between two expression levels to the corresponding edge in the network. Measures that are, among others, used in the context of GRNI are correlations and mutual information.

By construction, information-theoretic approaches reveal statistical dependencies, no causal relationships, and do not distinguish between direct and indirect regulatory relationships. With the use of symmetric measures, the inferred networks are intrinsic undirected. Information-theoretic approaches reveal the topological structure of GRNs. On account of the straightforward approach and low computational complexity, they are used in the study of large regulatory networks [53].

Algorithms, which use a local comparison between two-point dependencies, have been proposed to detect indirect interactions and reduce the number of false positive interactions which have a statistical dependence without a direct regular relationship [53, 54].

Bayesian network

A Bayesian network represents the GRN by a set of genes and their conditional dependencies via a directed acyclic graph. In contrast to an information-theoretic model, the Bayesian network intrinsically incorporates causal regulatory relationships between genes. Koller and Friedman provide a comprehensive description of Bayesian networks in their textbook [55].

Within a Bayesian network, one assumes gene expression to be a stochastic process, and random variables, x_i , represent gene expression level. The regulatory relationships are encoded in a conditional probability distribution, $p(x_i|\pi_i)$, where π_i defines a set of gene expression level of parental genes. Gene expression is described by a joint probability distribution

$$p(\mathbf{x}) = \prod_{i} p(x_{i}|\pi_{i}), \qquad (3.1)$$

which factorises on account of the assumption of acyclic regulatory relationships. Inference within a Bayesian network is based on two steps. The first step is to find an optimal set of parental genes. The second step is to estimate the functional dependencies that best describe the gene expression data. The inference based on a Bayesian network is a computationally complex problem. Still, under the assumption of acyclic regulatory relationships, it is possible to infer causal regulatory interactions within the framework of a Bayesian network [56, 57].

Differential equation model

Detailed chemical kinetics and spatial models of molecular dynamics increase the understanding of gene regulation [58, 59]. On account of complexity, these models are parametrised with many sensitive constants, even for small systems. Furthermore, parametrisation may depend substantially on the chemical environment and could be entirely different from dilute solution chemistry.

Instead of a detailed biochemical model, a system of differential equations is used as a generic model of a GRN. In those models time evolution of the gene expression level,

$$\frac{\mathrm{d}x_{\mathrm{i}}}{\mathrm{d}t} = f_{\mathrm{i}}(\mathbf{x}(t), u_{\mathrm{i}}(t)), \qquad (3.2)$$

is a function of the gene expression level, $\mathbf{x}(t)$, a set of model parameter, \mathbf{m} , and eventually an external perturbation, $u_i(t)$, which may be timedependent. One includes the regulatory interactions and constants of gene expression in the model parameter set, Θ . Inference within a differential equation model is based on a cost function that quantifies how the model describes a given data set and model parameter optimisation. It is possible to infer causal regulatory interactions using time series data [60] or steady state gene expression data [7] within a differential equation model.

Gaussian model

Within the framework of a Gaussian model, statistical relationships between gene expression levels are investigated. One combines the measurements of gene expression levels at one point of time and under one experimental condition into a vector, \mathbf{x} . This gene expression vector is assumed to be a Gaussian distributed multivariate random variable. Gene expression is modelled by a multivariate normal distribution,

$$p(\mathbf{x}|\mathbf{m},\chi) = \frac{1}{\sqrt{2\pi \det \chi}} \exp\left(-\frac{1}{2} \left(\mathbf{x} - \mathbf{m}\right)^{\mathrm{T}} \chi^{-1} \left(\mathbf{x} - \mathbf{m}\right)\right), \qquad (3.3)$$

with mean expression level, \mathbf{m} , and a covariance matrix, χ . Within the framework of Bayesian inference, one can infer the model parameter, \mathbf{m} and χ , and reveal topological information about the GRN encoded in the covariance matrix, χ . Based on gene expression data, undirected regulatory mechanisms can be reconstructed by integration of biological prior information about regulatory mechanisms and prior knowledge about network topology [61]. On account of the quadratic nature of the Gaussian model, this approach is limited to infer two-point regulatory relationships.

3.3. Stochastic model of gene regulation

We propose our stochastic model of gene regulation, a system of stochastic differential equations in this section.

The motivation to employ stochastic differential equations to model a GRN is that gene expression is an intrinsic stochastic process. Randomness in transcription and translation leads to significant fluctuations in mRNA and protein levels [9]. Moreover, we try to gain additional information out of correlations between fluctuations in gene expression levels for the network reconstruction. Within this study, we focus on GRNI based on steady state gene expression data. Therefore, we focus on the steady state distribution of our stochastic model of gene regulation.

Within subsection 3.3.1, we define the stochastic model of gene expression. We characterise the steady state of gene expression within our model in subsection 3.3.2. In the closing subsection 3.3.3, we point out that our proposed model has a similar steady state characteristic as the asymmetric Ising model. The gene expression levels and the spin moments obey similar algebraic relations in the steady state. On account of this connection, we extend in the following section methods developed in the context of the asymmetric Ising model to solve the forward problem of our stochastic model of gene regulation.

3.3.1. Stochastic model of gene regulation

We employ a stochastic version of a deterministic model developed in the context of GRNI by Nelander [7]. In this chapter the system of stochastic differential equation,

$$\frac{\mathrm{d}x_{\rm i}^{\mu}}{\mathrm{d}t} = a_{\rm i}\tanh{(h_{\rm i}^{\mu})} - b_{\rm i}x_{\rm i}^{\mu} + c_{\rm i}\xi_{\rm i}, \qquad (3.4)$$

is used to model the dynamics of gene regulation. The time evolution of gene expression levels is modelled with a system of non-linear differential equations, comparable to a Hopfield network in the field of artificial neural networks [62]. The network is composed of a set of nodes, x_i^{μ} , representing the gene expression level. The gene expression is regulated by an external field,

$$h_{\rm i}^{\mu} = \sum_{\rm k} \omega_{\rm ik} x_{\rm k}^{\mu} + \theta_{\rm i} + u_{\rm i}^{\mu} \,, \qquad (3.5)$$

which is called a synaptic field in the context of neuronal networks, and proportional to a gene specific expression coefficient, a_i . The external field

depends on the expression level, $x_{\mathbf{k}}^{\mu}$, via the interaction matrix, $\omega_{\mathbf{i}\mathbf{k}}$, a threshold vector, $\theta_{\mathbf{i}}$, and a perturbation vector, $u_{\mathbf{i}}^{\mu}$. The exponent μ labels different perturbations and corresponding trajectories of gene expression, $\mathbf{x}^{\mu}(t)$. Within this model, we assume an exponential decay of mRNA or protein concentration, which is quantified by decay constant $b_{\mathbf{i}} > 0$.

Gene expression and the decay of the gene products are intrinsic stochastic processes because of the relatively small number of specific gene products. Random fluctuations in transcription and translation cause significant cellto-cell variations in gene product concentrations [9, 63]. Mathematically we model the stochastic nature of gene expression with a rapidly fluctuating term, ξ_i , defined by zero mean (2.9) and correlation (2.10). The magnitude of the stochastic contribution is quantified by the constant $c_i > 0$.

We assume the external field to be linear in the gene expression level. Therefore, only two-gene interactions are considered. We do not incorporate the phenomenon of co-activation nor co-repression, which are three-gene interactions.

Nevertheless, the model can represent essential aspects of GRNs, including oscillatory states, saturation, and homeostasis of gene expression. In figure 3.3 a plot of the simulated time evolution of our stochastic model of gene regulation is shown.

There is no reasonable argument for symmetric interactions between individual genes. The amplification of gene expression by another gene does not imply a reverse amplification. Such that the interaction is typically asymmetric, and our model violates the detailed balance relation (2.20). Therefore, there is no reasonable argument for the assumption that gene expression levels are distributed according to a Boltzmann distribution. On account of the absence of a Boltzmann distribution, we base our network inference on the dynamic model 3.4.

To test our approach to GRNI, we simulate the time evolution according to generated model parameters and take random samples after the system has reached a steady state.

For the generation of the gene interaction matrix, ω , we set $\forall i \neq j : \omega_{ij} = \beta \zeta_{ij}$, where ζ_{ij} is a random variable drawn independently from a normal distribution with mean equal to zero and variance equal to the inverse of total numbers of genes, 1/N. Thus, assuming independent, x_j , the variance of the external field, h_i , would be N-independent.

The constant β , the inter-gene coupling strength, is a measure of the regulation strength between genes. Such a generated GRN is referred to as a fully connected network, in contrast to a sparse network where there are only a few non-zero matrix elements, ω_{ij} . In our proof-of-principle study, self-regulation is not considered, $\forall i : \omega_{ii} = 0$.

Figure 3.3.: Simulated time evolution of gene expression level for a fully connected GRN as small as 10 nodes, where the model parameter are are set to $a_i = 1, b_i = 1$ and $c_i = 0.1$. For the generation of the interaction matrix, ω , the coupling strenght, $\beta = 0.5$, is used.



3.3.2. Steady state of our stochastic model of gene regulation

We characterise the state of our stochastic model with a multivariate probability distribution of gene expression levels. Because of the focus on steady state gene expression data, we study in this subsection the steady state distribution of our model. We describe the distribution of gene expression level with mean vector, \mathbf{m}^{μ} , and covariance matrix, χ^{μ} , defined by

$$m_{\rm i}^{\mu} = \langle x_{\rm i}^{\mu} \rangle \quad \text{and} \quad \chi_{\rm ij}^{\mu} = \left\langle x_{\rm i}^{\mu} x_{\rm j}^{\mu} \right\rangle - \left\langle x_{\rm i}^{\mu} \right\rangle \left\langle x_{\rm j}^{\mu} \right\rangle \,.$$
 (3.6)

The calculation of mean, $\langle x_{i}^{\mu} \rangle$, and two-point correlation function, $\langle x_{i}^{\mu} x_{j}^{\mu} \rangle$, is based on the stochastic differential equation (3.4),

$$\frac{\mathrm{d}}{\mathrm{dt}}x_{\mathrm{i}}^{\mu} = F_{\mathrm{i}}(\mathbf{x}^{\mu}) + c_{\mathrm{i}}\xi_{\mathrm{i}}.$$
(3.7)

The first term on the right hand side of the equation (3.7) describes the deterministic contribution to dynamics,

$$F_{\rm i}(\mathbf{x}^{\mu}) = a_{\rm i} \tanh\left(\sum_{\rm k} \omega_{\rm ik} x^{\mu}_{\rm k} + \theta_{\rm i} + u^{\mu}_{\rm i}\right) - b_{\rm i} x^{\mu}_{\rm i} \,. \tag{3.8}$$

For the calculation of $\langle x_i^{\mu} \rangle$ the steady state average of the time derivative of gene expression (3.7) is set equal to zero and one receives the mean gene

expression in the steady state,

$$0 = \langle F_{\mathbf{k}}(\mathbf{x}^{\mu}) \rangle_{\mathbf{s}}$$
$$\iff \langle x_{\mathbf{i}}^{\mu} \rangle_{\mathbf{s}} = \frac{a_{\mathbf{i}}}{b_{\mathbf{i}}} \langle \tanh(h_{\mathbf{i}}^{\mu}) \rangle_{\mathbf{s}} , \qquad (3.9)$$

where $\langle x_{i}^{\mu} \rangle_{s} = \text{const.}$ and $\langle \xi_{i} \rangle_{s} \stackrel{(2.9)}{=} 0$ is used. For the calculation of $\langle x_{i}^{\mu} x_{j}^{\mu} \rangle$, one starts with the total time derivative of $f_{ij}(\mathbf{x}^{\mu}) = x_{i}^{\mu} x_{j}^{\mu}$. The many variables version of Ito's formula (2.18) is used and one obtains a total time derivative,

$$\frac{\mathrm{d}}{\mathrm{dt}}(x_{\mathbf{k}}^{\mu}x_{\mathbf{l}}^{\mu}) = F_{\mathbf{k}}(\mathbf{x}^{\mu})x_{\mathbf{l}}^{\mu} + F_{\mathbf{l}}(\mathbf{x}^{\mu})x_{\mathbf{k}}^{\mu} + \delta_{\mathbf{k}\mathbf{l}}c_{\mathbf{l}}^{2} + c_{\mathbf{k}}x_{\mathbf{l}}^{\mu}\xi_{\mathbf{k}} + c_{\mathbf{l}}x_{\mathbf{k}}^{\mu}\xi_{\mathbf{l}}.$$
 (3.10)

Similar to the calculation of mean gene expression one performs the steady state time average over the total time derivative (3.10) and yields

$$0 = \left\langle F_{\mathbf{k}}(\mathbf{x}^{\mu})x_{\mathbf{l}}^{\mu} + F_{\mathbf{l}}(\mathbf{x}^{\mu})x_{\mathbf{k}}^{\mu}\right\rangle_{\mathbf{s}} + \delta_{\mathbf{k}\mathbf{l}}c_{\mathbf{k}}^{2}$$

$$\iff \left\langle x_{\mathbf{i}}^{\mu}x_{\mathbf{j}}^{\mu}\right\rangle_{\mathbf{s}} = \frac{a_{\mathbf{i}}}{b_{\mathbf{i}} + b_{\mathbf{j}}}\left\langle \tanh\left(h_{\mathbf{i}}^{\mu}\right)x_{\mathbf{j}}^{\mu}\right\rangle_{\mathbf{s}} + \frac{1}{2}\frac{c_{\mathbf{i}}^{2}}{b_{\mathbf{i}} + b_{\mathbf{j}}}\delta_{\mathbf{i}\mathbf{j}} + \left(\mathbf{i}\leftrightarrow\mathbf{j}\right), \qquad (3.11)$$

where $\langle x_l^{\mu} \xi_k \rangle_s \stackrel{(2.15)}{=} 0$ for the non-anticipating variable x_l^{μ} is used. The term $(i \leftrightarrow j)$ indicates a summand with interchanges indices.

For the solution of the forward problem within mean field approximation, the three point correlation function, $\langle x_i^{\mu} x_j^{\mu} x_k^{\mu} \rangle_s$, turns out to be useful for the calculation of the covariance matrix, χ_{ij}^{μ} . The calculation of $\langle x_i^{\mu} x_j^{\mu} x_k^{\mu} \rangle_s$ is completely analogue to the calculation of $\langle x_i^{\mu} x_j^{\mu} \rangle$, thus the multivariable version of Ito's formula is applied to the function $f_{ijk}(\mathbf{x}^{\mu}) = x_i^{\mu} x_j^{\mu} x_k^{\mu}$ and the steady state average is taken. Collecting the results, the first three moments of the steady state distribution are given by the following set of equations

$$\begin{aligned} \left\langle x_{i}^{\mu} \right\rangle_{s} &= \frac{a_{i}}{b_{i}} \left\langle \tanh\left(h_{i}^{\mu}\right)\right\rangle_{s} \\ \left\langle x_{i}^{\mu} x_{j}^{\mu} \right\rangle_{s} &= \frac{a_{i}}{b_{i} + b_{j}} \left\langle \tanh\left(h_{i}^{\mu}\right) x_{j}^{\mu} \right\rangle_{s} + (i \leftrightarrow j) + \frac{c_{i}^{2}}{b_{i} + b_{j}} \delta_{ij} \\ \left\langle x_{i}^{\mu} x_{j}^{\mu} x_{k}^{\mu} \right\rangle_{s} &= \frac{a_{i}}{b_{i} + b_{j} + b_{k}} \left\langle \tanh\left(h_{i}^{\mu}\right) x_{j}^{\mu} x_{k}^{\mu} \right\rangle_{s} + \frac{c_{i}^{2}}{b_{i} + b_{j} + b_{k}} \delta_{ij} \left\langle x_{k}^{\mu} \right\rangle_{s} \\ &+ (i \leftrightarrow j \leftrightarrow k) , \end{aligned}$$

$$(3.12)$$

where $(i \leftrightarrow j \leftrightarrow k)$ indicates summands with cyclic permutations of indices.

3.3.3. Connection to the asymmetric Ising model

We begin this subsection with an outline of the symmetric Ising model before we focus on the asymmetric Ising model and choose a dynamics. The chosen dynamics leads to a non-equilibrium steady state. We show that spin moments in the steady state of the asymmetric Ising model are similar to the moments of steady state gene expression of our stochastic model of gene regulation.

The Ising model is named after Cologne-born Ernst Ising. He studied the symmetric model in his PhD thesis and solved it in one dimension, proving that there exists no phase transition in the one-dimensional Ising model [64]. His eventful life was marked by racist oppression by the National Socialists, but also by great kindness and joy in teaching students [65].

The free parameter of the symmetric Ising model are binary spin-variables, $s_i \in \{-1, 1\}$, and the Hamiltonian is given by

$$\mathcal{H}(\mathbf{s}) = -\frac{1}{2} \sum_{ij} J_{ij} s_i s_j - \sum_i \theta_i s_i$$
(3.13)

with local magnetic field, $\theta_{\rm i}$, and symmetric interactions, $J_{\rm ij} = J_{\rm ji}$. The equilibrium statistics of the symmetric Ising model are described by the Bolzmann distribution,

$$p(\mathbf{s}) = \frac{1}{Z} \exp\left(-\beta \mathcal{H}(\mathbf{s})\right) , \text{ where}$$

$$Z = \sum_{\mathbf{s}} \exp\left(-\beta \mathcal{H}(\mathbf{s})\right)$$
(3.14)

is the partition function of the system.

A common choice to introduce a time evolution is sequential Glauber dynamics. Within sequential Glauber dynamics one assumes discrete time steps and in every time step a random spin, s_i , is updated according to the probability distribution

$$p(s_{i}(t+1)|s(t)) = \frac{\exp\left(\beta s_{i}(t+1)h_{i}(t)\right)}{2\cosh\left(\beta h_{i}(t)\right)}.$$
(3.15)

Analogously to our stochastic model of gene expression a local effective field acting on spin i is defined as $h_i(t) = \sum_j J_{ij} s_j(t) + \theta_i$. The time evolution according to parallel Glauber dynamics is discussed in the Appendix section I.1.

An asymmetric interaction matrix characterises the asymmetric Ising model. For the asymmetric Ising model, both choices of dynamics (sequential as well as parallel Glauber dynamics) converge to a non-equilibrium steady state [66], which does not hold the detailed balance relation (2.20). The quantification of the non-equilibrium steady state of the asymmetric Ising model is a hard problem. However one can derive a set of self-consistent relations for the first moments of the steady state distribution. We obtain for the first moment of the steady state distribution the relation

$$\begin{aligned} \langle s_{\mathbf{i}}(t+1) \rangle &\stackrel{(3.15)}{=} \frac{1}{N} \langle \tanh(h_{\mathbf{i}}(t)) \rangle + \frac{N-1}{N} \langle s_{\mathbf{i}}(t) \rangle \\ \Rightarrow \langle s_{\mathbf{i}} \rangle_{\mathbf{s}} &= \langle \tanh(h_{\mathbf{i}}) \rangle_{\mathbf{s}} , \end{aligned}$$
(3.16)

where in any time step one out of N spins is updated according to the sequential update rule (3.15). The averages become time independent in the steady state. We obtain the analogous relation

$$\left\langle s_{i}(t+1)s_{j}(t+1)\right\rangle \stackrel{(3.15)}{=} \frac{1}{N} \left\langle \tanh(h_{i}(t))s_{j}(t)\right\rangle + \frac{1}{N} \left\langle \tanh(h_{j}(t))s_{i}(t)\right\rangle + \frac{N-2}{N} \left\langle s_{i}(t)s_{j}(t)\right\rangle \Rightarrow \left\langle s_{i}s_{j}\right\rangle_{s} = \frac{1}{2} \left\langle \tanh(h_{i})s_{j}\right\rangle_{s} + \frac{1}{2} \left\langle \tanh(h_{j})s_{i}\right\rangle_{s} .$$

$$(3.17)$$

for the two-point correlation function in the steady state.

An important property in the context of this project is that the asymmetric Ising model equipped with sequential Glauber dynamics has structurally similar steady state relations as our model of stochastic gene expression (3.12). This is also true for the first moment with parallel Glauber dynamics (I.2), whereas the second moment $\langle s_i s_j \rangle = \langle \tanh(h_i) \tanh(h_j) \rangle$ has a slightly different structure. The calculation for parallel Glauber dynamics is presented in the Appendix section I.1.

We will use the connection between these models in the next section. There, we extend ideas developed in the context of asymmetric Ising models to solve the forward problem of our stochastic model of gene expression.

3.4. Forward problem

The objective of this section is the forward problem of our stochastic model of gene regulation (3.4). We estimate mean, m_i^{μ} , and covariance of gene

expression, χ^{μ}_{ij} , in the steady state dependent on the model parameter including the interaction matrix.

For the forward problem, we employ an MFT in subsection 3.4.1 and a GT in subsection 3.4.2. We show that the GT provides a precise solution to the forward problem and outperforms the MFT in the regime of strong inter-gene coupling, β , within the last subsection 3.4.3.

3.4.1. Mean field theory

The general idea of MFT is to compute high-dimensional sums or integrals over random variables under the assumption that one can neglect dependencies between the variables. An excellent and comprehensive description of MFT can be found in the textbook by Opper and Saad [67], which covers the foundations of different approaches to MFT and demonstrates their application to various areas of probabilistic modelling. The textbook by Mézard, Parisi and Virasoro [68] contains a comprehensive and self-contained presentation of spin glass theory. The theory of spin glasses demonstrates a rich behaviour. Methods to analyse fluctuations around the mean field solution provide insights into other complex systems, such as our stochastic model of gene expression.

In the first paragraph of this subsection, we outline the general idea of MFT based on the free energy function. In the second paragraph, a variational approach to MF is introduced. This variational approach is even applicable in the absence of a free energy function, which is the case in our stochastic model for gene regulation. We solve the forward problem of our stochastic model of gene expression within a mean field approximation in the closing paragraph. To this end, we employ an approach based on an idea by Kappen and Spanjers in the context of asymmetric neural networks [69].

In the second paragraph, a variational approach to MF is introduced. This variational approach is even applicable in the absence of a free energy function, which is the case in our stochastic model for gene regulation.

General idea of MFT

The general idea of MFT is to approximate a large number of additive contributions to a system hamiltonian by a mean field. For this purpose one studies a system with a Hamiltonian,

$$\mathcal{H} = \mathcal{H}^0 + \mathcal{H}^{\text{int}} , \qquad (3.18)$$

composed of an interacting part, \mathcal{H}^{int} , and a non-interacting part, \mathcal{H}^{0} . One assumes that the system, which one tries to approximate, contains only pairwise interactions,

$$\mathcal{H}^{\text{int}} = \sum_{\text{ij}} h_{\text{ij}}(x_{\text{i}}, x_{\text{j}}) \,. \tag{3.19}$$

Within such a bipartite system the Bogoliubov inequality states that the free energy of the whole system, F, is bounded from above,

$$F \le F^0 := \left\langle \mathcal{H} \right\rangle_0 - TS^0 \,, \tag{3.20}$$

by the free energy of the non-interacting system, $F^0 = -k_B T \ln Z^0$ [70]. The statistical properties of the non-interacting system are described by the partition function, $Z^0 = \sum_{\mathbf{x}} e^{-\beta \mathcal{H}^0(\mathbf{x})}$. The non-interacting free energy, F^0 , is used as an approximation from above for the free energy of the entire system, F. To calculate $\langle \mathcal{H} \rangle_0$ and the entropy S^0 one employs the normalised Bolzmann distribution,

$$p^{0}(\mathbf{x}) = \frac{e^{-\beta \mathcal{H}^{0}(\mathbf{x})}}{Z^{0}} = \prod_{i} \underbrace{\frac{e^{-\beta h_{i}(x_{i})}}{Z_{i}^{0}}}_{=:p_{i}^{0}(x_{i})}, \qquad (3.21)$$

of the non-interacting system. The minimisation of F_0 with respect to the non-interacting distribution, $p_i^0(x_i)$, results in a set of self-consistent equations,

$$p_{\rm i}^0(x_{\rm i}) = \frac{{\rm e}^{-\beta h_{\rm i}^{\rm MF}(x_{\rm i})}}{Z_0}\,, \qquad (3.22)$$

in which the pairwise interactions are approximated by the field,

$$h_{i}^{\rm MF}(x_{i}) = \sum_{j} \int h_{ij}(x_{i}, x_{j}) p_{j}^{0}(x_{j}) \, \mathrm{d}x_{j} \,. \tag{3.23}$$

This field, $h_i^{\text{MF}}(x_i)$, is referred to as a MF because it incorporates the time averaged pairwise interaction energy between the system and the degree of freedom x_i .

Variational approach to MFT

We can not construct a free energy function within our stochastic model of gene regulation because there is no such Hamiltonian, \mathcal{H} . On account of the absence of a free energy function, the general idea of MFT, outlined in the

paragraph above, is not applicable within our model. Instead of minimising the free energy, one employs a variational approach based on the distribution of system states in this setting.

The system is characterised with a multivariate probability distribution, $p(\mathbf{x})$, and one minimises the Kullback – Leibler divergence,

$$D_{\rm KL}\left(p(\mathbf{x})|q(\mathbf{x})\right) = \sum_{\mathbf{x}} q(\mathbf{x}) \frac{q(\mathbf{x})}{p(\mathbf{x})} = \left\langle \ln \frac{q(\mathbf{x})}{p(\mathbf{x})} \right\rangle_q, \qquad (3.24)$$

between $p(\mathbf{x})$ and a variational distribution, $q(\mathbf{x})$. The probability distribution, $p(\mathbf{x})$, is characterised by the Boltzmann distribution,

$$p(\mathbf{x}) = \frac{\mathrm{e}^{-\mathcal{H}(\mathbf{x})}}{Z}, \qquad (3.25)$$

with the partition function, $Z = \sum_{\mathbf{x}} e^{-\mathcal{H}(\mathbf{x})}$, as a normalisation constant. Thus, the Kullback – Leibler divergence,

$$D_{\mathrm{KL}}\left(p(\mathbf{x})|q(\mathbf{x})\right) = \ln Z + \underbrace{E[q] - S[q]}_{=:F[q]}, \qquad (3.26)$$

can be expressed in terms of the variational energy, $E[q] = \sum_{\mathbf{x}} q(\mathbf{x}) \mathcal{H}(\mathbf{x})$, and the entropy, $S[q] = -\sum_{\mathbf{x}} q(\mathbf{x}) \ln q(\mathbf{x})$. Such that a minimal Kullback – Leibler divergence, D_{KL} , corresponds to a minimum in the variational free energy, F[q]. The space of distributions is restricted to factorising distributions,

$$q(\mathbf{x}) = \prod_{i} q_i(x_i) , \qquad (3.27)$$

accordingly to the general idea of MFT to approximate the system interaction by a meanfield.

Within the MFT of our stochastic model of gene regulation in the following paragraph, we require the factorising distribution, $q(\mathbf{x})$, to retain the mean expectation value, $\langle x_i \rangle_p = \langle x_i \rangle_q$. This requirement is motivated by the Ising model, where the minimal Kullback – Leibler divergence corresponds to retaining mean expectation values.

To find this result we study the Ising model with binary spin variables, $x_i \in \{-1, 1\}$, as introduced in subsection (3.3.3). The Hamiltonian is given by

$$\mathcal{H}(\mathbf{x}) = -\sum_{ij} J_{ij} x_i x_j - \sum_i x_i \theta_i \,. \tag{3.28}$$

For the factorising distribution, one sets the couplings equal to zero, $J_{ij}^* = 0$, and variates the external fields, θ_i^* , to minimise D_{KL} . The most general form of a factorizing distribution is a product of Bernoulli distributions,

$$q(\mathbf{x}) = \prod_{i} \left(\frac{1 + m_i^*(\theta_i^*) x_i}{2} \right) , \qquad (3.29)$$

where $m_i^*(\theta_i^*)$ is the expectation value of variable x_i . Setting the derivative of the Kullback – Leibler divergence with respect to the θ_i^* equal to zero,

$$0 \stackrel{!}{=} \frac{\partial}{\partial \theta_{\mathbf{i}}^{*}} \sum_{\mathbf{x}} p(\mathbf{x}) \sum_{\mathbf{k}} \ln\left(\frac{1 + m_{\mathbf{k}}^{*}(\theta_{\mathbf{k}}^{*})x_{\mathbf{k}}}{2}\right)$$
$$= \left(\frac{p(x_{\mathbf{i}} = +1)}{1 + m_{\mathbf{i}}^{*}(\theta_{\mathbf{i}}^{*})} - \frac{p(x_{\mathbf{i}} = -1)}{1 - m_{\mathbf{i}}^{*}(\theta_{\mathbf{i}}^{*})}\right) \left(1 - m_{\mathbf{i}}^{*2}(\theta_{\mathbf{i}}^{*})\right)$$
$$= \langle x_{\mathbf{i}} \rangle_{\mathbf{s}} - m_{\mathbf{i}}^{*}(\theta_{\mathbf{i}}^{*})$$
(3.30)

we find the requirement of retaining the mean expectation values.

MFT of our stochastic model of gene regulation

In this paragraph, we employ MFT to solve the forward problem of our stochastic model of gene regulation. The approach is based on an idea developed by Kappen and Spanjers in the context of asymmetric neural networks [69]. Kappen and Spanjers derived self-consistent equations up to the second order in the couplings to calculate mean firing rates and their correlations.

The neural network used by Kappen and Spanjers in their mean field calculation is an asymmetric Ising model equipped with sequential Glauber dynamics. The Ising model is in contrast to our stochastic model of gene regulation discrete. Nevertheless, the steady state averages for the first two moments, $\langle x_i^{\mu} \rangle_s$ and $\langle x_i^{\mu} x_j^{\mu} \rangle_s$, are similar as we showed in subsection 3.3.3.

We follow the arguments by Kappen and Spanjers. The general idea is to approximate the intractable distribution of gene expression level, $p(\mathbf{x}^{\mu})$, with an optimal factorising distribution, $q^*(\mathbf{x}^{\mu})$. Further we expand mean gene expression, m_i^{μ} , and covariance of gene expression, χ_{ij}^{μ} , up to second order in the interaction matrix, ω , around $q^*(\mathbf{x}^{\mu})$. Within our stochastic model of gene regulation we derive a system of self-consistent equations to calculate m_i^{μ} and χ_{ij}^{μ} .

The calculation is based on the steady state averages defined in equation (3.12). We set the model parameter a_i, b_i , and c_i equal to one and substitute

the perturbation, u_i^{μ} , into a perturbation dependent expression threshold, $\theta_i^{\mu} = \theta_i + u_i^{\mu}$. The mean and covariance,

$$\begin{split} m_{\rm i}^{\mu} &\stackrel{(3.12)}{=} \left\langle \tanh\left(h_{\rm i}^{\mu}\right)\right\rangle_{\rm s} \quad \text{with } h_{\rm i}^{\mu}(\mathbf{x}^{\mu}) = \sum_{\rm k} \omega_{\rm ik} x_{\rm k}^{\mu} + \theta_{\rm i}^{\mu} \\ \chi_{\rm ij}^{\mu} &\stackrel{(3.12)}{=} \frac{1}{2} \left\langle (x_{\rm i}^{\mu} - m_{\rm i}^{\mu}) \left(\tanh\left(h_{\rm j}^{\mu}\right) - m_{\rm j}^{\mu}\right) \right\rangle_{\rm s} + ({\rm i} \leftrightarrow {\rm j}) + \frac{1}{2} \delta_{\rm ij} \,, \end{split}$$
(3.31)

are expanded around the factorising probability distribution, $q^*(\mathbf{x}^{\mu}|\theta^{*\mu})$. The construction of $q^*(\mathbf{x}^{\mu}|\theta^{*\mu})$ is based on non-interacting gene expression, $\omega_{ij}^{*\mu} = 0$, and a mean field, $\theta_i^{*\mu}$. We require to recover the mean gene expression of the full gene regulatory model within the factorising distribution,

$$\langle x_{\mathbf{i}}^{\mu} \rangle_{q} \stackrel{!}{=} \langle x_{\mathbf{i}}^{\mu} \rangle_{q^{*}} = \tanh\left(\theta_{\mathbf{i}}^{*\mu}\right) \,.$$
 (3.32)

Such that the external field, $h_i^{\mu}(\mathbf{x}^{\mu})$, is approximated by a mean field, $\theta_i^{*\mu}$. The expansion of m_i^{μ} and χ_{ij}^{μ} at $q^*(\mathbf{x}^{\mu}|\theta^{*\mu})$ is performed in the variables $\delta\omega_{ij}$ and $\delta\theta_i^{\mu}$, which are defined as

$$\begin{aligned}
\omega_{ij} &= 0 + \delta \omega_{ij} \\
\theta_i^{\mu} &= \theta_i^{*\mu} + \delta \theta_i^{\mu} .
\end{aligned}$$
(3.33)

The Taylor expansion of the mean gene expression, m_i^{μ} , up to quadratic order in $\delta \omega_{ij}$ and $\delta \theta_i^{\mu}$ around the factorising distribution,

$$\begin{split} m_{i}^{\mu}\left(\delta\theta,\delta\omega\right) &= m_{i}^{*\mu} + \sum_{j} \left.\frac{\partial m_{i}^{\mu}}{\partial\theta_{j}}\right|_{q^{*}} \delta\theta_{j} + \sum_{jk} \left.\frac{\partial m_{i}^{\mu}}{\partial\omega_{jk}}\right|_{q^{*}} \delta\omega_{jk} \\ &+ \frac{1}{2} \sum_{jk} \left.\frac{\partial^{2} m_{i}^{\mu}}{\partial\theta_{j}\partial\theta_{k}}\right|_{q^{*}} \delta\theta_{j}\delta\theta_{k} \\ &+ \sum_{jkl} \left.\frac{\partial^{2} m_{i}^{\mu}}{\partial\theta_{j}\partial\omega_{kl}}\right|_{q^{*}} \delta\theta_{j}\delta\omega_{kl} \\ &+ \frac{1}{2} \sum_{jklm} \left.\frac{\partial^{2} m_{i}^{\mu}}{\partial\omega_{jk}\partial\omega_{lm}}\right|_{q^{*}} \delta\omega_{jk}\delta\omega_{lm} + \sigma\left(\delta^{3}\right) \,, \end{split}$$
(3.34)

is the basis of our mean field calculation. To calculate first order terms in the Taylor expansion of $m_{\rm i}^{\mu} = \int \mathrm{d}x \left[p\left(\mathbf{x}|\theta\omega\right) \tanh\left(h_{\rm i}^{\mu}\right) \right]$ one employs partial differentiation at q^* ,

$$\tanh\left(h_{\rm i}^{\mu}\right)|_{q^{*}} = m_{\rm i}^{*\mu} \tag{3.35}$$

$$\frac{\partial \tanh\left(h_{i}^{\mu}\right)}{\partial \theta_{j}}\Big|_{q^{*}} = \left(1 - \left(m_{i}^{\mu q}\right)^{2}\right)\delta_{ij}$$

$$(3.36)$$

$$\left. \frac{\partial \tanh\left(h_{\rm i}^{\mu}\right)}{\partial \omega_{\rm jk}} \right|_{q^*} = \left(1 - \left(m_{\rm i}^{\mu q}\right)^2\right) \delta_{\rm ij} x_{\rm k}^{\mu} \,. \tag{3.37}$$

With these partial derivatives and integration over gene expression levels, \mathbf{x} , we obtain first order Taylor coefficients,

$$\begin{aligned} \frac{\partial m_{i}^{\mu}}{\partial \theta_{j}}\Big|_{q^{*}} &= \int d\mathbf{x} \left[\frac{\partial p\left(\mathbf{x}|\theta\omega\right)}{\partial \theta_{j}} \underbrace{\tanh\left(h_{i}^{\mu}\right)}_{\substack{(3.35)\\ =} m_{i}^{\mu q}} \right]_{q^{*}} \\ &+ \int d\mathbf{x} \left[p\left(\mathbf{x}|\theta\omega\right) \underbrace{\frac{\partial \tanh\left(h_{i}^{\mu}\right)}{\partial \theta_{j}}}_{\substack{(3.36)\\ =} \left(1 - \left(m_{i}^{\mu q}\right)^{2}\right) \delta_{ij}} \right]_{q^{*}} \end{aligned}$$
(3.38)

$$\begin{split} \frac{\partial m_{i}^{\mu}}{\partial \omega_{jk}}\Big|_{q^{*}} &= \int d\mathbf{x} \left[\frac{\partial p\left(\mathbf{x} | \theta \omega\right)}{\partial \omega_{jk}} \tanh\left(h_{i}^{\mu}\right) \right]_{q^{*}} \\ &+ \int d\mathbf{x} \left[p\left(\mathbf{x} | \theta \omega\right) \underbrace{\frac{\partial \tanh\left(h_{i}^{\mu}\right)}{\partial \omega_{jk}}}_{\substack{(3.37)\\(1-(m_{i}^{\mu q})^{2})\delta_{ij}x_{k}^{\mu}}} \right]_{q^{*}} \\ &= 0 + \left(1 - (m_{i}^{\mu q})^{2}\right) m_{k}^{\mu} \delta_{ij} \,. \end{split}$$
(3.39)

Summing over all contributions linear in $\delta\omega$ and $\delta\theta^{\mu}$ one obtains first order corrections to mean gene expression,

$${}^{\theta}m_{i}^{\mu q^{*}} := \sum_{j} \left. \frac{\partial m_{i}^{\mu}}{\partial \theta_{j}} \right|_{q^{*}} \delta \theta_{j} = \left(1 - \left(m_{i}^{\mu q} \right)^{2} \right) \delta \theta_{i} \quad \text{and} \tag{3.40}$$

$${}^{\omega}m_{\mathbf{i}}^{\mu q^*} := \sum_{\mathbf{jk}} \left. \frac{\partial m_{\mathbf{i}}^{\mu}}{\partial \omega_{\mathbf{jk}}} \right|_{q^*} \delta\theta_{\mathbf{j}} = \left(1 - \left(m_{\mathbf{i}}^{\mu q} \right)^2 \right) \sum_k \delta\omega_{\mathbf{ik}} m_{\mathbf{k}}^{\mu} \,. \tag{3.41}$$

Based on the condition to recover mean gene expression of the full regulatory system in equation (3.32) one gets a condition,

$$0 \stackrel{!}{=} {}^{\theta}m_{i}^{\mu q^{*}} + {}^{\omega}m_{i}^{\mu q^{*}} + \sigma\left(\delta^{2}\right), \qquad (3.42)$$

on ${}^{\theta}m_{i}^{\mu q^{*}}$ and ${}^{\omega}m_{i}^{\mu q^{*}}$. Using this condition and the results (3.40) and (3.41) one can calculate an analytic expression for $\delta\theta$ in first order mean field approximation,

$$\delta\theta_{\rm i}^{\rm MFT1} = -\sum_k \delta\omega_{\rm ik} m_{\rm k}^{\mu} \,. \tag{3.43}$$

With the analytic expression for $\delta \theta_i^{MFT1}$ we finally calculate mean gene expression for our stochastic model in first order MFT,

$$m_{\rm i}^{\mu\,\rm MFT1} = \tanh\left(\sum_{k}\omega_{\rm ik}m_{\rm k}^{\mu} + \theta_{\rm i}\right)\,,\tag{3.44}$$

using the relation (3.32) and definition of $\delta\theta_i$ in equation (3.33). The mean field approximation for the mean gene expression in first order, $m_i^{\mu MFT1}$, turns out to be equivalent to the mean field approximation of the Ising model.

In the Appendix section I.2.1 we calculate corrections in second order by evaluation of partial derivatives at the factorising distribution q^* and integration over gene expression, **x**. The calculation is conceptually analogue to the first order calculation.

We find the repetitive pattern that in the calculation of higher order corrections one can identify lower order derivertives (such as $\partial m_n^{\mu}/\partial \omega_{jk}|_{q^*}$ and $\partial m_n^{\mu}/\partial \theta_i|_{q^*}$), which we have already calculated in first order MFT. Summing over all contributions quadratic in $\delta \omega$ and $\delta \theta^{\mu}$ one gets second order corrections in the mean gene expression,

$${}^{\theta\theta}m_{\mathbf{i}}^{\mu q^*} := \sum_{\mathbf{jk}} \left. \frac{\partial^2 m_{\mathbf{i}}^{\mu}}{\partial \theta_{\mathbf{j}} \partial \theta_{\mathbf{k}}} \right|_{q^*} \delta\theta_{\mathbf{j}} \delta\theta_{\mathbf{k}} = (-2)m_{\mathbf{i}}^{\mu q} \left(1 - \left(m_{\mathbf{i}}^{\mu q}\right)^2 \right) \delta\theta_{\mathbf{i}}^2 \qquad (3.45)$$

$$\begin{split} {}^{\omega\omega}m_{i}^{\mu q^{*}} &:= \sum_{jklm} \left. \frac{\partial^{2}m_{i}^{\mu}}{\partial \omega_{jk} \partial \omega_{lm}} \right|_{q^{*}} \delta \omega_{jk} \delta \omega_{lm} \\ &= \left(1 - \left(m_{i}^{\mu q} \right)^{2} \right) \left(\left(-2 \right) m_{i}^{\mu q} \left(-\theta_{i}^{MFT1} \right)^{2} \right. \\ &+ \left(-2 \right) m_{i}^{\mu q} \frac{1}{2} \sum_{j} \left(\delta \omega_{ij} \right)^{2} \\ &+ 2 \sum_{m} \delta \omega_{im} \left(1 - \left(m_{m}^{\mu q} \right)^{2} \right) \left(-\delta \theta_{m}^{MFT1} \right) \right) . \end{split}$$

$$\end{split}$$

$$(3.47)$$

Based on the condition to recover the mean gene expression also in second order an equation for the corrections up to quadratic order,

$$0 \stackrel{!}{=} {}^{\theta}m_{i}^{\mu q^{*}} + {}^{\omega}m_{i}^{\mu q^{*}} + \frac{1}{2}{}^{\theta\theta}m_{i}^{\mu q^{*}} + {}^{\theta\omega}m_{i}^{\mu q^{*}} + \frac{1}{2}{}^{\omega\omega}m_{i}^{\mu q^{*}} + \sigma\left(\delta^{3}\right), \quad (3.48)$$

is derived. Using second order corrections in equation (3.45), (3.46), and (3.47), and the previous results (3.40) and (3.41), an analytical expression for $\delta\theta$ in second order mean field approximation,

$$\delta\theta_{i}^{\rm MFT2} = \delta\theta_{i}^{\rm MFT1} + \frac{1}{2}m_{i}^{\mu}\sum_{k}\left(\delta\omega_{ik}\right)^{2}, \qquad (3.49)$$

is derived. Analogue to the first order approximation we finally calculate the mean gene expression of our stochastic model in second order MFT,

$$m_{\rm i}^{\mu\,{\rm MFT2}} = \tanh\left(\sum_{k}\omega_{\rm ik}m_{\rm k}^{\mu} + \theta_{\rm i} - \frac{1}{2}m_{\rm i}^{\mu}\sum_{k}\left(\omega_{\rm ik}\right)^{2}\right)\,,$$
 (3.50)

using the equations (3.32) and (3.33). The basis to MFT approximation of the covariance is the Taylor expansion of χ^{μ}_{ij} around the factorising distribution, q^* ,

$$\chi_{ij}^{\mu} \left(\delta\theta, \delta\omega\right) = \chi_{i}^{\mu q^{*}} + \sum_{k} \frac{\partial \chi_{ij}^{\mu}}{\partial \theta_{k}} \Big|_{q^{*}} \delta\theta_{k} + \sum_{kl} \frac{\partial \chi_{ij}^{\mu}}{\partial \omega_{kl}} \Big|_{q^{*}} \delta\omega_{kl} + \frac{1}{2} \sum_{jk} \frac{\partial^{2} \chi_{ij}^{\mu}}{\partial \theta_{k} \partial \theta_{l}} \Big|_{q^{*}} \delta\theta_{k} \delta\theta_{l} + \sum_{klm} \frac{\partial^{2} \chi_{ij}^{\mu}}{\partial \theta_{k} \partial \omega_{lm}} \Big|_{q^{*}} \delta\theta_{k} \delta\omega_{lm} + \frac{1}{2} \sum_{klmn} \frac{\partial^{2} \chi_{ij}^{\mu}}{\partial \omega_{kl} \partial \omega_{mn}} \Big|_{q^{*}} \delta\omega_{kl} \delta\omega_{mn} + \sigma \left(\delta^{2}\right) .$$

$$(3.51)$$

To expand the term $\frac{1}{2} \int d\mathbf{x} \left[p(\mathbf{x}|\theta\omega) \left(x_{i}^{\mu} - m_{i}^{\mu}\right) (\tanh\left(h_{j}^{\mu}\right) - m_{j}^{\mu}) \right]$ of the covariance (3.31) partial derivatives at q^{*} ,

$$\frac{\partial \left(\tanh\left(h_{\mathbf{i}}^{\mu}\right) - m_{\mathbf{i}}^{\mu}\right)}{\partial \theta_{\mathbf{k}}} \bigg|_{q*} \stackrel{(3.37)(3.39)}{=} 0 \tag{3.52}$$

$$\frac{\partial \left(\tanh\left(h_{i}^{\mu}\right) - m_{i}^{\mu}\right)}{\partial \omega_{kl}} \bigg|_{q*} \stackrel{(3.39)}{=} \left(1 - \left(m_{i}^{\mu q}\right)^{2}\right) \delta_{ik} \left(x_{l}^{\mu} - m_{l}^{\mu}\right) , \qquad (3.53)$$

are used. With equation (3.52) and (3.53) one obtains first order Taylor coefficients,

$$\begin{split} \frac{\partial \chi_{ij}^{\mu}}{\partial \theta_{k}}\Big|_{q^{*}} &= \frac{1}{2} \int d\mathbf{x} \left[\frac{\partial p\left(\mathbf{x} | \theta \omega\right) \left(x_{j}^{\mu} - m_{j}^{\mu} \right)}{\partial \theta_{k}} \underbrace{\left(\tanh\left(h_{i}^{\mu}\right) - m_{i}^{\mu}\right)}_{\substack{(3.35) \\ = 0}} \right]_{q^{*}} \\ &+ \frac{1}{2} \int d\mathbf{x} \left[p\left(\mathbf{x} | \theta \omega\right) \left(x_{j}^{\mu} - m_{j}^{\mu} \right) \underbrace{\frac{\partial \left(\tanh\left(h_{i}^{\mu}\right) - m_{i}^{\mu}\right)}{\partial \theta_{k}}}_{\substack{(3.52) \\ = 0}} \right]_{q^{*}} \\ &+ \left(i \leftrightarrow j \right) \\ &= 0 \end{split} \qquad (3.54)$$

$$\begin{split} \frac{\partial \chi_{ij}^{\mu}}{\partial \omega_{kl}}\Big|_{q^{*}} &= \frac{1}{2} \int d\mathbf{x} \left[p\left(\mathbf{x} | \theta \omega\right) \left(x_{j}^{\mu} - m_{j}^{\mu} \right) \underbrace{\frac{\partial \left(\tanh\left(h_{i}^{\mu}\right) - m_{i}^{\mu}\right)}{\partial \omega_{kl}}}_{\substack{(3.53) \left(1 - (m_{i}^{\mu q})^{2} \right) \delta_{ik} \left(x_{l}^{\mu} - m_{l}^{\mu} \right) \right]}_{q^{*}} \\ &+ \frac{1}{2} \int d\mathbf{x} \left[\frac{\partial p\left(\mathbf{x} | \theta \omega\right) \left(x_{j}^{\mu} - m_{j}^{\mu} \right)}{\partial \omega_{kl}} \underbrace{\frac{(\tanh\left(h_{i}^{\mu}\right) - m_{i}^{\mu}\right)}{\overset{(3.35)}{=} 0}}_{q^{*}} \right]_{q^{*}} \\ &+ (i \leftrightarrow j) \\ &= \frac{1}{4} \left(1 - (m_{i}^{\mu q})^{2} \right) \delta_{ik} \delta_{jl} + (i \leftrightarrow j) \,. \end{split}$$
(3.55)

Only terms with at least one partial derivative with respect to ω are non-vanishing. This result we find also within second order MFT approximation

of the covariance. Accordingly there is no contribution linear in θ ,

$${}^{\theta}\chi_{ij}^{\mu q^*} := \sum_{\mathbf{k}} \left. \frac{\partial \chi_{ij}^{\mu}}{\partial \theta_{\mathbf{k}}} \right|_{q^*} \delta \theta_{\mathbf{k}} = 0.$$
(3.56)

Summing over all contributions linear in $\delta \omega$ one gets a first order correction to the mean gene expression,

$${}^{\omega}\chi_{ij}^{\mu q^*} := \sum_{\mathbf{k}} \left. \frac{\partial \chi_{ij}^{\mu}}{\partial \omega_{\mathbf{k}}} \right|_{q^*} \delta \omega_{\mathbf{k}} = \frac{1}{4} \left(1 - \left(m_{\mathbf{i}}^{\mu q} \right)^2 \right) \delta \omega_{\mathbf{ij}} + \left(\mathbf{i} \leftrightarrow \mathbf{j} \right) \,. \tag{3.57}$$

Plugging in the first order contribution (3.56) and (3.57) into the Taylor expansion (3.51) of χ^{μ}_{ij} one obtains the covariance of gene expression in first order MFT,

$$\chi_{ij}^{\mu \,\text{MFT1}} = \frac{1}{2} \delta_{ij} + \frac{1}{4} \left(1 - \left(m_i^{\mu q} \right)^2 \right) \omega_{ij} + (i \leftrightarrow j) \,. \tag{3.58}$$

The approximation up to second order is conceptually equal to the first order calculation. The calculation leading to the result,

$$\begin{split} \chi_{ij}^{\mu \,\text{MFT2}} &= +\frac{1}{2} \delta_{ij} + \frac{1}{4} \left(1 - (m_i^{\mu q})^2 \right) \omega_{ij} - \frac{1}{2} \left(1 - \left(m_j^{\mu q} \right)^2 \right) \omega_{ji} \theta_j^{\text{MFT1}} \\ &+ \frac{1}{8} \left(1 - (m_i^{\mu q})^2 \right) \left(1 - \left(m_j^{\mu q} \right)^2 \right) \sum_{l} \omega_{il} \omega_{jl} \\ &+ \frac{1}{8} \left(1 - \left(m_j^{\mu q} \right)^2 \right) \sum_{k} \omega_{jk} \left(1 - (m_k^{\mu q})^2 \right) \omega_{ki} \\ &- \frac{m_j^{\mu q}}{6} \left(1 - \left(m_j^{\mu q} \right)^2 \right) \left(m_i^{\mu q} \left(\omega_{ij} \right)^2 \\ &- 2\omega_{ij} \delta \theta_i^{\text{MFT1}} - \frac{m_j^{\mu q}}{2} \sum_{k} \left(\omega_{ik} \right)^2 \right) \\ &+ \left(i \leftrightarrow j \right) \,, \end{split}$$
(3.59)

is given in the Appendix section I.2.2.

To solve the forward problem within MFT in first order, we iteratively calculate \mathbf{m}^{μ} and χ^{μ} according to the set of self-consistent equation (3.44) and (3.58). For the calculation within second order we use respectively the equation (3.50) and (3.59). A simple iterative procedure with initial values $\mathbf{m}_{\text{init}}^{\mu} = \mathbf{0}$ and $\chi_{\text{init}}^{\mu} = \mathbf{1}$ converges and thus we obtain a solution to the forward problem of our stochastic model of gene expression within mean field approximation.

3.4.2. Gaussian theory

In this subsection, we employ a continuous version of the GT developed by Mézard and Sakallariou in the context of the asymmetric Ising model [31] to solve the forward problem of the stochastic model of gene regulation.

GT within the asymmetric Ising model

The asymmetric Ising model, which we introduced in subsection 3.3.3, is characterised by the Hamiltonian

$$\mathcal{H}(\mathbf{s}) = -\frac{1}{2} \sum_{ij} J_{ij} s_i s_j - \sum_i \theta_i s_i \,. \tag{3.60}$$

The model parameter are local external fields, θ_i , and exchange couplings, J_{ij} . To solve the forward problem, Mézard and Sakallariou introduce parallel Glauber dynamics for the asymmetric Ising model, which is outlined in the Appendix section I.1. To calculate the local magnetisation and equal-time correlation,

$$m_{i} \stackrel{(I.2)}{=} \langle \tanh(h_{i}) \rangle$$

$$C_{ij} \stackrel{(I.3)}{=} \langle (\tanh(h_{i}) - m_{i})(\tanh(h_{j}) - m_{j}) \rangle ,$$
(3.61)

Mézard and Sakallariou approximate the sum over system interaction with spin i, $\sum_{j} J_{ij} s_{j}$, with a Gaussian distribution. Thus the system interaction is characterised by mean,

$$g_{i} = \sum_{j} J_{ij} m_{j}, \text{ and variance,}$$

$$\Delta_{i} = \sum_{j} J_{ij}^{2} \left(1 - m_{j}^{2}\right). \qquad (3.62)$$

For the calculation of the variance Mézard and Sakallariou check selfconsistently that the typical correlation, $\langle s_i s_j \rangle - m_i m_j$, is of order 1/N, where N is the system size.

Within this approximation a set of self consistent equations for $m_i(t)$ and $C_{ij}(t)$, which become exact in the limit of a large system size, is derived. The GT provides a precise solution to the forward problem of the asymmetric Ising model in the regime of strong interactions [31].

GT within our stochastic model of gene regulation

To derive a set of self-consistent equations for m_i^{μ} and χ_{ij}^{μ} , we assume the local effective field, h_i^{μ} , defined within our stochastic model of gene expression in equation (3.5) to be Gaussian distributed. A deterministic contribution, $g_i^{\mu}(\mathbf{m}^{\mu})$, and a probabilistic contribution, η_i , to the local effective field,

$$h_{\rm i}^{\mu}(\mathbf{m}^{\mu}) = g_{\rm i}^{\mu}(\mathbf{m}^{\mu}) + \eta_{\rm i}$$
 with $g_{\rm i}^{\mu}(\mathbf{m}^{\mu}) = \sum_{\rm k} \omega_{\rm ik} m_{\rm k}^{\mu} + \theta_{\rm i} + u_{\rm i}^{\mu}$, (3.63)

are defined. The probability distribution of η_i^{μ} is given by a multivariate normal distribution $\mathcal{P}(\eta_i^{\mu})$, which is characterised by zero mean and covariance matrix Δ^{μ} .

The argument for the Gaussian nature of $h_i^{\mu}(\mathbf{m}^{\mu})$ is that in a sizeable system the local effective field, $h_i^{\mu}(\mathbf{m}^{\mu})$, is the sum of a large number of stochastic variables. The requirement of uncorrelated variables in the central limit theorem (CLT) only holds in a tree-like graph. Regarding regulatory motifs like feed forward loops, we can not assume the GRN to be generally tree-like, and we must consider correlated variables. Under rather technical restrictions, there are versions of the standard CLT that allow for correlations between the random variables that are being summed [71, 72]. It is unknown whether these restrictions are valid for the effective local field in our stochastic model of gene regulation. Nevertheless, simulations provide numerical evidence for the Gaussian nature of the local effective field for fully connected systems with as few as ten nodes.

To obtain a set of self consistent equations for m_i^{μ} we start with the steady state relation for the mean, $m_i^{\mu} = a_i/b_i \langle \tanh(h_i^{\mu}) \rangle_s$, which we derived in equation (3.12). Performing the steady state average one obtains the integral equation

$$m_{\rm i}^{\mu} = \frac{1}{\sqrt{2\pi}} \frac{a_{\rm i}}{b_{\rm i}} \int \mathrm{d}\eta \exp\left(\frac{\eta^2}{2\Delta_{\rm ii}^{\mu}}\right) \tanh\left(g_{\rm i}^{\mu}(\mathbf{m}^{\mu}) + \eta\right) \,. \tag{3.64}$$

To solve this integral equation, one has to calculate the variance of the probabilistic contribution to the external field, Δ^{μ}_{ii} . Based on the definition of the external field in equation (3.63), the relation between covariance of external fields, Δ^{μ}_{ij} , and covariance of gene expression level, χ^{μ}_{ij} ,

$$\begin{split} \Delta_{ij}^{\mu} &= \operatorname{cov}(\eta_{i}^{\mu}, \eta_{j}^{\mu}) = \operatorname{cov}(h_{i}^{\mu}, h_{j}^{\mu}) \\ &\stackrel{(3.63)}{=} \left\langle \sum_{k} \omega_{ik} \left(x_{k}^{\mu} - m_{k}^{\mu} \right) \sum_{l} \omega_{jl} \left(x_{l}^{\mu} - m_{l}^{\mu} \right) \right\rangle_{s} \quad (3.65) \\ &= \sum_{kl} \omega_{ik} \omega_{jl} \chi_{kl}^{\mu} = \left[\omega \chi^{\mu} \omega^{\top} \right]_{ij} \,, \end{split}$$

is calculated. Equipped with the equations (3.64) and (3.65) it is left over to find an expression for χ^{μ}_{ij} to solfe the forward problem within the Gaussian approximation. To derive also a set of self consistent equations for χ^{μ}_{ij} , we calculate the steady state averages $\langle \tanh(h^{\mu}_{i}) x^{\mu}_{j} \rangle_{s}$. Therefore the interaction matrix,

$$x_{\rm i}^{\mu} = \sum_{\rm k} \omega_{\rm ik}^{-1} (h_{\rm k}^{\mu} - \theta_{\rm k} - u_{\rm k}^{\mu}) \,, \tag{3.66}$$

is inverted to replace the expression levels, x_{i}^{μ} , with the external fields, h_{k}^{μ} . With the inverted equation (3.66) one can rewrite the steady state average,

$$\left\langle \tanh\left(h_{i}^{\mu}\right)x_{j}^{\mu}\right\rangle_{s} \stackrel{(3.66)}{=} \sum_{k} \omega_{ik}^{-1} \left\langle \left(h_{k}^{\mu}-u_{k}^{\mu}\right) \tanh\left(h_{j}^{\mu}\right)\right\rangle_{s}, \qquad (3.67)$$

as a sum over steady state averages containing only external fields, which we assumed to be Gaussian. To calculate the steady state averages in (3.67) the correlation coefficient,

$$\rho_{jk} = \frac{\operatorname{cov}(h_{j}^{\mu}, h_{k}^{\mu})}{\sqrt{\Delta_{jj}^{\mu} \Delta_{kk}^{\mu}}}, \qquad (3.68)$$

is defined. The multivariate normal distribution of the external field is expanded in linear order in $\rho_{\rm jk},$

$$\mathcal{P}(\eta_{j}^{\mu},\eta_{k}^{\mu}) = \mathcal{P}(\eta_{j}^{\mu})\mathcal{P}(\eta_{k}^{\mu}) \left(1 + \frac{\eta_{j}^{\mu}}{\sqrt{\Delta_{jj}^{\mu}}} \frac{\eta_{k}^{\mu}}{\sqrt{\Delta_{kk}^{\mu}}} \rho_{jk}\right) + \sigma(\rho_{jk}^{2}).$$
(3.69)

Using the expanded probability distribution (3.69), an expression for the steady state averages on the right hand side of (3.67),

$$\left\langle h_{\mathbf{k}}^{\mu} \tanh\left(h_{\mathbf{j}}^{\mu}\right)\right\rangle_{\mathbf{s}} \stackrel{(3.69)}{=} g_{\mathbf{k}}^{\mu} m_{\mathbf{j}}^{\mu} + \lambda_{\mathbf{j}}^{\mu} \operatorname{cov}(h_{k}^{\mu} h_{j}^{\mu}), \qquad (3.70)$$

is calculated. The factor λ_j^{μ} , which is a measure for sensitivity of the first moments to fluctuations in the expression levels, is given by the integral equation

$$\lambda_{j}^{\mu} = \frac{1}{\sqrt{2\pi}} \int d\eta \exp\left(\frac{\eta^{2}}{2\Delta_{ii}^{\mu}}\right) \left(1 - \tanh^{2}\left(g_{j}^{\mu} + \eta\right)\right) \,. \tag{3.71}$$

Using the definition of χ^{μ}_{ij} in equation (3.6), the steady state averages (3.67) and (3.70), one gets an expression for the covariance

$$\begin{split} \chi_{ij}^{\mu} &\stackrel{(3.6)}{=} \frac{a_{i}}{b_{i} + b_{j}} \left\langle \tanh\left(h_{i}^{\mu}\right) x_{j}^{\mu} \right\rangle_{s} + (i \leftrightarrow j) + \frac{c_{i}^{2}}{2b_{i}} \delta_{ij} - \left\langle x_{i}^{\mu} \right\rangle \left\langle x_{j}^{\mu} \right\rangle \\ &\stackrel{(3.67)(3.70)}{=} \frac{a_{i}}{b_{i} + b_{j}} \sum_{k} \omega_{ik}^{-1} \left(g_{k}^{\mu} m_{j}^{\mu} + \lambda_{j}^{\mu} \operatorname{cov}(h_{k}^{\mu} h_{j}^{\mu}) - u_{k}^{\mu} m_{j}^{\mu} \right) \\ &+ (i \leftrightarrow j) + \frac{c_{i}^{2}}{2b_{i}} - m_{i}^{\mu} m_{j}^{\mu} \,. \end{split}$$
(3.72)

With the definition of $g_i^{\mu}(\mathbf{m}^{\mu})$ in equation (3.63) and the covariance of the external field in equation (3.65) one finally obtains a set of self-consistent equations for the covariance matrix of the expression level,

$$\chi^{\mu}_{ij} = \frac{a_j}{b_i + b_j} \left(\chi^{\mu} \omega^{\mathrm{T}} \right)_{ij} \lambda^{\mu}_j + (i \leftrightarrow j) + \frac{c_i^2}{2b_i} \delta_{ij} \,. \tag{3.73}$$

To obtain a solution of the forward problem within GT we solve iteratively the set of self-consistent equations (3.64) and (3.73). In each iteration step the integrals in (3.64) and in (3.71) are solved numerically using adaptive Gauss-Kronrod quadrature [73]. We find that a simple iterative procedure with initial values $\mathbf{m}_{\text{init}}^{\mu} = \mathbf{0}$ and $\chi_{\text{init}}^{\mu} = \mathbb{1}$ converges. Finally, this is our solution to the forward problem within the mean field approximation.

3.4.3. Comparison of mean field theory and Gaussian theory

In this subsection, we discuss the results for m_i^{μ} and χ_{ij}^{μ} within the first two orders MFT and GT. We investigate the precision of MFT and GT depending on the coupling strength, β , which we defined with our stochastic model in section 3.3. Furthermore, we show that GT outperforms MFT in the regime of strong inter-gene couplings.

In figure 3.4, we show scatter plots for the mean, m_i^{μ} , and covariance, χ_i^{μ} , of gene expression obtained by first and second order MFT, and GT. All methods give good results in the regime of weak interaction. For large inter-gene couplings, we find an overestimation of mean gene expression obtained by first order MFT. This behaviour is because local fluctuations in the network are neglected, and the system order is overestimated within the mean field approximation. Correspondingly the variance, $\chi_{ii}^{\mu MFT1} = 1/2$, is underestimated within first order MFT. We find that in second order MFT, the variance is overestimated.

To quantify the precision of MFT and GT we use the relative quadratic error between simulated and predicted values of $m_{\rm i}^{\mu}$ and $\chi_{\rm ij}^{\mu}$,

$$r = \sqrt{\frac{\sum_{i} (q_{i}^{\rm sim} - q_{i}^{\rm pre})^{2}}{\sum_{i} (q_{i}^{\rm sim})^{2}}}.$$
 (3.74)

For mean and covariance, we plot the β -dependence of the relative quadratic error in figure 3.6.

The regime of weak inter-gene regulation corresponds to small values of β . Within this regime, one can see that the MFT in first and second order, as well as the GT, give an accurate estimate of mean gene expression, m_i^{μ} , and covariance of gene expression, χ_{ij}^{μ} . For large β -values, corresponding to the regime of strong inter-gene regulation, MFT breaks down. The reason for this breakdown is that the mean field approximation is based on an expansion of m_i^{μ} and χ_{ij}^{μ} at the factorizing distribution, which is characterised by $\omega_{ij}^* = 0$ corresponding to $\beta = 0$.

We find that for higher coupling strength, the error in the quadratic contribution outweighs the error made in the linear contribution. Thus first order MFT outperforms second order MFT in the regime of strong couplings. The approximation of the external field, h_i^{μ} , with a Gaussian distribution is also valid for strong couplings. The relative quadratic error of m_i^{μ} and χ_{ij}^{μ} within the GT slightly decreases with increasing β . Therefore GT outperforms MFT in a regime of strong inter-gene coupling.

Figure 3.4.: Scatter plot of mean, m_i^{μ} , and covariance, χ_i^{μ} , for first order MFT \bullet , second order MFT \bullet , and GT \bullet . For the generation of the interaction matrix a medium inter-gene coupling strength, $\beta = 0.5$, and a large inter-gene coupling strength, $\beta = 1.0$, are used.



Figure 3.6.: Relative quadratic error of mean and the covariance of gene expression for first order MFT ●, second order MFT ●, and the GT ●.



3.5. Inverse problem

The inverse problem of GRNI is solved within our stochastic model of gene regulation. We employ two generic approaches, least squares fits and our MLM. We find that our likelihood-based approach outperforms the least squares fits in the regime of strong stochastic contribution to the system dynamics.

The least squares fits, which are introduced in subsection 3.5.1, are based on simple quadratic cost functions. These cost functions quantify the accurateness of a given set of model parameters to describe steady state data. The cost functions are based on the exact relations of the first and second moment of the steady state distribution (3.12) and the results of the GT. The inference within least squares methods takes place without a concrete assumption on a statistical model of gene expression fluctuation.

Within our MLM, which is introduced in subsection 3.5.2, we employ the framework of statistical inference. We assume the gene expression level to be distributed according to a multivariate Gaussian distribution characterised by the results of the forward problem, $\mathbf{m}^{\mu}(\omega)$ and $\chi^{\mu}(\omega)$. Under the presumption of no prior information, we obtain an interaction matrix with a maximum likelihood estimate.

For the minimisation within the least squares fits and respectively the likelihood maximisation based on simulated data, we tested local and global optimisation algorithms implemented within the NLopt library for non-linear optimisation [74]. We find no clear evidence for the benefit of global optimisation over a local optimisation starting at the origin, $\forall i, j : \omega_{ij}^{init} = 0$.

Within our study, we employ bound optimisation by quadratic approximation (BOBYQA) [75]. The BOBYQA algorithm is a realisation of a local gradient-free optimisation algorithm. In each iteration step of the algorithm a local quadratic approximation, q, of the function to maximise, f, is computed. The interpolation points are updated by minimizing the Frobenius norm of the second derivative matrix of q. No derivatives of f are required explicitly, thus an implementation of a gradient is not necessary.

The results of GRNI based on simulated data within least squares fits and our MLM are discussed in subsection 3.5.3.

3.5.1. Least squares methods

Least squares fits are used in the reconstruction of GRNs as big as 100 nodes based on perturbation data [76, 8]. The inferred GRNs confirm and extend the knowledge about biological pathways and accurately predict the outcome of untested perturbations.

We employ a least squares fits based on the steady state relation of the first moment (3.12). To obtain an estimate for the interaction matrix,

$$\omega_{\rm ms1o} = \underset{\omega}{\arg\min} \left(\sum_{\mu i} \left[\langle x_i^{\mu} \rangle - \frac{a_i}{b_i} \langle \tanh\left(h_i^{\mu}(\mathbf{x}^{\mu}, \omega)\right) \rangle \right]^2 \right) , \qquad (3.75)$$

a quadratic cost function is minimised with respect to ω . The averages in the cost function are performed over the steady state data. Discrepancies between mean gene expression, $\langle x_i^{\mu} \rangle$, and expected expression in the steady state, $a_i/b_i \langle \tanh(h_i^{\mu}(\mathbf{x}^{\mu}, \omega)) \rangle$, are penalised quadratically. Therefore, a lower cost interaction matrix represents the data more accurately.

The solution of the forward problem within GT intrinsically incorporates the second moment of the gene expression distribution. Therefore, we also examine the information contained in the second moments of the samples with a least squares method. The steady state relation of first and second moments (3.12) is used to infer the interaction matrix, ω_{ms2o} , within a mean square approach in second order,

$$\omega_{\rm ms2o} = \arg\min_{\omega} \left(\sum_{\mu i} \left[\langle x_i^{\mu} \rangle - \frac{a_i}{b_i} \langle \tanh(h_i^{\mu}(\mathbf{x}^{\mu}, \omega)) \rangle \right]^2 + \frac{1}{N} \sum_{\mu i j} \left[\langle x_i^{\mu} x_j^{\mu} \rangle - C_{i j}^{\mu} (h_i^{\mu}(\mathbf{x}^{\mu}, \omega)) \right]^2 \right).$$
(3.76)

In the second line the steady state estimate of the two-point correlation function,

$$C_{ij}^{\mu} (h_{i}^{\mu}(\mathbf{x}^{\mu}, \omega)) \stackrel{(3.12)}{=} \frac{1}{2} \frac{a_{i}}{b_{i} + b_{j}} \left\langle \tanh\left(h_{i}^{\mu}(\mathbf{x}^{\mu}, \omega)\right) x_{j}^{\mu}\right\rangle + (i \leftrightarrow j) + \frac{1}{2} \delta_{ij} c_{j}^{2}, \qquad (3.77)$$

is used. On account of a system-size independent relation between the contributions of first and second moments, we chose a relative factor of the inverse number of nodes, 1/N.

We additionally use the results of GT developed in section 3.4.2, $m_{\rm i}^{\mu}(\omega)$ and $C_{\rm ij}^{\mu}(\omega) = \chi_{\rm ij}^{\mu}(\omega) + m_{\rm i}^{\mu}(\omega)m_{\rm j}^{\mu}(\omega)$, in a least squares approach. Thus we infer the interaction matrix, $\omega_{\rm msGt}$, making use of the results of GT for the first two moments,

$$\omega_{\rm msGt} = \arg\min_{\omega} \left(\sum_{\mu i} \left[\langle x_i^{\mu} \rangle - m_i^{\mu}(\omega) \right]^2 + \frac{1}{N} \sum_{\mu i j} \left[\langle x_i^{\mu} x_j^{\mu} \rangle - C_{ij}^{\mu}(\omega) \right]^2 \right).$$
(3.78)

3.5.2. Maximum likelihood method

We try to get an unbiased estimate about the model parameter based on gene expression data. Therefore we employ a maximum likelihood estimate. The obstacle of our likelihood-based approach to GRN inference is the likelihood calculation.

We quantify the conditional probability, $p({\mathbf{x}^{\mu}} | \omega, \mathbf{a}, \mathbf{b}, \mathbf{c})$, of measuring a set of gene expression values, ${\mathbf{x}^{\mu}}$, given a complete set of model parameters. Therefore, we employ GT to solve the forward problem. We obtain mean, $\mathbf{m}^{\mu}(\omega)$, and covariance of gene expression, $\chi^{\mu}(\omega)$, given the model parameter. Within our MLM we assume the data to be drawn from a multivariate Gaussian distribution specified by the forward problem results. This assumption corresponds to the Gaussian model approach introduced in subsection 3.2.2. Finally we obtain an estimate about the gene regulatory relationships,

$$\omega_{\text{mlGt}} = \arg\max_{\omega} \left(-\frac{1}{2} \sum_{\mu k} \left(\mathbf{x}_{k}^{\mu} - \mathbf{m}^{\mu}(\omega) \right) \chi^{\mu}(\omega)^{-1} \left(\mathbf{x}_{k}^{\mu} - \mathbf{m}^{\mu}(\omega) \right) -\frac{1}{2} \ln \det \left(\chi^{\mu}(\omega) \right) \right),$$
(3.79)

by maximisation of the likelihood function, $p({\mathbf{x}^{\mu}} | \omega, \mathbf{a}, \mathbf{b}, \mathbf{c})$, with respect to ω .

3.5.3. Comparison of least squares and maximum likelihood method

Within this subsection, we compare the results of the least squares methods and our MLM. For the comparison we employ scatter plots of generated and reconstructed interactions as well as the reconstruction error as a quantitative precision measure. We find that the likelihood-based approach outperforms the other methods in the regime of a strong stochastic contribution to the system dynamics.

In this subsection, data is generated by drawing independent random samples from the steady state of our stochastic model of gene regulation. For inference, we use samples from the steady state of different perturbations, μ .

Inference based on our GT is not scalable to large system sizes due to an iterative calculation with numerical integration in each iteration step. Therefore we focus on the inference of small subnetworks. Such small subnetworks are realised by regulatory pathways, which are of great importance in the formation and spread of cancer [36, 40].

To visualise the precision of reconstruction, the difference between reconstructed, ω_{ij}^{rec} , and generated interaction, ω_{ij}^{gen} , we employ scatter plots. In figure 3.8 we show scatter plots of an inferred interaction matrix in a regime of a significant stochastic contribution to system dynamics, $c_i = 1.0$, based on 50 samples per perturbation. All four methods tend to overestimate the gene regulatory interactions. In the case of our MLM, an overestimation is to be expected because of the Gaussian distribution used for the generation of the interaction matrix and the flat prior distribution used in the likelihood-based inference. In summarising section 5.1, we discuss the use of prior information in the context of GRNI. For $c_i = 1.0$, we find that the likelihood-based approach provides a significantly more accurate reconstruction of the generated interaction matrix compared to the least squares methods.

In the Appendix figure I.1 we show scatter plots of an inferred interaction matrix in a regime of a small stochastic contribution to system dynamics, $c_{\rm i} = 0.1$. We find that even with only 5 samples per perturbation the inference in the regime of a small stochastic contribution is significantly more precise. For $c_{\rm i} = 0.1$ wo do not find a significant difference in reconstruction

accuracy between the four employed approaches.

To quantify the reconstruction precision, we employ a quadratic reconstruction error,

$$r = \sqrt{\frac{\sum_{ij} \left(\omega_{ij}^{\text{gen}} - \omega_{ij}^{\text{rec}}\right)^2}{\sum_{ij} \left(\omega_{ij}^{\text{gen}}\right)^2}},$$
(3.80)

analogously to the forward problem. The reconstruction error is plotted against the samples per perturbation in figure 3.9. In subplot 3.9b the inference takes place in a regime of significant stochastic contribution to the system dynamics, $c_{\rm i} = 1.0$, whereas in subplot 3.9a there is only a small stochastic contribution, $c_{\rm i} = 0.1$.

As it is to be expected all methods yield a precise reconstruction in the limit of a large number of samples. The reconstruction is for all four employed methods significantly less precise in the regime of a large stochastic contribution to system dynamics.

Nevertheless, GT within the likelihood-based approach consistently outperforms the other approaches in the regime of a significant stochastic contribution, $c_{\rm i} = 1.0$. With a smaller stochastic contribution, $c_{\rm i} = 0.1$, there is no significant difference between the approaches.

In the regime of a large stochastic contribution, $c_i = 1.0$, we find significant information for the network reconstruction in the second moments of the samples. This information is visualised in the difference of least squares approach in first- and second-order in figure 3.9b. The impact of the probabilistic samples model in the regime of a large stochastic contribution can be estimated by comparing the results of the GT within the likelihood-based approach and the least squares approach based on the GT 3.9b. We find that the gain in precision based on the information in second moments and the impact of the probabilistic model in the likelihood function is of the same order of magnitude.

The least square methods in first and second order are only exact in the limit of a large number of samples per perturbation, because they are based on steady state averages. The maximum likelihood estimate is also valid for a small number of samples per perturbation. In the limit of only a few samples per perturbation, we expect the likelihood-based approach to outperform the least squares methods. We do not find such a significant difference as can be seen in figure 3.9a. It is not clear why approaches based on the steady state averages work in the limit of only a few samples per steady state without a significant loss of accuracy compared to our MLM.

Figure 3.8.: Scatter plot of generated and reconstructed interactions for a fully connected system as small as 10 nodes. The model parameter are set to $a_i = 1$, $b_i = 1$, and $c_i = 1$. For the inference 10 distinct single drug perturbations with 50 samples per perturbation are used.



Figure 3.9.: The reconstruction error is plotted against the samples per perturbation in a fully connected system as small as 10 nodes and set of model parameter $a_i = 1$, $b_i = 1$, and $c_i = 0.1$ in subfigure 3.9a such as $c_i = 1.0$ in subfigure 3.9b. For the reconstruction we use 10 distinct single drug perturbations. We employ the least squares fit in first order, \bigcirc , in second order, \bigcirc , and based on the GT, \bigcirc . We employ furthermore our MLM, \bigcirc .



3.6. Inference and response prediction in a melanoma cell line

We focus in this section on the inference based on experimental perturbation data. We employ the four approaches introduces in the previous section based on measurements of signalling activity in the melanoma cell line SK-MEL-133. We find that a prediction based on inferred model parameter gives an accurate estimate of unknown signalling activity. For our likelihood-based approach, we show a graph representation of the inferred interaction matrix and compare the reconstructed signalling network to the literature.

Before we start with the response prediction we outline the experimental setup in subsection 3.6.1. We divide the perturbation data up into training and prediction set. Based on inference within the training set we predict signalling activity for the prediction set in subsection 3.6.2. In the closing subsection 3.6.3 we employ our likelihood-based approach on the whole dataset and discuss the inferred regulatory interactions.

3.6.1. Perturbation experiment

The perturbation experiment in the SK-MEL-133 cell line is published in the supporting information of the publication [76]. The SK-MEL-133 melanoma cell line has functional mutations within the MAPK pathway, which is discussed in the motivational section 3.1 as an important pathway in the development and spread of melanoma.

The SK-Mel-133 cells are singularly and pairwise perturbed using a set of 8 inhibitors, which predominantly target the PI3K/AKT and MAPK pathway. As an inhibitory concentration the IC_{40} value is used to generate a gentle perturbation with a measurable effect. The IC_{40} value is the minimal concentration required to reduce the activity of the targeted protein by 40%, which can be estimated from a dose-response curve.

Within a set of 16 proteins, protein phosphorylation levels are measured. This set of proteins does not include the perturbed proteins, which is a fundamental difference to the inference based on simulated data in the previous section. For each of the 44 perturbations, three independent biological replicates are measured and response is quantified by a logarithmic ratio between perturbed and unperturbed phosphorylation levels.

Phosphorylation levels are measured with reverse-phase protein arrays (RPPA). RPPA is a microarray technology. DNA microarray technologies are outlined in section 3.2.1. The RPPA technology is designed for the simultaneous measurement of protein concentration and phosphorylation state in a large number of biological samples [77].

Apart from the measurement of phosphorylation level, the cell viability after drug perturbation is quantified 72 hours after perturbation in a resazurin assay. Resazurin is a weakly blue fluorescent, cell-permeable substance. It is irreversibly reduced to the pink-coloured and highly fluorescent resorufin by metabolic cell activity. Therefore the reduction of resazurin is a widely used indicator of cell proliferation and viability.

3.6.2. Response prediction

We divide the 44 perturbations up into training sets with 33 perturbations and corresponding prediction sets consisting out of 11 perturbations. Such that there are in total four pairs of distinct training and prediction sets. Based on training data we infer the gene interaction matrix, ω , and the model parameter **a**, **b**, and **c**.

We employ the four inference approaches introduced in the previous section. Compared to the previous section, we enlarge the parameter space by the set of model parameter **a**, **b**, and **c**. Matrix elements, ω_{ij} , are set equal to zero that quantify the regulatory effect on the perturbed proteins and for
which there are no measurements of signalling activity. Such that $\omega_{ij} = 0$ for all i labelling the perturbed proteins.

There exist regions of the parameter space where the GT does not converge. Moreover, we find that the success of our MLM depends on the starting point of the local optimisation. We obtain reliable predictions by using the result of the least squares fit based on GT as a starting point for our MLM.

Given the inferred interaction matrix and model parameter we predict mean phosphorylation level, \mathbf{m}^{μ} , within the prediction set using a numerical simulation based on our stochastic model.

For the first pair of training and prediction set, the predicted signalling activities are compared to the measured ones in figure 3.10. The predictions based on the other training sets are depicted in the Appendix section I.3.

We find that phosphorylation levels can be predicted based on our stochastic model. Based on all four sets of training and prediction data, we obtain for the least squares method in first order a mean square error of $r_{\rm ls1o} = 0.8 \pm 1.2$ and in second order of $r_{\rm ls2o} = 1.3 \pm 2.1$. Whereas the mean square error of $r_{\rm lsGt} = 0.086 \pm 0.020$ for the least squares Gaussian theory approach and our likelihood based method $r_{\rm mlGt} = 0.083 \pm 0.014$ are significant smaller.

3.6.3. Network reconstruction

For the network reconstruction, we employ our likelihood-based method, which we introduce in subsection 3.5.2, based on the whole dataset. To incorporate the knowledge that most of the proteins do not interact, we use the heavy-tailed Laplace distribution, as a sparsity-favouring prior distribution. We employ the Laplace distribution with mean, $\nu = 0$, and variance, $2\sigma^2 = 1/2$ as a prior for the matrix elements, ω_{ij} . In this way, ω_{ij} is set equal to zero whenever there is no clear evidence in the data for an interaction. On the other hand ω_{ij} which are required are allowed to be sizeable, because the Laplace distribution decreases slightly with $\exp(-|x|)$, whereas the normal distribution decreases with $\exp(-x^2)$. We do not take additional prior information for the parameter **a**, **b**, and **c** into account and assume a flat prior distribution for these parameters. Finally, we obtain a maximum Figure 3.10.: Scatter plot of predicted and measured signalling activity. The inference is based on a training set, which does not include the predicted perturbations. We employ the least squares fit in first order, ●, in second order, ●, and based on the GT, ●. We employ furthermore our MLM, ●. We give the mean squared distances between predicted and measured signalling activity.



Table 3.1.: Confusion table between interaction from the Reactome pathway database [78] and inferred interaction based on our MLM. The reference interaction are classified as activating, non-interacting, inhibiting, and undirected interaction.

ref. inf.	act.	non- int.	inh.	und. int.
$\omega_{ij} \ge 0.2$	0	12	0	2
$-0.2 \geq \omega_{\rm ij} \geq 0.2$	26	175	2	20
$\omega_{\rm ij} \leq -0.2$	1	14	1	3

posterior estimate based on the GT,

$$\{\omega, \mathbf{a}, \mathbf{b}, \mathbf{c}\}_{\mathrm{mpGt}} = \underset{\{\omega, \mathbf{a}, \mathbf{b}, \mathbf{c}\}}{\mathrm{arg\,max}} \left(-\frac{1}{2} \sum_{\mu \mathbf{k}} \left(^{\mathrm{data}} \mathbf{x}_{\mathbf{k}}^{\mu} - \mathbf{m}^{\mu}(\omega, \mathbf{a}, \mathbf{b}, \mathbf{c}) \right)^{\mathrm{T}} \\ (\chi^{\mu}(\omega, \mathbf{a}, \mathbf{b}, \mathbf{c}))^{-1} \left(^{\mathrm{data}} \mathbf{x}_{\mathbf{k}}^{\mu} - \mathbf{m}^{\mu}(\omega, \mathbf{a}, \mathbf{b}, \mathbf{c}) \right) \\ -\frac{1}{2} \ln \det \left(\chi^{\mu}(\omega, \mathbf{a}, \mathbf{b}, \mathbf{c}) \right) - \sum_{ij} \frac{|\omega_{ij} - \nu|}{\sigma} \right),$$
(3.81)

for the model parameter.

In figure 3.11 we give a graph representation of the inferred subnetwork of the measured proteins. We employ a lower limit of 0.2 on the absolute value for the representation of an regulatory interaction.

Our network is more complex than most signalling cascades in the literature. Nonetheless, we compare our findings, the inferred signalling network, to the manually curated, open access Reactome pathway database [78]. The regulatory interactions within the database are represented in figure 3.12. In table 3.1 we compare our inferred interactions with the data base. We find that the inferred regulatory interactions do not represent known regulatory relationships.

The high ratio between the number of proteins and the sample size, and probably also the experimental noise, make the inference of biochemical interactions without integration of prior knowledge infeasible. Figure 3.11.: The graph represents our maximum posterior estimate for the signalling network. We plot inferred regulatory relationships with $|\omega_{ij}| \ge 0.2$. The graph representation is created with Cytoscape [79].







4. Cancer immunotherapy response prediction

Science, for me, gives a partial explanation for life. In so far as it goes, it is based on fact, experience and experiment.

Rosalind Franklin

We aim to identify patients who will likely benefit from checkpoint blockade immunotherapy (CBI) to support clinical decision-making. Within this chapter, we address the question of whether one can predict cancer immunotherapy response based on frameshift mutations which result in random peptide sequences that are entirely different from self-peptides.

We introduce cancer immunotherapy by CBI in section 4.1 and review a wide range of approaches to CBI response prediction. To investigate the information about CBI response contained within frameshift mutation, we employ statistical classification and survival analysis. We outline the mathematical foundations of the employed statistical methods in section 4.3 and introduce hypothetical frameshift-based response determinants in section 4.2.

Response evaluation criteria are the basis for our statistical classification in section 4.4. We quantify the predictive power of frameshift mutation on the overall survival of patients treated with CBI in section 4.5. For this purpose, we sort the patients into two groups based on their mutational profile and compare the survival rates.

We summarise our findings on the immunogenic potential of frameshiftderived peptides in section 4.6. We find slight evidence that frameshift mutations are related to immunotherapy response. Nonetheless, our statistical analysis revealed that frameshift-derived peptides are not significantly associated with immunotherapy response. Our findings are compatible with a hidden factor, e.g. the mutation rate, that increases the number of unknown immunogenic mutations and the number of frameshift mutations. Still, there is no evidence that the frameshift mutations are causal for immunotherapy response.

4.1. Cancer immunotherapy by immune checkpoint blockade

An immune checkpoint is a pathway that regulates the immune response. This regulation is central to immune tolerance, e.g. insensitivity to the microbiome, immune tolerance in pregnancy and tolerance to non-self peptides in food. During tumour evolution, some tumours acquire the ability to stimulate immune checkpoints that prevent an immune response. CBI blocks inhibitory checkpoints such as PD-1, PD-L1, and CTLA-4 to restore the protective function of the immune system [80].

Extensive clinical trials show that within a minority of patients, checkpoint inhibition has outstanding benefits [12]. Although immunotherapy generally has fewer adverse effects than chemotherapy, CBI can cause severe adverse reactions by altering immunologic self-tolerance. Because only a small patient group receives durable clinical benefits, the high variability in patient response limits the clinical use. Therefore the understanding of determinants that drive immune response, resistance, and adverse side effects is a key scientific issue in the field of immuno-oncology [2].

In the following, we provide an overview of CBI response determinants that have been studied in the literature. We refer to the expression of tumour suppressor genes, the tumour genome, host germline genetics and the tumour microenvironment.

Expression of the immune suppressor gene PD-L1 is a predictive biomarker of CBI response. Improved efficacy of CBI over chemotherapy, with fewer adverse effects than in chemotherapy, was found in patients with advanced non-small lung cancer (NSCLC) and PD-L1 expression [81]. Within another phase III clinical trial, CBI was not associated with significantly longer survival times than chemotherapy among NSCLC patients with PD-L1 expression [82]. Although a metastudy indicates that PD-L1 expression is a predictive biomarker, the PD-L1 expression has limitations, and further determinants have to be carefully investigated [3].

Tumour genomes contain self-antigens that are a hypothetical determinant for CBI response [83]. These non-mutated antigens are based on tumour overexpressed genes or genes, only expressed in cells without MHC presentation within healthy tissue. Thus, self-antigens are potential targets for immune recognition.

Besides non-mutated self-antigens, tumour genomes have antigens based on non-synonymous mutations. Studies show that there is a correlation between the overall number of mutations and the response to CBI. This correlation is likely linked to these somatic mutations [4]. The overall number of mutations is called tumour mutational burden (TMB). TMB-based biomarkers are in clinical use. However, in a pan-cancer dataset of more than 2500 CBI-treated patients, Mirny et al. found little evidence that TMB is a predictive biomarker. The statistical analysis of this dataset suggests that previously reported correlations are based on confounding cancer subtypes and incorrect statistical testing [5].

Within the tumour genome, microsatellite instability (MSI) is a further response determinant. MSI is a predisposition to mutations in repetitive DNA sequences due to a non-functioning DNA repair mechanism. Thus, the MSI consequently leads to a high TMB. Efficacy of PD-1 blockade therapy is found in cancer patients with MSI across 12 different tumour types [84].

Frameshift mutation derived peptides are highly distinct from self-peptides. Thus frameshift mutations in the tumour genome are hypothetically a rich source of immunogenic antigens. Across three melanoma studies, Swanton et al. found that the number of frameshift mutations is significantly associated with response to CBI [6].

The MHC molecules are encoded by the human leukocyte antigen (HLA) complex. The HLA system is the most heterogeneous gene complex within the human genome. MHC diversity is a critical component for immune defence. Therefore, a more diverse set of MHC molecules is a hypothetical response determinant. An analysis of more than 1500 CBI-treated patients gives evidence that a more diverse set of MHC molecules within a patients genome is related with clinical benefit [85]. The ability to present a broader range of peptides via MHC molecules on the cell surface may explain this finding.

The expression of inhibitory checkpoints is associated with the density and distribution of $CD8^+$ T cells within the tumour microenvironment. The presence of $CD8^+$ T cells located at the tumour margin before starting therapy may determine CBI response in metastatic melanoma [86].

Cancer is a systemic disease that causes far-reaching immune system dysfunction. Thus, the entire immune system can play a crucial role in the tumour macroenvironment. Accumulating evidence indicates that CBI drives new immune responses rather than enhancing pre-existing ones [87]. The understanding of systematic immunity in cancer may be an essential step for reliable response prediction.

In summary, there are numerous hypothetical determinants of response to CBI. Despite further clinical studies, we need careful analysis of genomic and systemic characteristics of CBI response.

4.2. Hypothetical frameshift-based response determinants

In our study, we consider that a high number of frameshift mutations is not the whole story. Their expression is a precondition for immune recognition. Frameshift derived peptides must be trimmed for MHC presentation and bind to an MHC molecule. Finally, the immune system must recognise the peptides as non-self for a frameshift mutation to be immunogenic. In addition to the number of frameshift mutations, we focus on additional hypothetical frameshift-related determinants and, in particular, on frame shift-derived peptides. We consider observables without free parameters and search for a clear response signal.

We introduce the safeguard mechanisms nonsense-mediated decay (NMD) and non-stop decay (NSD) that prevent the translation of erroneous mRNA in subsection 2.2.2. NMD prevents the translation of frameshift sequences that contain a stop codon 50 nucleotides before an exon-exon junction. NSD safeguards against frameshift-derived mRNA sequences that do not contain a proper stop codon. Taking into account these safeguard mechanisms, we consider the number of frameshift-derived peptides that are not affected by NMD and NSD.

We expect the clonal structure of the tumour to be relevant for CBI response. A clonal mutation is likely to cause a more significant immune response than a non-clonal mutation. In the case of an immune reaction due to a non-clonal mutation, some tumour clones are resistant to the immune reaction. A study within a mouse model indicates that neoantigens with a low cancer cell fraction (CCF) do not lead to immune-mediated cell rejection and that the CCF threshold for an immune response is antigen-dependent [88]. We assume that the relative measure of a CCF is relevant for immune response. We define the accumulated frequency of frameshift mutations,

$$f_{\rm i} = \sum_{\rm m \in FM_{\rm i}} {\rm CCF_{\rm m}} \,, \tag{4.1}$$

where FM_i is the set of frameshift mutations within the tumour genome of individual i and CCF_m is the CCF of mutation m.

The expression of a frameshift mutation is a candidate for a CBI response determinant. We investigate the accumulated expression of frameshift mutations,

$$e_{\rm i} = \sum_{\rm m \in FM_{\rm i}} {\rm GE_m} \,, \tag{4.2}$$

where GE_m is the expression of gene m.

MHC presentation relies on trimming proteins into peptides of canonical 8 – 10 amino acids. Peptide processing significantly shapes the MHC immunopeptidome, the set of peptides presented on MHC molecules. However, the specificity of this process is not well understood, and the prediction of peptide processing is not entirely reliable [89]. We, therefore, do not consider the specificity of peptide trimming.

We investigate the length of frameshift-derived peptides,

$$l_{i} = \sum_{m \in FM_{i}} L_{m}, \qquad (4.3)$$

where L_m denotes the length of a frameshift derived peptide. In Appendix section II.1, we explain our bioinformatic method for obtaining frameshift-derived peptide sequences.

The presentation of frameshift-derived peptides on MHC molecules on the cell surface is another hypothetical determinant of CBI response. Therefore we investigate the number of MHC-binding frameshift-derived peptides,

$$b_{i} = \sum_{m \in FM_{i}} B_{m} , \qquad (4.4)$$

where we estimate peptide-MHC (pMHC) binding via NetMHC [90, 91]. We restrict the prediction to peptides with a lenght of 9 amino acids. Because we do not have information about a patient-specific MHC repertoire, we base our pMHC binding prediction on supertype representatives. The first set of supertype representatives was defined in 1999 [92]. Nowadays, it is possible to cover the binding properties of almost all known MHC molecules based on functional binding specificities of a few supertype representatives [93].

The adaptive immune response depends on an interaction between the T cell receptor (TCR) and the pMHC complex. An essential feature in the characterisation of the TCR-pMHC interaction is the binding affinity, the electronic property by which the TCR and the pMHC are prone to form a chemical compound. However, to predict the T cell recognition with the TCR-pMHC affinity is far from conclusive [94]. Maybe the information transmitted to the T Cell, measured by the entropy of TCR-pMHC binding dynamics, is decisive for recognition [95]. Thus, we do not consider TCR-pMHC affinity within our study.

4.3. Mathematical foundations of statistical classification and survival analysis

We base our CBI response prediction on information about genetic alterations within the patient's genome. For this purpose, we employ statistical classification and survival analysis. In subsection 4.3.1, we introduce the framework of statistical classification. We employ statistical classification in the following section to build a stochastic model predicting binary CBI response. We try to divide the set of patients into responders and non-responders based on information about their genetic alterations.

In subsection 4.3.2, we introduce the framework of survival analysis, which is a branch of statistics for the analysis of time-to-event data, such as the overall survival after the start of therapy.

4.3.1. Statistical classification

Before we introduce the method of Bayesian logistic regression, we formulate the statistical classification problem for a binary output variable. The classification problem and the method of Bayesian logistic regression are introduced in Bishop's excellent textbook "Pattern Recognition and Machine Learning" [96]. To quantify the performance of a classifier, we define a quality measure in the closing paragraph.

Binary classification problem

Binary classification in the context of our CBI response prediction is the task of classifying patients into responders and non-responders based on tumour genetic information. Cancer tissue samples are successfully classified by their gene expression patterns for clinical decision-making [97].

Within our CBI response prediction, we base the classification of a patient on an input variable, \mathbf{x}^n , with patient index n, and a classification rule. The \mathbf{x}^n characterise the mutational landscape of a patient with a real-valued vector. The classification rule is a mapping to a binary output variable, $\mathbf{x}^n \mapsto k^n$ with $k^n \in \{1, 0\}$. To construct a classification rule, we employ the framework of Bayesian logistic regression, which is introduced in the following paragraph.

The mapping can be characterised by the true positive rate (TPR) and the false positive rate (FPR). The TPR is a measure of sensitivity. It is the ratio between the number of patients correctly classified as responders and the total number of patients classified as responders. The FPR is the probability of falsely classifying an actual non-responder. Thus the FPR is calculated as the ratio between the number of non-responders wrongly categorised as responders and the total number of actual non-responder.

Bayesian logistic regression

Bayesian logistic regression is based on data consisting of independent input variables, \mathbf{x}^{n} , and dependent binary output variables, $r^{n} \in \{1, 0\}$.

It is assumed that the dependent outcome variable is Bernoulli distributed, such that

$$p(r^{n}) = (p^{n})^{r^{n}} (1 - p^{n})^{1 - r^{n}} .$$
(4.5)

For the response variable, p^n , one assumes a linear relationship between the input variables and the log-odds of p^n ,

$$\ln\left(\frac{p^{n}}{1-p^{n}}\right) = \alpha + \sum_{i} \beta_{i} x_{i}^{n} .$$

$$(4.6)$$

To solve the inverse problem, to get an unbiased estimate for the model parameters α and β_i , we employ statistical inference.

As a prior for the model parameters, we assume a normal distribution centred at zero with variance σ^2 . We employ an inverse gamma distribution as a conjugate prior for the variance of the normal distribution as explained in subsection 2.1.3. The posterior distribution of the model parameters is determined by the stochastic model

$$\begin{aligned} \sigma^{2} &\sim \text{InverseGamma}(2,3) \\ \alpha &\sim \text{Normal}(0,\sigma^{2}) \\ \beta_{i} &\sim \text{Normal}(0,\sigma^{2}) \\ r^{n} &\sim \text{Bernoulli}\left(p^{n}\right) \text{ with } p^{n} \stackrel{(4.5)}{=} \frac{1}{1 + \exp\left(-\alpha - \sum_{i} \beta_{i} x_{i}^{n}\right)}. \end{aligned}$$
(4.7)

Based on the data, we obtain a maximum a posteriori probability (MAP) estimate for the model parameter, α^* and β_i^* . To this end, we draw a sequence of random samples with a Hamiltonian Monte Carlo algorithm [98] implemented within the probabilistic programming library Turing.jl [99]. The sample sequence converges to be distributed according to the posterior distribution.

To solve the forward problem, the classification of a patient based on mutational information, we define a score function

$$s^{n}(\mathbf{x}^{n}) = \alpha^{\star} + \sum_{i} \beta_{i}^{\star} x_{i}^{n} \,. \tag{4.8}$$

Within this score function we use our maximum posterior estimate of the model parameter. Our response classification, $\mathbf{x}^{n} \mapsto k^{n}$ with

$$k^{n} = \begin{cases} 1, & \text{if } s^{n}(\mathbf{x}^{n}) \ge \delta \\ 0, & \text{else} \end{cases},$$
(4.9)

depends on the score function, $s^{n}(\mathbf{x}^{n})$, and a variable threshold parameter, δ .

Receiver operating characteristic

The receiver operating characteristic (ROC) is a method to quantify the performance of a threshold dependent binary classifier. The ROC was developed in the context of signal detection and gained application in various quantitative research fields. Nowadays, it is frequently used in medical decision making and as a performance measure in machine learning. Within this chapter, we employ the ROC curve as a graphical representation of the performance of a classifier. The area under this curve is a quantitative measure for classification precision. A rigorous analysis of the ROC can be found in the literature [100].

The FPR and the TPR, introduced with the binary classification problem, span the two-dimensional ROC space. Both rates characterise a binary classifier and depend on the threshold of the classifier, δ . One generates the ROC curve by plotting the TPR against FPR for each value of δ . Figure 4.1 shows the ROC curves of four classifiers with different predictive power.

To quantify the quality of a binary classifier, one can employ the area under the curve (AUC). One assumes that the value of the score function s is distributed according to the probability density $f_{\rm R}(s)$ for patients with clinical response and according to $f_{\rm N}(s)$ respectively for patients without response. Under this hypothesis, TPR and FPR can be represented as an integral,

$$\begin{aligned} \text{TPR}(\delta) &= \int_{\delta}^{\infty} f_{\text{R}}(s) \text{d}s \\ \text{FPR}(\delta) &= \int_{\delta}^{\infty} f_{\text{N}}(s) \text{d}s \,. \end{aligned} \tag{4.10}$$

For a randomly chosen responder, \mathbf{x}^{n} , and a randomly chosen non-responder, \mathbf{x}^{m} , the AUC is equal to the probability that the inequality $s^{n} > s^{m}$ holds,

$$\begin{aligned} \text{AUC} &= \int_{0}^{1} \text{TPR} \left(\text{FPR}^{-1}(x) \right) \mathrm{d}x \\ & \stackrel{\text{FPR}^{-1}(x) = :\delta}{=} -\int_{-\infty}^{\infty} \text{TPR}(\delta) \frac{\partial \text{FPR}}{\partial \delta}(\delta) \, \mathrm{d}\delta \\ & \stackrel{(4.10)}{=} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f_{\text{R}}(s) f_{\text{N}}(\delta) \theta(s-\delta) \, \mathrm{d}s \, \mathrm{d}\delta = p(s^{\text{n}} > s^{\text{m}}) \,. \end{aligned}$$
(4.11)

Figure 4.1.: Examples of ROC curves for different predictive classifier. Random classifier \bigcirc , weak classifier \bigcirc , strong classifier \bigcirc , perfect classifier \bigcirc .



4.3.2. Survival analysis

Survival analysis is a branch of statistics to analyse the timespan until one event occurs, such as a death in CBI trials. At the end of observation time, either an event occurred, or the patient was censored. Censoring takes place if no event is observed and nothing is known about the patient after observation. Within the context of our CBI response prediction, the observation time starts with treatment. Censoring occurs when the patient leaves the clinical study on account of any reason without observing an event. An event may or may not occur after censoring.

We define common terms in survival analysis and introduce the Kaplan-Meier estimator to obtain patient groups' survival rates. In the closing paragraph, we outline the method of Bayesian linear regression, which we employ for our survival time prediction.

Common terms in survival analysis

We refer to survival data as a dataset consisting of observation times, $t^{n} \in \mathbb{R}^{+}$, and information about censoring, $c^{n} \in \{1, 0\}$, where $c^{n} = 1$ marks patients leaving the clinical study. The survival function,

$$s^{n}(t) = P(t^{n} > t),$$
 (4.12)

gives the probability that a patient will survive beyond time t.

To estimate the survival function of a patient group, we employ a Kaplan-

Meier estimator. The Kaplan-Meier estimator of survival rates

$$\hat{s}(t) = \prod_{i:t_i < t} (1 - \frac{d_i}{n_i})$$
(4.13)

depends on d_i , the number of deaths that happened at t_i , and n_i , the patients at risk which have not yet died or been censored up to time t_i . By comparing $\hat{s}(t)$ between patient groups, one gets evidence for differential efficacy of therapy. The Kaplan-Meier estimator is a widely used method for demonstrating clinical benefit.

The interested reader can find an extensive mathematical description of survival analysis in the lecture notes [101].

Bayesian linear regression

To quantify the relation between mutational information, \mathbf{x}^n , and survival time, t^n , we employ the framework of Bayesian linear regression. On account of unknown event times we have to deal with right-censored data. Rightcensored time spans end later than a certain point in time but it is unknown by how much. The general assumption of our approach is a linear relation between x_i^n and t^n . We employ a stochastic model,

$$\begin{split} &\sigma^2 \sim \text{InverseGamma}(2,3) \\ &\alpha \sim \text{Normal}(0,\sigma^2) \\ &\beta_{\text{i}} \sim \text{Normal}(0,\sigma^2) \\ &t^{\text{n}} \sim \begin{cases} \text{Normal}(\mu^{\text{n}},1), & \text{if } c^{\text{n}} = 0 \\ \text{CCDFNormal}(\mu^{\text{n}},1), & \text{if } c^{\text{n}} = 1 \end{cases}, \end{split}$$

$$\end{split}$$

where the inverse-gamma distribution is used as a prior for the variance analogously to our stochastic model of logistic regression. We assume the uncensored survival times to be distributed according to a normal distribution with mean

$$\mu^{\mathbf{n}}(\mathbf{x}^{\mathbf{n}}) = \alpha + \sum_{\mathbf{i}} \beta_{\mathbf{i}} x_{\mathbf{i}}^{\mathbf{n}}$$
(4.15)

and unit variance. We employ the complementary cumulative distribution function (CCDF) of the normal distribution for the right-censored data. We generate samples with a Hamiltonian Monte Carlo algorithm using the probabilistic programming library Turing.jl [99]. If we do not fix the variance, our stochastic model is more unstable, and the Hamilton Monte Carlo algorithm does not converge. For survival time prediction we employ equation (4.15) with a MAP estimate for the model parameter, α^* and β_i^* .

Table 4.1.: Map from RECIST to a binary outcome variable

	Miao	Liu
response	clinical benefit	complete response & partial response
no-response	stable disease & no clinical benefit	stable disease & progressive disease

4.4. Statistical classification

We employ binary classification to investigate the information content in hypothetical CBI response determinants. For this purpose, we map the discrete response evaluation criteria in solid tumours (RECIST) to a binary outcome variable, response and no-response. The mapping is shown in table 4.1.

In this and the following section, we analyse datasets published by Miao et al. [102] and Liu et al. [103]. The Miao dataset is based on whole-exome sequencing of 249 tumours and healthy tissue from patients with clinically evaluated responses to CBI across multiple cancer types. The Liu dataset contains a cohort of 144 melanoma patients treated with anti-PD1 CBI. These datasets contain patient-specific information about genetic mutations, clinical information about therapy's effects, and the overall survival after the start of CBI. The information about various types of genetic mutations are based on variant calling, the identification of mutations from sequencing data [104]. We chose these datasets because of the relatively large number of patients, their public accessibility, and the included information about frameshift mutations.

Within the first subsection of this section, we use the ROC to investigate hypothetical response determinants. To predict the outcome of CBI, we employ the framework of Bayesian logistic regression within the second subsection.

4.4.1. Hypothetical response determinants

Given the binary response, we compare the predictive value of frameshift mutations with other mutation types. Therefore, we study the number of mutations as a classifier within a series of ROC curves, which is depicted in figure 4.2. We do not find a significant difference in the predictive value between the number of missense, nonsense, frameshift, and silent mutations. This finding is also valid within a second melanoma dataset in the Appendix

Figure 4.2.: ROC curve based on the number of a specific mutation type as classifier. We examine the number of missense ●, nonsense ●, frameshift ●, and silent mutation ●.



subfigure II.1a. In comparison between cancer types, we find that the number of mutations has a higher predictive value for patients with lung cancer than melanoma patients.

All ROCs in figure 4.2 are comparable and are in particular similar with silent and nonsense mutations. Thus, the data are compatible with a hidden factor, e.g. the mutation rate, which increases the number of all mutations independent of whether they are a hypothetical target for an immune response. The silent and nonsense mutations are not immunogenic, but they still have a signal via this implicit mechanism.

In the second series of ROC curves, we investigate hypothetical response determinants related to frameshift mutations, which we defined in section 4.2. The ROC curves are shown in figure 4.3. In addition to the underlying signal of the number of frameshift mutations, we do not find an increased predictive value within the investigated frameshift-related response determinants. We also find no response signal in the second melanoma data set, shown in the attached subfigure II.1b.

4.4.2. Response prediction

We use statistical inference to make a statement about the significance of the classification based on the response determinants studied. For this purpose, we employ our stochastic response model defined in subsection 4.3.1. We focus on the number of missense and frameshift mutations as input variables for our binary response prediction. Within the lung cancer dataset, we found

Figure 4.3.: ROC curve for frameshift mutation based classifier. We examine the number of frameshift-derived peptides, ●, the accumulated frequency ●, the accumulated expression ●, the length of frameshift peptides ●, and number of MHC-binding frameshift-derived peptides ●.



a small response signal. Therefore we focus on the response prediction within this dataset.

In subfigure 4.4a we show a scatter plot of the drawn samples for the regression coefficients, β_1 and β_2 , corresponding to the missense and frameshift mutations. The number of missense mutations is generally much larger. Thus the axis scaling is not comparable.

We do not find a significant response signal. The null hypothesis, $\forall i : \beta_i = 0$, is within the 1σ credible region. Therefore, there is no evidence at the 1σ level that CBI response depends on the number of missense and frameshift mutations.

We employ our stochastic response model to predict binary CBI response based on our MAP estimate of the model parameter. For this purpose, we employ cross-validation. We divide the data into two sets, one prediction set, consisting of three individuals, and a training set, with the other 54 patients. The response score for each patient relies on a maximum likelihood estimate based on the corresponding training set. The resulting cross-validated ROC curve is plotted in subfigure 4.4b. A comparison with subfigure 4.2a shows that the prediction based on our stochastic response model is not significantly better than a classification based on the number of frameshift mutations. This finding is compatible with the hypothesis that a generally high mutation right increases both the number of unknown immunogenic mutations and the number of frameshifts. Still, there is no evidence that the

- Figure 4.4.: We employ Bayesian logistic regression for CBI response prediction based on the number of missense and frameshift mutations.
- (a) Samples from the posterior, ●,
 MAP estimate, ●, 1σ credible region, ●, 2σ credible region, ●.
- (b) Cross-validated ROC curve based on a MAP estimate within our stochastic model.



frameshift mutations are causal for CBI response.

4.5. Survival analysis

We employ survival analysis to investigate hypothetical determinants of overall survival after the start of CBI. Analogously to the previous section, we explore the predictive power of hypothetical survival determinants in subsection 4.5.1 and try to predict the overall survival in subsection subsection: Survival prediction.

4.5.1. Hypothetical survival determinants

To quantify the predictive power of hypothetical survival determinants, we divide the set of patients according to a real-valued score function, representing a hypothetical response determinant, into a high-score and a low-score subset. Within our datasets, we find that about 30% of the patients are classified as responders. We use this prior knowledge for the subset size. We calculate Kaplan-Meier estimators and visualise the estimated survival rates

in a Kaplan-Meier plot.

We explore the predictive value of different mutation types within a lung cancer and melanoma dataset in the figures 4.5 and 4.6. We show the corresponding Kaplan-Meier curves for a second melanoma dataset in the appended figure II.2.

We find that the patient group with more frameshift mutations has a higher survival rate within the Miao lung cancer data. There is no significant difference in overall survival after the start of CBI between the high-score and the low-score subset in the Miao melanoma data. It is unclear whether the number of frameshift mutations is the cause of CBI response or whether there is a common cause for the frameshift mutation and the overall survival. In contrast to the binary response prediction, we do not find a response signal of silent mutations in overall survival within the lung cancer dataset.

We investigate whether the number of frameshift mutations has a predictive value apart from the overall number of mutations. We, therefore, divide the whole dataset into two subsets of equal size according to the overall mutation number and divide each of these subsets again into two subsets of equal size according to the frameshift mutation number. The corresponding ROC curves are shown in figure 4.7 and for the second melanoma dataset in the Appendix figure II.3. Within the lung cancer dataset, we find slight evidence that the number of frameshift mutations in the subset characterised by a high overall number of mutations has some predictive value. The melanoma datasets also support this minor trend.





Figure 4.6.: Kaplan-Meier curves based on hypothetical survival determinants for melanoma patients within the Miao dataset [102], high score subset, ●, low score subset, ●.



Figure 4.7.: Kaplan-Meier curves for patient groups characterised by high overall and high frameshift, ●, high overall and low frameshift, ●, low overall and high frameshift, ●, low overall and low frameshift, ● number of mutation.



4.5.2. Survival prediction

We try to predict the overall survival after the start of CBI within the lung cancer dataset analogously to the binary response prediction. For this purpose, we employ our stochastic survival model defined in subsection 4.3.1 to divide the set of patients into a high-score and a low-score subset. We based our prediction on the number of missense and frameshift mutations.

The posterior distribution of the model parameter is characterised within subfigure 4.8a. Based on our MAP estimate, we employ our stochastic survival model to predict overall survival. We use the same cross-validation scheme as in our binary response prediction. The resulting Kaplan-Meier curves, which are based on our MAP estimate, are plotted in subfigure 4.8b. We find that the high-score patient group has a higher survival rate than the low-score group. Notwithstanding, this prediction is not significantly better than a prediction solely based on the number of frameshift mutations.

- Figure 4.8.: We employ Bayesian linear regression for CBI survival prediction based on the number of missense and frameshift mutations.
 - (a) Samples from the posterior, \bigcirc , MAP estimate, \bigcirc , 1σ credible region, \bigcirc , 2σ credible region, \bigcirc .
- (b) Kaplan-Meier curve based on a MAP survival prediction within our linear model.



4.6. Conclusion about the immunogenic potential of frameshift mutation

Our statistical analysis of CBI data revealed that the number of frameshift mutations is not significantly related to the response class and the overall survival across three clinical cohorts. Furthermore, we also investigated the number of frameshift-derived peptides concerning safeguard mechanisms (NMD and NSD), accumulated frequency based on CCF, accumulated gene expression, length of frameshift-derived peptides, and the number of MHCbinding frameshift-derived peptides. We also found no significant CBI response signal in these hypothetical frameshift-related determinants.

Nonetheless, we find little evidence that the number of frameshift mutations is associated with response classes and overall survival after the start of CBI. We can make no statement about a causal relationship between the emergence of frameshift mutation and immunotherapy success. It is beyond the scope of the investigated datasets whether frameshift-derived peptides have a significant physical impact on CBI response. We find some evidence that a hidden factor, e.g. the mutation rate, increases both the number of unknown immunogenic mutations and the number of frameshift mutations. Still, there is no evidence that the frameshift mutations are causal for immunotherapy response.

Our stochastic models, based on frameshift and missense mutations, do neither improve cross-validated response nor survival prediction. We find no evidence of additional information from multiple mutation types, which is compatible with a common hidden factor.

5. Conclusion

Science never solves a problem without creating ten more.

George Bernard Shaw

Within this chapter, we conclude our research, combining statistical mechanics and information theory with high-throughput technologies, to address open questions in systems biology.

In section 5.1, we summarise our findings on network inference and biological information processing. We review our results on the frameshift-based prediction of CBI response in section 5.2.

5.1. Likelihood-based gene regulatory network inference

Our proposed system of stochastic differential equations based on an interaction network (3.4) proves to be a promising model for the research on biological information processing. We employ our dynamic model to construct a maximum likelihood estimate for the interaction network based on gene expression data.

Based on exact steady state relations (3.12) we can solve the forward problem within the MFT and our GT. We find that the GT outperforms MFT in the regime of strong inter-gene couplings, where the mean field approximation breaks down.

Based on the GT, we propose an MLM to solve the inverse problem. In the regime of a sizeable stochastic contribution to the system dynamics, we find that our MLM outperforms standard least squares fits, which are successfully used for the reconstruction of GRNs [76, 8]. Within our MLM, we implicitly assume a flat prior distribution, which causes an overestimation of interactions. We find this overestimation also within the least squares methods. In the case of the GT, the use of a prior distribution can solve this straightforwardly. Whereas in the case of the least squares methods, one needs an additional term in the cost function that penalises the network complexity [76]. We can predict unknown signalling activity levels in the SK-MEL-133 cell line with all investigated methods. We show that the approaches based on the GT yield a significantly smaller mean squared error than the least squares methods.

We find that the inferred regulatory interactions do not represent known regulatory relationships. This issue constitutes a lack of interpretability because we constructed an accurate prediction system with hidden internal logic.

A decisive difference between simulated and experimental perturbation data is the missing availability and specificity of drugs, such that one can not systematically perturb all genes within an experiment. Even with sufficient drugs, one is limited by the financial and temporal cost of testing many drug combinations with multiple biological replicates per perturbation. Moreover, in perturbation experiments, the number of perturbations is typically as large as the system size with only a few samples per perturbation [7, 76, 8].

A high ratio between the number of genes and the sample size makes the inference of biological knowledge without integrating prior knowledge infeasible. Therefore, taking present-day data availability into account, biological prior information that may be incorrect in a particular physical context is crucial for GRNI. We note that quantitative GRNI is usually based on information from the literature or interaction databases [105, 106, 107].

Our MLM should be helpful in open questions about gene regulation and biological information processing. One can extend our approach into several future research directions.

Research literature and online databases offer a vast amount of biological knowledge about gene regulation and signalling networks. The integration of biological knowledge is challenging due to programmatic access, various gene name conventions, and quantification of prior information. Nonetheless, our MLM offers a straightforward way to incorporate prior knowledge to achieve reliable network reconstruction.

The investigation of the posterior landscape could be another promising starting point for further research. Based on the posterior distribution, one could quantify the uncertainty of inferred regulatory relationships and propose promising drug combinations to test uncertain interactions.

We used protein concentration and phosphorylation data from a cell-line perturbation experiment. However, incorporating additional information on mRNA concentration could be a starting point for the inference of multi-level gene regulation. By combining mRNA and protein data, one can expect to generate whole-cell models of gene regulation and signalling pathways that support the design of clinical trials.

Finally, the investigation of hidden nodes within the GRN could be a future research direction. We base our network inference on the convenient but unrealistic hypothesis that all the relevant gene expression levels are measured. In practice, it is impossible to be sure that there are no other interacting genes. At the cost of a larger parameter space, one could extend our MLM by allowing some hidden nodes.

5.2. Frameshift-based cancer immunotherapy response prediction

We address the potential of statistical inference for clinical decision-making in the context of cancer immunotherapy. Our statistical analysis of CBI data revealed that the number of frameshift mutations is not significantly associated with response to CBI. Nonetheless, we find slight evidence that frameshift mutations are related to CBI response. Our findings are compatible with a hidden factor, e.g. the mutation rate, that increases both the number of unknown immunogenic mutations and the number of frameshifts. Still, there is no evidence that the frameshift mutations are causal for CBI response.

We found that cross-validation is essential within the investigation of hypothetical determinants to CBI response. Free parameters tend to improve the descriptiveness at the cost of predictiveness. Based on noisy data within a small dataset, even a statistical model with few model parameters is prone to overfitting. Such a model will probably describe a data set well but will usually fail in making predictions. In addition to the model parameter selection, this also applies to a subset choice based on cancer subspecies or other clinical factors. There is evidence that confounding of cancer subtypes and incorrect statistical tests lead to previously reported response determinants [5].

During the timespan of our research on CBI response prediction, a promising marker for melanoma is found. The protein midkine (MDK) is a driver of an inflamed, but immune-evasive tumour microenvironment that is correlated with resistance to CBI in melanoma patients [108]. The study [108] links MDK expression with poor CBI outcome and points out a promising combined MDK immune checkpoint inhibition. Moreover, the study stresses the importance of a systemic investigation of CBI response.

The understanding of molecular and cellular drivers of immune escape is one of the biggest challenges to move the field of cancer immunotherapy forward [109]. The identification of immunogenic mutations is an open research question. Extensive CBI trials and statistical inference can provide knowledge about predictive biomarkers to improve clinical decision-making.

Bibliography

- [1] M. Hecker et al. "Gene regulatory network inference: Data integration in dynamic models - A review". In: *BioSystems* 96.1 (2009).
- [2] J. J. Havel, D. Chowell, and T. A. Chan. "The evolving landscape of biomarkers for checkpoint inhibitor immunotherapy". In: *Nature Reviews Cancer* 19.3 (2019).
- [3] A. A. Davis and V. G. Patel. "The role of PD-L1 expression as a predictive biomarker." In: *Journal for ImmunoTherapy of Cancer* 7 (2019).
- [4] T. Chan et al. "Development of tumor mutation burden as an immunotherapy biomarker: utility for the oncology clinic". In: Annals of Oncology 30.1 (2019).
- [5] C. Gurjao et al. "Limited evidence of tumour mutational burden as a biomarker of response to immunotherapy". In: *bioRxiv* 260265 (2020).
- [6] C. Swanton et al. "Insertion-and-deletion-derived tumour-specific neoantigens and the immunogenic phenotype: a pan-cancer analysis." In: *The Lancet. Oncology* 18.8 (2017).
- S. Nelander et al. "Models from experiments: Combinatorial drug perturbations of cancer cells". In: *Molecular Systems Biology* 4.216 (2008).
- [8] A. Korkut et al. "Perturbation biology nominates upstream downstream drug combinations in RAF inhibitor resistant melanoma cells". In: *eLife* 8.4 (2015).
- [9] A. Raj and A. van Oudenaarden. "Stochastic gene expression and its consequences". In: *Cell* 135.2 (2008).
- [10] V. Chandrasekaran, N. Srebro, and P. Harsha. "Complexity of Inference in Graphical Models". In: arXiv 1206.3240 (2012).
- [11] S. Klamt and J. Stelling. "Combinatorial Complexity of Pathway Analysis in Metabolic Networks". In: *Molecular Biology Reports* 29.1 (2002).
- [12] A. Ribas and J. D. Wolchok. "Cancer immunotherapy using checkpoint blockade". In: *Science* 359.6382 (2018).

- [13] R. Weinberg. The Biology of Cancer. W.W. Norton, 2013.
- [14] L. Sompayrac. How the Immune System Works. The How it Works Series. Wiley, 2019.
- [15] F. S. Gnesotto et al. "Broken detailed balance and non-equilibrium dynamics in living systems: a review". In: *Reports on Progress in Physics* 6 (2018).
- [16] H. C. Nguyen, R. Zecchina, and J. Berg. "Inverse statistical problems: from the inverse Ising problem to data science". In: Advances in Physics 66.3 (2017).
- [17] E. T. Jaynes. "Information Theory and Statistical Mechanics". In: *Phys. Rev.* 106 (1957).
- [18] E. Jaynes and G. Bretthorst. Probability Theory: The Logic of Science. Cambridge University Press, 2003.
- [19] C. Shannon. "A mathematical theory of communication". In: Bell System Technical Journal 27.3 (1949).
- [20] D. Applebaum. "Lévy Processes—From Probability to Finance and Quantum Groups". In: Notices of the American Mathematical Society 51 (2004).
- [21] J. Blath, P. Imkeller, and S. Rœlly. Surveys in Stochastic Processes. EMS series of congress reports. European Mathematical Society, 2011.
- [22] P. Langevin. "Sur la théorie du mouvement brownien". In: *Comptes rendus de l'Académie des sciences* (1908).
- [23] K. Itô. "On Stochastic Differential Equations". In: American Mathematical Society (1951).
- [24] C. Gardiner. Handbook of stochastic methods for physics, chemistry, and the natural sciences. Springer, 1985.
- [25] F. Kelly. *Reversibility and Stochastic Networks*. Wiley Series in Tracts on Probability and Statistics. J. Wiley, 1979.
- [26] D. J. C. Mackay. Information Theory, Inference and Learning Algorithms. Cambridge University Press, 2003.
- [27] F. F. Costa. "Non-coding RNAs: Lost in translation?" In: Gene 386.1 (2007).
- [28] G. Cooper. *Regulation of Transcription in Eukaryotes.* The Cell: A Molecular Approach. Sinauer Associates, 2000.

- [29] R. G. H. Lindeboom et al. "The impact of nonsense-mediated mRNA decay on genetic disease, gene editing and cancer immunotherapy". In: *Nature Genetics* 51.11 (2019).
- [30] BioRender. 639 Queen Street West, Toronto. 2021. URL: https://biorender.com/.
- [31] M. Mézard and J. Sakellariou. "Exact mean-field inference in asymmetric kinetic Ising systems". In: *Journal of Statistical Mechanics* 7 (2011).
- [32] C. Collinet and T. Lecuit. "Programmed and self-organized flow of information during morphogenesis". In: Nature Reviews Molecular Cell Biology 22.4 (2021).
- [33] P. W. H. Holland. "The future of evolutionary developmental biology". In: *Nature* 402.6761 (1999).
- [34] F. Emmert-Streib, M. Dehmer, and B. Haibe-Kains. "Gene regulatory networks and their applications: understanding biological and medical problems in terms of networks". In: *Frontiers in Cell and Developmental Biology* 2 (2014).
- [35] WHO. "Monitoring Health for the Sustainable Sevelopment Goals". In: World Health Statistics (2020).
- [36] D. Hanahan and R. A. Weinberg. "The hallmarks of cancer." In: *Cell* 100.1 (2000).
- [37] K. Wang, S. I. Grivennikov, and M. Karin. "Implications of anticytokine therapy in colorectal cancer and autoimmune diseases". In: *Annals of the Rheumatic Diseases* 72.2 (2013).
- [38] J. Downward. "Targeting RAS signalling pathways in cancer therapy". In: *Nature Reviews Cancer* 3.1 (2003).
- [39] G. L. Semenza. "Targeting HIF-1 for cancer therapy". In: *Nature Reviews Cancer* 3.10 (2003).
- [40] D. Hanahan and R. A. Weinberg. "Hallmarks of cancer: The next generation". In: *Cell* 144.5 (2011).
- [41] D. Schadendorf and A. Hauschild. "Melanoma the run of success continues". In: *Nature Reviews Clinical Oncology* 11 (2014).
- [42] M. Adams et al. "Complementary DNA Sequencing: Expressed Sequence Tags and Human Genome Project". In: *Science* 252 (1991).
- [43] A. Athar et al. "ArrayExpress update from bulk to single-cell expression data". In: *Nucleic Acids Research* 47.D1 (2018).

- [44] A. D. Yates et al. "Ensembl 2020". In: Nucleic Acids Research 48.D1 (2019).
- [45] A. Schulze and J. Downward. "Navigating gene expression using microarrays - a technology review Upstream considerations: microarray technology". In: *Nature Cell Biology* 3.8 (2001).
- [46] J. D. Hoheisel. "Microarray technology: Beyond transcript profiling and genotype analysis". In: *Nature Reviews Genetics* 7.3 (2006).
- [47] Z. Wang, M. Gerstein, and M. Snyder. "RNA-Seq: a revolutionary tool for transcriptomics". In: *Nature Reviews Genetics* 10.1 (2010).
- [48] Y. Chu and D. R. Corey. "RNA Sequencing: Platform Selection, Experimental Design, and Data Interpretation". In: *Nucleic Acid Therapeutics* 22.4 (2018).
- [49] M. Wu and A. K. Singh. "Single-cell protein analysis". In: Current Opinion in Biotechnology 23.1 (2012).
- [50] I. C. Macaulay, C. P. Ponting, and T. Voet. "Single-Cell Multiomics: Multiple Measurements from Single Cells". In: *Trends in Genetics* 33.2 (2017).
- [51] S. Ambardar et al. "High Throughput Sequencing: An Overview of Sequencing Chemistry". In: *Indian journal of microbiology* 56.4 (2016).
- [52] A. Mortazavi et al. "Mapping and quantifying mammalian transcriptomes by RNA-Seq". In: *Nature Methods* 5.7 (2008).
- [53] P. E. Meyer et al. "Information-theoretic inference of large transcriptional regulatory networks". In: *Eurasip Journal on Bioinformatics* and Systems Biology 10.1155 (2007).
- [54] A. A. Margolin et al. "ARACNE: An Algorithm for the Reconstruction of Gene Regulatory Networks in a Mammalian Cellular Context". In: *BMC Bioinformatics* 7.1 (2006).
- [55] D. Koller and N. Friedman. Probabilistic Graphical Models: Principles and Techniques - Adaptive Computation and Machine Learning. The MIT Press, 2009.
- [56] S. M. Hill et al. "Bayesian Inference of Signaling Network Topology in a Cancer Cell Line". In: *Bioinformatics* 28.21 (2012).
- [57] S. Nee et al. "Inferring Cellular Networks". In: *Science* 303.2 (2004).
- [58] B. Schoeberl et al. "Computational modeling of the dynamics of the MAP kinase cascade activated by surface and internalized EGF receptors". In: *Nature Biotechnology* 20.4 (2002).

- [59] B. B. Aldridge et al. "Physicochemical modelling of cell signalling pathways". In: *Nature Cell Biology* 8.11 (2006).
- [60] R. Bonneau et al. "The Inferelator: an algorithm for learning parsimonious regulatory networks from systems-biology data sets de novo". In: Genome Biology 7.5 (2006).
- [61] Y. Li and S. A. Jackson. "Gene Network Reconstruction by Integration of Prior Biological Knowledge". In: Genes Genetics 5.6 (2015).
- [62] J. Hertz, A. Krogh, and R. G. Palmer. Introduction to the Theory of Neural Computation. Redwood City: Addison-Wesley, 1991.
- [63] M. B. Elowitz et al. "Stochastic gene expression in a single cell". In: Science 297.5584 (2002).
- [64] E. Ising. "Beitrag zur Theorie des Ferromagnetismus". In: Zeitschrift für Physik 31 (1925).
- [65] S. Kobe. "Ernst Ising 1900-1998". In: Brazilian Journal of Physics 30 (2000).
- [66] A. C. C. Coolen. "Statistical Mechanics of Recurrent Neural Networks I. Statics". In: arXiv cond-mat.0006010 (2000).
- [67] M. Opper and D. Saad. Advanced Mean Field Methods: Theory and Practice. Neural information processing series. MIT Press, 2001.
- [68] M. Mezard, G. Parisi, and M. Virasoro. Spin Glass Theory And Beyond: An Introduction To The Replica Method And Its Applications. World Scientific Lecture Notes In Physics. World Scientific Publishing Company, 1987.
- [69] H. J. Kappen and J. J. Spanjers. "Mean field theory for asymmetric neural networks". In: *Physical Review E* 61.5 (2000).
- [70] B. H. Callen. Thermodynamics and an Introduction to Thermostatistics. Wiley, 1985.
- [71] S. Umarov, C. Tsallis, and S. Steinberg. "On a q-central limit theorem consistent with nonextensive statistical mechanics". In: *Milan Journal of Mathematics* 76.1 (2008).
- [72] H. J. Hilhorst. "Central limit theorems for correlated variables: some critical remarks". In: *Brazilian Journal of Physics* 39.2A (2009).
- [73] D. Laurie. "Calculation of Gauss-Kronrod quadrature rules". In: Math. Comput. 66 (1997).
- [74] Steven G. Johnson. The NLopt nonlinear-optimization package. 2021. URL: http://github.com/stevengj/nlopt.

- [75] M. Powell. "The BOBYQA algorithm for bound constrained optimization without derivatives". In: *NA Report* 6 (2009).
- [76] E. J. Molinelli et al. "Perturbation Biology: Inferring Signaling Networks in Cellular Systems". In: *PLoS Computational Biology* 9.12 (2013).
- [77] R. Tibes et al. "Reverse phase protein array: Validation of a novel proteomic technology and utility for analysis of primary leukemia specimens and hematopoietic stem cells". In: *Molecular Cancer Therapeutics* 5 (2006).
- [78] J. Griss et al. "ReactomeGSA Efficient Multi-Omics Comparative Pathway Analysis". In: *Molecular & Cellular Proteomics* 19.12 (2020).
- [79] Cytoscape Consortium. 2021. URL: https://cytoscape.org/.
- [80] D. M. Pardoll. "The blockade of immune checkpoints in cancer immunotherapy". In: *Nature Reviews Cancer* 12 (2012).
- [81] M. Reck et al. "Pembrolizumab versus Chemotherapy for PD-L1 Positive Non – Small-Cell Lung Cancer". In: New England Journal of Medicine 375.19 (2016).
- [82] D. P. Carbone et al. "First-Line Nivolumab in Stage IV or Recurrent Non-Small-Cell Lung Cancer". In: New England Journal of Medicine 376.25 (2017).
- [83] M. F. Gjerstorff, M. H. Andersen, and H. J. Ditzel. "Oncogenic cancer/testis antigens: prime candidates for immunotherapy". In: Onco-Target 6.18 (2015).
- [84] D. T. Le et al. "Mismatch repair deficiency predicts response of solid tumors to PD-1 blockade". In: *Science* 357.6349 (2017).
- [85] D. Chowell et al. "Patient HLA class I genotype influences cancer response to checkpoint blockade immunotherapy". In: Science 359.6375 (2018).
- [86] P. C. Tumeh et al. "PD-1 blockade induces responses by inhibiting adaptive immune resistance". In: *Nature* 515.7528 (2014).
- [87] K. J. Hiam-Galvez, B. M. Allen, and M. H. Spitzer. "Systemic immunity in cancer". In: *Nature Reviews Cancer* 21.6 (2021).
- [88] R. S. Gejman et al. "Rejection of immunogenic tumor clones is limited by clonal fraction". In: *eLife* 7 (2018).
- [89] M. Y. Lee et al. "Antigen processing and presentation in cancer immunotherapy". In: *Journal for immunotherapy of cancer* 8.2 (2020).
- [90] M. Andreatta and M. Nielsen. "Gapped sequence alignment using artificial neural networks: application to the MHC class I system". In: *Bioinformatics* 32.4 (2016).
- [91] M. Nielsen et al. "Reliable prediction of T-cell epitopes using neural networks with novel sequence representations". In: *Protein science : a publication of the Protein Society* 12.5 (2003).
- [92] A. Sette and J. Sidney. "Nine major HLA class I supertypes account for the vast preponderance of HLA-A and-B polymorphism". In: *Immunogenetics* 50.3 (1999).
- [93] J. Sidney et al. "HLA class I supertypes: a revised and updated classification". In: *BMC Immunology* 9.1 (2008).
- [94] J. Gálvez, J. J. Gálvez, and P. García-Peñarrubia. "Is TCR/pMHC Affinity a Good Estimate of the T-cell Response? An Answer Based on Predictions From 12 Phenotypic Models". In: *Frontiers in immunology* 10 (2019).
- [95] J. R. Egan, T. Elliott, and B. D. MacArthur. "Fluctuations in TCR and pMHC interactions regulate T cell activation". In: *bioRxiv* 430441 (2021).
- [96] C. Bishop. *Pattern Recognition and Machine Learning*. Information Science and Statistics. Springer, 2006.
- [97] A. C. Tan et al. "Simple Decision Rules for Classifying Human Cancers from Gene Expression Profiles". In: *Bioinformatics Oxford jour*nal 21 (2005).
- [98] M. Betancourt. "A Conceptual Introduction to Hamiltonian Monte Carlo". In: *arXiv* 1701.02434 (2017).
- [99] H. Ge, K. Xu, and Z. Ghahramani. "Turing: a language for flexible probabilistic inference". In: International Conference on Artificial Intelligence and Statistics 21 (2018).
- [100] T. Fawcett. "An introduction to ROC analysis". In: *Pattern Recogni*tion Letters 27.8 (2006).
- [101] D. Bakry, R. D. Gill, and S. A. Molchanov. Lectures on Probability Theory. Ecole d'Ete de Probabilites de Saint-Flour, 1992.
- [102] D. Miao et al. "Genomic correlates of response to immune checkpoint blockade in microsatellite-stable solid tumors". In: *Nature Genetics* 50 (2018).
- [103] D. Liu et al. "Integrative molecular and clinical modeling of clinical outcomes to PD1 blockade in patients with metastatic melanoma". In: *Nature Medicine* 25 (2019).

- [104] R. Nielsen et al. "Genotype and SNP calling from next-generation sequencing data". eng. In: *Nature reviews. Genetics* 12.6 (2011).
- [105] P. Phillip, A. Bahl, and L. Ungar. "Using Prior Knowledge to Improve Genetic Network Reconstruction from Microarray Data". In: *In silico biology* 4 (2004).
- [106] M. Ghanbari, J. Lasserre, and M. Vingron. "Reconstruction of gene networks using prior knowledge". In: BMC Systems Biology 9.1 (2015).
- [107] D. Altarawy, F.-E. Eid, and L. S. Heath. "PEAK: Integrating Curated and Noisy Prior Knowledge in Gene Regulatory Network Inference". In: Journal of Computational Biology 24.9 (2017).
- [108] D. Cerezo-Wallis et al. "Midkine rewires the melanoma microenvironment toward a tolerogenic and immune-resistant state". In: *Nature Medicine* 26.12 (2020).
- [109] P. S. Hegde and D. S. Chen. "Top 10 Challenges in Cancer Immunotherapy". In: *Immunity* 52.1 (2020).
- [110] Genome Reference Consortium Human Build 37 patch release 13. 2013. URL: https://www.ncbi.nlm.nih.gov/assembly.
- [111] S. Durinck et al. "Mapping identifiers for the integration of genomic datasets with the R/Bioconductor package biomaRt". In: *Nature Protocols* 4.8 (2009).
- [112] K. L. Howe et al. "Ensembl 2021". In: Nucleic Acids Research 49.D1 (2020).

I. Gene regulatory network inference

Within this Appendix chapter, we explain technical details and show supplementary figures regarding out network inference.

I.1. Parallel Glauber dynamics

Within parallel Glauber dynamics the spins, \mathbf{s} , evolve in discrete time with synchronous updates. Dependent on the spin configuration, s(t), the spin configuration in the next time step, s(t + t), is a random sample according to the probability distribution

$$p(s(t+1) \mid s(t)) = \prod_{i} \frac{\exp\left(\beta s_{i}(t+1)h_{i}(t)\right)}{2\cosh\left(\beta h_{i}(t)\right)}.$$
 (I.1)

One obtains for the first moment of the steady state distribution the relation

$$\begin{aligned} \langle s_{\mathbf{i}}(t+1) \rangle \stackrel{(\mathrm{I.1})}{=} \langle \tanh(h_{\mathbf{i}}(t)) \rangle \\ \Rightarrow \langle s_{\mathbf{i}} \rangle_{\mathbf{s}} &= \langle \tanh(h_{\mathbf{i}}) \rangle_{\mathbf{s}} , \end{aligned}$$
(I.2)

where in any time step all spins are updated according to (I.1). The averages become time independent in the steady state as it is the case for sequential Glauber dynamics. Analogously, the relationship for the second moment is as follows

$$\left\langle s_{i}(t+1)s_{j}(t+1)\right\rangle \stackrel{(I.1)}{=} \left\langle \tanh(h_{i}(t))\tanh(h_{j}(t))\right\rangle$$

$$\Rightarrow \left\langle s_{i}s_{j}\right\rangle_{s} = \left\langle \tanh(h_{i})\tanh(h_{j})\right\rangle_{s}.$$
 (I.3)

I.2. Mean field theory

In this Appendix section, we calculate second order contributions to mean gene expression m_i^{μ} and covariance of gene expression χ_{ij}^{μ} within the mean field approximation.

I.2.1. Mean gene expression in second order MFT

To calculate the quadratic terms in the Taylor expansion of $m_i^{\mu} = \int dx \left[p\left(x|\theta\omega\right) \tanh\left(h_i^{\mu}\right) \right]$, we use partial derivatives,

$$\frac{\partial^2 \tanh\left(h_{\rm i}^{\mu}\right)}{\partial \theta_{\rm j} \partial \theta_{\rm k}} \bigg|_{q^*} = -2m_{\rm i}^{\mu q} \left(1 - \left(m_{\rm i}^{\mu q}\right)^2\right) \delta_{\rm ij} \delta_{\rm ik} \tag{I.4}$$

$$\frac{\partial^2 \tanh\left(h_{\rm i}^{\mu}\right)}{\partial \omega_{\rm kl} \partial \theta_{\rm j}}\Big|_{q^*} = -2m_{\rm i}^{\mu q} \left(1 - \left(m_{\rm i}^{\mu q}\right)^2\right) \delta_{\rm ij} \delta_{\rm ik} x_{\rm l}^{\mu} \tag{I.5}$$

$$\frac{\partial^2 \tanh\left(h_{\rm i}^{\mu}\right)}{\partial \omega_{\rm jk} \partial \omega_{\rm lm}} \bigg|_{q^*} = -2m_{\rm i}^{\mu q} \left(1 - \left(m_{\rm i}^{\mu q}\right)^2\right) \delta_{\rm ij} \delta_{\rm il} x_{\rm k}^{\mu} x_{\rm m}^{\mu} \,, \tag{I.6}$$

evaluated at the factorizing distribution q^* . With partial differentiation and integration over the space of gene expression we obtain second order Taylor coefficients,

$$\begin{split} \frac{\partial m_{i}^{\mu}}{\partial \theta_{j} \partial \theta_{k}} \bigg|_{q^{*}} &= \int dx \left[\frac{\partial^{2} p\left(x | \theta \omega\right)}{\partial \theta_{j} \partial \theta_{k}} \tanh\left(h_{i}^{\mu}\right) \right]_{q^{*}} \\ &+ \int dx \left[p\left(x | \theta \omega\right) \underbrace{\frac{\partial^{2} \tanh\left(h_{i}^{\mu}\right)}{\partial \theta_{j} \partial \theta_{k}}}_{\substack{(I.4) \\ = -2m_{i}^{\mu q}\left(1 - \left(m_{i}^{\mu q}\right)^{2}\right) \delta_{ij} \delta_{ik}} \right]_{q^{*}} \\ &+ \int dx \left[\frac{\partial p\left(x | \theta \omega\right)}{\partial \theta_{j}} \frac{\partial \tanh\left(h_{i}^{\mu}\right)}{\partial \theta_{j}} \right]_{q^{*}} + (j \leftrightarrow k) \\ &= 0 - 2m_{i}^{\mu q} \left(1 - \left(m_{i}^{\mu q}\right)^{2}\right) \delta_{ij} \delta_{ik} + 0 \end{split}$$
(I.7)

104

$$\begin{split} \frac{\partial^2 m_i^{\mu}}{\partial \theta_j \partial \omega_{kl}} \bigg|_{q^*} &= \int dx \left[\frac{\partial^2 p \left(x | \theta \omega \right)}{\partial \theta_j \partial \omega_{kl}} \tanh \left(h_i^{\mu} \right) \right]_{q^*} \\ &+ \int dx \left[p \left(x | \theta \omega \right) \underbrace{\frac{\partial^2 \tanh \left(h_i^{\mu} \right)}{\partial \theta_j \partial \omega_{kl}}}_{\substack{(1.5) - 2m_i^{\mu q} \left(1 - (m_i^{\mu q})^2 \right) \delta_{il} \delta_{ik} x_m^{\mu}}} \right]_{q^*} \\ &+ \int dx \left[\frac{\partial p \left(x | \theta \omega \right)}{\partial \theta_j} \underbrace{\frac{\partial \tanh \left(h_i^{\mu} \right)}{\partial \omega_{kl}}} \right]_{q^*} \\ &+ \int dx \left[\frac{\partial p \left(x | \theta \omega \right)}{\partial \omega_{kl}} \underbrace{\frac{\partial \tanh \left(h_i^{\mu} \right)}{\partial \theta_j}} \right]_{q^*} \\ &= 0 - 2m_i^{\mu q} \left(1 - (m_i^{\mu q})^2 \right) m_i^{\mu} \delta_{ij} \delta_{ik} \\ &+ \left(1 - (m_i^{\mu q})^2 \right) \left(1 - (m_1^{\mu q})^2 \right) \delta_{ik} \delta_{lj} + 0 \end{split}$$

$$\begin{aligned} \frac{\partial^2 m_i^{\mu}}{\partial \omega_{jk} \partial \omega_{lm}} \bigg|_{q^*} &= \int dx \left[\frac{\partial^2 p \left(x | \theta \omega \right)}{\partial \omega_{jk} \partial \omega_{lm}} \tanh \left(h_i^{\mu} \right) \right]_{q^*} \\ &+ \int dx \left[\frac{p \left(x | \theta \omega \right)}{\partial \omega_{jk} \partial \omega_{lm}} \tanh \left(h_i^{\mu} \right) \right]_{q^*} \\ &+ \int dx \left[\frac{\partial p \left(x | \theta \omega \right)}{\partial \omega_{jk} \partial \omega_{lm}} \frac{\partial^2 \tanh \left(h_i^{\mu} \right)}{\partial \omega_{lm}} \right]_{q^*} + \left(jk \leftrightarrow lm \right) \\ &= 0 + \left(-2 \right) m_i^{\mu q} \left(1 - (m_i^{\mu q})^2 \right) \delta_{il} \delta_{il} \left(m_k^{\mu q} m_m^{\mu q} + \frac{1}{2} \delta_{km} \right) \\ &+ \left(1 - (m_i^{\mu q})^2 \right) \left(1 - (m_m^{\mu q})^2 \right) \delta_{il} \delta_{il} m_k^{\mu \mu} + \left(jk \leftrightarrow lm \right) . \end{aligned}$$

In the third term of (I.9) we identify the first order derivertive $\frac{\partial m_{\rm m}^{\mu}}{\partial \omega_{\rm jk}}\Big|_{q^*}$, which we have already calculated in equation (3.39). This a repetitive pattern, in the calculation of higher order MFT corrections one can identify contributions of lower order. Summing over all contributions in quadratic order in $\delta \omega$ and $\delta \theta^{\mu}$, we obtain the second order corrections in mean gene expression,

$${}^{\theta\theta}m_{\mathbf{i}}^{\mu q^*} := \sum_{\mathbf{jk}} \left. \frac{\partial^2 m_{\mathbf{i}}^{\mu}}{\partial \theta_{\mathbf{j}} \partial \theta_{\mathbf{k}}} \right|_{q^*} \delta\theta_{\mathbf{j}} \delta\theta_{\mathbf{k}} = (-2)m_{\mathbf{i}}^{\mu q} \left(1 - \left(m_{\mathbf{i}}^{\mu q}\right)^2 \right) \delta\theta_{\mathbf{i}}^2 \qquad (\mathbf{I}.10)$$

105

$$\begin{split} ^{\omega\omega}m_{i}^{\mu q^{*}} &:= \sum_{jklm} \left. \frac{\partial^{2}m_{i}^{\mu}}{\partial \omega_{jk} \partial \omega_{lm}} \right|_{q^{*}} \delta \omega_{jk} \delta \omega_{lm} \\ &= \left(1 - \left(m_{i}^{\mu q}\right)^{2} \right) \left(\left(-2\right)m_{i}^{\mu q} \left(-\theta_{i}^{MFT1}\right)^{2} + \left(-2\right)m_{i}^{\mu q} \frac{1}{2} \sum_{j} \left(\delta \omega_{ij}\right)^{2} \right. \\ &+ 2 \sum_{m} \delta \omega_{im} \left(1 - \left(m_{m}^{\mu q}\right)^{2} \right) \left(-\delta \theta_{m}^{MFT1}\right) \right) \,. \end{split}$$
(I.12)

Based on the condition to recover the mean gene expression also in second order we obtain the requirement

$$0 \stackrel{!}{=} {}^{\theta}m_{i}^{\mu q^{*}} + {}^{\omega}m_{i}^{\mu q^{*}} + \frac{1}{2}{}^{\theta\theta}m_{i}^{\mu q^{*}} + {}^{\theta\omega}m_{i}^{\mu q^{*}} + \frac{1}{2}{}^{\omega\omega}m_{i}^{\mu q^{*}} + \sigma\left(\delta^{3}\right) \qquad (I.13)$$

on the corrections up to quadratic order. Using the second order corrections (3.45), (3.46) and (3.47) as well as the previous results (3.40) and (3.41), we obtain an analytical expression for $\delta\theta$ in second order mean field approximation,

$$\delta\theta_{i}^{\rm MFT2} = \delta\theta_{i}^{\rm MFT1} + \frac{1}{2}m_{i}^{\mu}\sum_{k}\left(\delta\omega_{ik}\right)^{2}.$$
 (I.14)

Finally, analogous to the first-order approximation, we calculate the mean gene expression in second order MFT,

$$m_{\rm i}^{\mu\,\rm MFT2} = \tanh\left(\sum_k \omega_{\rm ik} m_{\rm k}^{\mu} + \theta_{\rm i} - \frac{1}{2} m_{\rm i}^{\mu} \sum_k \left(\omega_{\rm ik}\right)^2\right)\,.\tag{I.15}$$

This result is based on the equations (3.32) and (3.33).

106

I.2.2. Covariance of gene expression in second order MFT

For the calculation of second order contributions to the covariance, we evaluate partial derivatives,

$$\frac{\partial^{2} \left(\tanh \left(h_{i}^{\mu} \right) - m_{i}^{\mu} \right)}{\partial \theta_{k} \partial \omega_{lm}} \bigg|_{q*} \stackrel{(I.8)(I.5)}{=} - 2m_{i}^{\mu q} \left(1 - \left(m_{i}^{\mu q} \right)^{2} \right) \delta_{ik} \delta_{il} x_{m}^{\mu q} x_{m}^{\mu q} - \frac{\partial^{2} m_{i}^{\mu}}{\partial \theta_{k} \partial \omega_{lm}} \bigg|_{q*} \tag{I.16}$$

$$\frac{\partial^{2} \left(\tanh\left(h_{i}^{\mu}\right) - m_{i}^{\mu}\right)}{\partial \omega_{kl} \partial \omega_{mn}} \Big|_{q*} \stackrel{(I.6)(I.9)}{=} - 2m_{i}^{\mu q} \left(1 - \left(m_{i}^{\mu q}\right)^{2} \right) \delta_{ik} \delta_{im} \left(x_{l}^{\mu q} x_{n}^{\mu q}\right) \\
- \frac{\partial^{2} m_{i}^{\mu}}{\partial \omega_{kl} \partial \omega_{mn}} \Big|_{q*},$$
(I.17)

at the factorizing distribution q^* . As in the previous calculation, we get second order Taylor coefficients,

$$\begin{split} \frac{\partial^2 \chi_{ij}^{\mu}}{\partial \theta_k \partial \theta_l} \Big|_{q^*} &= \frac{1}{2} \int dx \left[\frac{\partial^2 p\left(x | \theta \omega \right) \left(x_j^{\mu} - m_j^{\mu} \right)}{\partial \theta_k \partial \theta_l} \left(\tanh\left(h_i^{\mu} \right) - m_i^{\mu} \right) \right]_{q^*} \\ &+ \frac{1}{2} \int dx \left[p\left(x | \theta \omega \right) \left(x_j^{\mu} - m_j^{\mu} \right) \underbrace{\frac{\partial^2 \left(\tanh\left(h_i^{\mu} \right) - m_i^{\mu} \right)}{\partial \theta_k \partial \theta_l}}_{(1.6)(1.9)_0} \right]_{q^*} \\ &+ \frac{1}{2} \int dx \left[\frac{\partial p\left(x | \theta \omega \right) \left(x_j^{\mu} - m_j^{\mu} \right)}{\partial \theta_l} \frac{\partial \left(\tanh\left(h_i^{\mu} \right) - m_i^{\mu} \right)}{\partial \theta_k} \right]_{q^*} \\ &+ \frac{1}{2} \int dx \left[\frac{\partial p\left(x | \theta \omega \right) \left(x_j^{\mu} - m_j^{\mu} \right)}{\partial \theta_k} \frac{\partial \left(\tanh\left(h_i^{\mu} \right) - m_i^{\mu} \right)}{\partial \theta_l} \right]_{q^*} + (i \leftrightarrow j) \\ &= 0 \end{split}$$
(I.18)

$$\begin{split} \frac{\partial^2 \chi_{ij}^{\mu}}{\partial \theta_k \partial \omega_{lm}} \bigg|_{q^*} &= \frac{1}{2} \int dx \left[\frac{\partial^2 p \left(x | \theta \omega \right) \left(x_j^{\mu} - m_j^{\mu} \right)}{\partial \theta_k \partial \omega_{lm}} \left(\tanh \left(h_i^{\mu} \right) - m_i^{\mu} \right) \right]_{q^*} \\ &+ \frac{1}{2} \int dx \left[p \left(x | \theta \omega \right) \left(x_j^{\mu} - m_j^{\mu} \right) \underbrace{\frac{\partial^2 \left(\tanh \left(h_i^{\mu} \right) - m_i^{\mu} \right)}{\partial \theta_k \partial \omega_{lm}}}_{\left[\frac{-2m_i^{\mu} q \left(1 - \left(m_i^{\mu q} \right)^2 \right) \delta_{ik} \delta_{il} x_m^{\mu}}{-\frac{\partial^2 x_j^{\mu}}{\partial \theta_k \partial \omega_{lm}}} \right]_{q^*} \\ &+ \frac{1}{2} \int dx \left[\frac{\partial p \left(x | \theta \omega \right) \left(x_j^{\mu} - m_j^{\mu} \right)}{\partial \omega_{lm}} \underbrace{\frac{\partial \left(\tanh \left(h_i^{\mu} \right) - m_i^{\mu} \right)}{\partial \theta_k}}_{\left[\frac{-2m_i^{\mu} q \left(1 - \left(m_i^{\mu q} \right)^2 \right) \delta_{il} \left(\frac{\partial \mu \left(x | \theta \omega \right) \left(x_j^{\mu} - m_j^{\mu} \right)}{\partial \theta_k} \right)}{\left(1 - \left(m_j^{\mu q} \right)^2 \right) \delta_{jl} \delta_{jk} \delta_{im}} + \frac{1}{2} \left(1 - \left(m_j^{\mu q} \right)^2 \right) \delta_{jl} \left(\frac{\partial \chi_{im}^{\mu}}{\partial \theta_k} \right)}_{q^*} \\ &+ \left(i \leftrightarrow j \right) \\ &= -\frac{1}{2} m_j^{\mu q} \left(1 - \left(m_j^{\mu q} \right)^2 \right) \delta_{jl} \delta_{jk} \delta_{im}} + \frac{1}{2} \left(1 - \left(m_j^{\mu q} \right)^2 \right) \delta_{jl} \left(\frac{\partial \chi_{im}^{\mu}}{\partial \theta_k} \right)}{\left(1.19 \right)} \end{split}$$

$$\begin{split} \frac{\partial^{2}\chi_{ij}^{\mu}}{\partial \omega_{kl}\partial \omega_{mn}}\Big|_{q^{*}} &= \frac{1}{2}\int dx \left[\frac{\partial^{2}p\left(x|\theta\omega\right)\left(x_{j}^{\mu}-m_{j}^{\mu}\right)}{\partial \omega_{kl}\partial \omega_{mn}} \left(\tanh\left(h_{i}^{\mu}\right)-m_{i}^{\mu}\right) \right]_{q^{*}} \\ &+ \frac{1}{2}\int dx \left[p\left(x|\theta\omega\right)\left(x_{j}^{\mu}-m_{j}^{\mu}\right) \underbrace{\frac{\partial^{2}\left(\tanh\left(h_{i}^{\mu}\right)-m_{i}^{\mu}\right)}{\partial \omega_{kl}\partial \omega_{mn}}}_{\left(\overset{(1.0)(1.9)}{=}-\frac{\partial^{2}m_{i}^{\mu}\left(1-\left(m_{i}^{\mu q}\right)^{2}\right)\delta_{ik}\delta_{im}\left(x_{i}^{\mu q}x_{i}^{kq}\right)}{-\frac{\partial^{2}m_{kl}^{\mu}\partial \omega_{mn}}q_{*}} \right]_{q^{*}} \\ &+ \frac{1}{2}\int dx \left[\frac{\partial p\left(x|\theta\omega\right)\left(x_{j}^{\mu}-m_{j}^{\mu}\right)}{\partial \omega_{mn}}\frac{\partial\left(\tanh\left(h_{i}^{\mu}\right)-m_{i}^{\mu}\right)}{\partial \omega_{kl}}\right]_{q^{*}} \\ &+ \frac{1}{2}\int dx \left[\frac{\partial p\left(x|\theta\omega\right)\left(x_{j}^{\mu}-m_{j}^{\mu}\right)}{\partial \omega_{kl}}\frac{\partial\left(\tanh\left(h_{i}^{\mu}\right)-m_{i}^{\mu}\right)}{\partial \omega_{kl}}\right]_{q^{*}} \\ &+ \frac{1}{2}\int dx \left[\frac{\partial p\left(x|\theta\omega\right)\left(x_{j}^{\mu}-m_{j}^{\mu}\right)}{\partial \omega_{kl}}\frac{\partial\left(\tanh\left(h_{i}^{\mu}\right)-m_{i}^{\mu}\right)}{\partial \omega_{kl}}\right]_{q^{*}} \\ &+ \frac{1}{2}\left(1-\left(m_{j}^{\mu q}\right)^{2}\right)\delta_{jm}\frac{\partial\chi_{in}^{\mu}}{\partial\omega_{kl}}\right]_{q^{*}} + \frac{1}{2}\left(1-\left(m_{j}^{\mu q}\right)^{2}\right)\delta_{jk}\frac{\partial\chi_{il}^{\mu}}{\partial\omega_{mn}}\right]_{q^{*}} \\ &- m_{j}^{\mu q}\left(1-\left(m_{j}^{\mu q}\right)^{2}\right)\delta_{jl}\delta_{jm}\left(\left(x_{i}^{\mu}x_{m}^{\mu}x_{i}^{\mu}\right)_{q^{*}} - m_{i}^{\mu}\left(x_{m}^{\mu}x_{i}^{\mu}\right)_{q^{*}}\right) \\ &+ \left(\left(1-\left(m_{j}^{\mu q}\right)^{2}\right)\delta_{jm}\delta_{ln} + \left(1-\left(m_{i}^{\mu q}\right)^{2}\right)\delta_{il}\delta_{ln}\right) \\ &- m_{j}^{\mu q}\left(1-\left(m_{j}^{\mu q}\right)^{2}\right)\delta_{jk}\delta_{jm} \\ &\left(\frac{m_{i}^{\mu q}}{2}\delta_{in}\delta_{ln} + \frac{m_{i}^{\mu q}}{3}\delta_{il} + \frac{m_{i}^{\mu q}}{3}\delta_{in} - \frac{m_{i}^{\mu q}}{6}\delta_{nl}\right) \\ &+ \left(i\leftrightarrow j\right), \end{split}$$

by partial differentiation and integration over the space of gene expression levels. There is no contribution quadratic in θ ,

$${}^{\theta\theta}\chi_{i}^{\mu q^{*}} = 0. \qquad (I.21)$$

Summing over all contributions with a contribution in $\delta \omega$, one gets the second order corrections in the covariance of gene expression,

$${}^{\theta\omega}\chi_{i}^{\mu q^{*}} = -\frac{1}{2}\left(1 - \left(m_{j}^{\mu q}\right)^{2}\right)\delta\omega_{ji}\delta\theta_{j} + (i \leftrightarrow j)$$
(I.22)

Finally, analogous to the first-order approximation, we calculate the covariance of gene expression in the second-order MFT,

$$\begin{split} \chi_{ij}^{\mu \,\text{MFT2}} &= +\frac{1}{2} \delta_{ij} + \frac{1}{4} \left(1 - (m_i^{\mu q})^2 \right) \omega_{ij} \\ &- \frac{1}{2} \left(1 - \left(m_j^{\mu q} \right)^2 \right) \omega_{ji} \theta_j^{\text{MFT1}} \\ &+ \frac{1}{8} \left(1 - (m_i^{\mu q})^2 \right) \left(1 - \left(m_j^{\mu q} \right)^2 \right) \sum_l \omega_{il} \omega_{jl} \\ &+ \frac{1}{8} \left(1 - \left(m_j^{\mu q} \right)^2 \right) \sum_k \omega_{jk} \left(1 - (m_k^{\mu q})^2 \right) \omega_{ki} \\ &- \frac{m_j^{\mu q}}{6} \left(1 - \left(m_j^{\mu q} \right)^2 \right) \left(m_i^{\mu q} \left(\omega_{ij} \right)^2 - 2\omega_{ij} \delta \theta_i^{\text{MFT1}} - \frac{m_j^{\mu q}}{2} \sum_k (\omega_{ik})^2 \right) \\ &+ (i \leftrightarrow j) \;, \end{split}$$
(I.24)

by plugging in the first order corrections, (3.56) and (3.57), and second order corrections, (I.22) and (I.23), into the Taylor expansion of χ^{μ}_{ij} , (3.51).

I.3. Network inference and response prediction

Within this Appendix section, a scatter plot of generated and reconstructed interaction is shown, and we compare predicted gene expression level to the measured ones for three additional pairs of training and prediction set. Figure I.1.: Scatter plot of generated and reconstructed interactions for a fully connected system as small as 10 nodes. The model parameter are set to $a_i = 1$, $b_i = 1$, and $c_i = 0.1$. For the inference 10 distinct single drug perturbations with 50 samples per perturbation are used.





Figure I.2.: Scatter plot of predicted and measured gene expression level.



Figure I.3.: Scatter plot of predicted and measured gene expression level.

- (a) Least squares 1st order
- (b) Least squares 2nd order



Figure I.4.: Scatter plot of predicted and measured gene expression level.

II. Cancer immunotherapy response prediction

Within this Appendix chapter we describe some technical details about the prediction of frameshift-derived peptide sequences and show supplementary figures in the context of CBI response prediction.

II.1. Frameshift-derived peptide sequences

Based on information about frameshift mutations, we predict the frameshiftderived peptide sequence using a reference genome and transcriptional information. As a human reference genome we use the GRCh37.p13 [110] and we obtain transcriptional information from the BioMart Ensemble database [111, 112].

For the vast majority of genes, different transcript variants exist. We used the most biologically relevant transcript within our data analysis, which is called the principal isoform in the Ensemble database.

The prediction of the frameshift-derived peptide sequence is based on the information about the frameshift mutation and the nucleotide sequence within the reference genome. To predict NMD and NSD, we check whether a frameshift mutation occurs in the last exon based on transcriptional information. If the mutation does not occur in the last exon, we also calculate the distance of the frameshift derived-stop codon to the next exon-exon junction.

We consider the following exon according to the principal isoform for frameshift-derived peptides that do not have a stop codon within the first exon.

II.2. Response classification and survival analysis

We show supplementary figures concerning the response classification and survival analysis in this Appendix section.

Figure II.1.: ROC curves based on the melanoma dataset published by Liu in the year 2019 [103].

(a) Number of mutation, missense mutations ●, nonsense mutations ●, silent mutations ●.





Figure II.2.: Kaplan-Meier curves based on hypothetical survival determinants for melanoma patients within the Liu dataset [103].



- Figure II.3.: Kaplan-Meier curves for patient groups characterised by high overall and high frameshift ●, high overall and low frameshift ●, low overall and high frameshift ●, low overall and low frameshift number of mutation.
 - (a) Lung Miao

