

---

# THE TENSOR NETWORK APPROACH TO EFFICIENT AND EXPLAINABLE AI

---

RESEARCH NOTES IN THE ENEXA PROJECT

Alex Goessmann, DATEV eG

June 24, 2025

## ABSTRACT

### 1 Abstract

Recent models in artificial intelligence, despite performance breakthroughs in large language models, suffer from limited efficiency and explainability, which prevents them from unlocking their full application potential for economic and trustworthy use. We in this work leverage the mathematical structure of tensor networks, which has been eminent in artificial intelligence from the beginning, to achieve the goals of efficiency and explainability.

While tensors appear naturally in artificial intelligence as factored representations of systems, their decompositions into tensor networks is necessary to avoid the curse of dimensionality. Since the curse of dimensionality prevents feasible generic representations, logical and probabilistic reasoning approaches trade off efficiency and generality. This work presents these tradeoffs in the tensor network formalism and formulates feasible reasoning algorithms involving tensor network contractions. We review the classical logical and probabilistic approaches to reasoning in the first part and develop applications in neuro-symbolic AI in the second part. In the third part we investigate in more detail schemes to exploit tensor network contractions for calculus.

## Contents

<b>1</b>	<b>Abstract</b>	<b>1</b>
<b>2</b>	<b>Introduction</b>	<b>10</b>
2.1	Background . . . . .	10
2.1.1	Logic and Explainability in AI . . . . .	10
2.1.2	Tensor Networks in AI . . . . .	11
2.1.3	Representation Schemes of Systems . . . . .	12
2.2	Structure of the work . . . . .	12
2.2.1	Part I: Classical Approaches . . . . .	12
2.2.2	Part II: Neuro-Symbolic Approaches . . . . .	12
2.2.3	Part III: Contraction Calculus . . . . .	13
2.2.4	Focus I: Representation . . . . .	14
2.2.5	Foucs II: Reasoning . . . . .	14
<b>3</b>	<b>Notation and Basic Concepts</b>	<b>14</b>

3.1	Categorical Variables and Representations . . . . .	14
3.2	Tensors . . . . .	15
3.3	Contractions . . . . .	16
3.3.1	Graphical Illustrations . . . . .	16
3.3.2	Tensor Product . . . . .	17
3.3.3	Generic Contractions . . . . .	18
3.3.4	Decompositions . . . . .	19
3.3.5	Directed Tensors and normalizations . . . . .	20
3.4	Function encoding schemes . . . . .	20
3.4.1	Basis encodings . . . . .	20
3.4.2	Tensor-valued functions . . . . .	21
<b>I</b>	<b>Classical Approaches</b>	<b>21</b>
<b>4</b>	<b>Introduction into Part I</b>	<b>21</b>
4.1	Representation of Factored Systems . . . . .	21
4.2	Mechanisms of tensor network decompositions . . . . .	22
<b>5</b>	<b>Probability Distributions</b>	<b>22</b>
5.1	Classical Properties of Distributions . . . . .	22
5.1.1	Probability Tensors . . . . .	22
5.1.2	Base measures . . . . .	23
5.1.3	Marginal Distribution . . . . .	24
5.1.4	Conditional Probabilities . . . . .	25
5.1.5	Bayes Theorem and the Chain Rule . . . . .	26
5.1.6	Independence . . . . .	27
5.2	Sufficient Statistics and Exponential Families . . . . .	30
5.2.1	Sufficient Statistics . . . . .	30
5.2.2	Exponential families . . . . .	31
5.2.3	Tensor Network Representation . . . . .	32
5.3	Graphical Models . . . . .	35
5.3.1	Markov Networks . . . . .	35
5.3.2	Bayesian Networks . . . . .	36
5.3.3	Bayesian Networks as Markov Networks . . . . .	38
5.3.4	Hidden Markov Models . . . . .	38
5.3.5	Markov Networks as Exponential Families . . . . .	39
5.3.6	Representation of generic distributions . . . . .	41
5.4	Polytopes of mean parameters . . . . .	41
5.4.1	Representation by convex hulls . . . . .	41
5.4.2	Representation as intersecting half-spaces . . . . .	42

5.4.3	Characterization of the interior . . . . .	43
5.4.4	Characterization of the boundary by faces . . . . .	44
5.5	Discussion and Outlook . . . . .	46
<b>6</b>	<b>Probabilistic Inference</b>	<b>46</b>
6.1	Queries . . . . .	46
6.1.1	Querying by functions . . . . .	46
6.1.2	Mode Queries . . . . .	47
6.1.3	Energy representations . . . . .	47
6.2	Sampling . . . . .	48
6.2.1	Exact Methods . . . . .	48
6.2.2	Gibbs Sampling . . . . .	49
6.2.3	Simulated Annealing . . . . .	50
6.3	Maximum Likelihood Estimation . . . . .	50
6.3.1	Empirical Distributions . . . . .	50
6.3.2	Likelihood Loss . . . . .	52
6.3.3	Entropic Interpretation . . . . .	53
6.4	Variational Inference in Exponential Families . . . . .	55
6.4.1	Forward and Backward Mappings . . . . .	55
6.4.2	Variational Formulation . . . . .	56
6.5	Maximum entropy distributions . . . . .	58
6.5.1	Entropy maximization problem . . . . .	58
6.5.2	Tensor Network Representation . . . . .	59
6.5.3	Modes of maximum entropy distributions . . . . .	59
6.5.4	Base measure refinement . . . . .	63
6.5.5	Mean parameters by soft and hard constraints . . . . .	65
6.6	Forward Mapping in Exponential Families . . . . .	66
6.6.1	Mode queries by annealing . . . . .	66
6.7	Mean Field Methods . . . . .	66
6.7.1	Naive Mean Field Method . . . . .	67
6.7.2	Structured Variational Approximation . . . . .	68
6.8	Backward Map in Exponential Families . . . . .	70
6.8.1	Variational Formulation . . . . .	70
6.8.2	Interpretation as a moment projection . . . . .	70
6.8.3	Approximation by alternating algorithms . . . . .	71
6.8.4	Second order Methods . . . . .	72
6.9	Discussion . . . . .	72
<b>7</b>	<b>Propositional Logic</b>	<b>73</b>
7.1	Encoding of Booleans . . . . .	73

7.1.1	Representation by coordinates . . . . .	73
7.1.2	Representation by basis vectors . . . . .	74
7.1.3	Coordinate and Basis Calculus . . . . .	74
7.2	Semantics of Propositional Formulas . . . . .	75
7.2.1	Formulas . . . . .	75
7.2.2	Basis encoding of formulas . . . . .	75
7.3	Syntax of Propositional Formulas . . . . .	76
7.3.1	Atomic Formulas . . . . .	77
7.3.2	Syntactical combination of formulas . . . . .	77
7.3.3	Syntactical decomposition of formulas . . . . .	79
7.4	Outlook . . . . .	81
<b>8</b>	<b>Logical Inference</b>	<b>81</b>
8.1	Entailment in Propositional Logics . . . . .	81
8.1.1	Deciding Entailment by contractions . . . . .	81
8.1.2	Deciding Entailment by partial ordering . . . . .	82
8.1.3	Redundancy of entailed formulas . . . . .	82
8.1.4	Contraction Knowledge Base . . . . .	83
8.2	Formulas as Random Variables . . . . .	84
8.2.1	Probabilistic queries by formulas . . . . .	84
8.2.2	Uniform distributions of the models . . . . .	84
8.2.3	Probability of a formula given a Knowledge Base . . . . .	85
8.2.4	Knowledge Bases as Base Measures for Probability Distributions . . . . .	86
8.3	Constraint Satisfaction Problems . . . . .	87
8.3.1	Deciding entailment on Markov Networks . . . . .	87
8.3.2	Categorical Constraints . . . . .	88
8.4	Deciding Entailment by local contractions . . . . .	91
8.4.1	Monotonicity of entailment . . . . .	91
8.4.2	Knowledge Cores . . . . .	91
8.4.3	Knowledge Propagation . . . . .	92
8.4.4	Applications . . . . .	94
8.4.5	Mimiking Inference Rules by Propagation . . . . .	95
8.5	Discussion . . . . .	95
<b>II</b>	<b>Neuro-Symbolic Approaches</b>	<b>95</b>
<b>9</b>	<b>Introduction into Part II</b>	<b>95</b>
<b>10</b>	<b>Formula Selecting Networks</b>	<b>96</b>
10.1	Construction schemes . . . . .	96

10.1.1	Connective Selecting Maps . . . . .	96
10.1.2	Variable Selecting Maps . . . . .	97
10.2	State Selecting Maps . . . . .	97
10.3	Composition of formula selecting maps . . . . .	98
10.3.1	Formula Selecting Neuron . . . . .	99
10.3.2	Formula Selecting Neural Network . . . . .	99
10.4	Application of Formula Selecting Networks . . . . .	101
10.4.1	Representation of selection encodings . . . . .	101
10.4.2	Efficient Representation of Formulas . . . . .	101
10.4.3	Batch contraction of parametrized formulas . . . . .	102
10.4.4	Average contraction of parametrized formulas . . . . .	102
10.5	Examples of formula selecting neural networks . . . . .	102
10.5.1	Correlation . . . . .	102
10.5.2	Conjunctive and Disjunctive Normal Forms . . . . .	103
10.6	Extension to variables of larger dimension . . . . .	103
<b>11</b>	<b>Logical Network Representation</b>	<b>104</b>
11.1	Markov Logic Networks . . . . .	104
11.1.1	Markov Logic Networks as Exponential Families . . . . .	104
11.1.2	Tensor Network Representation . . . . .	104
11.1.3	Expressivity . . . . .	107
11.1.4	Examples . . . . .	108
11.2	Hard Logic Networks . . . . .	109
11.2.1	The limit of hard logic . . . . .	110
11.2.2	Tensor Network Representation . . . . .	111
11.3	Hybrid Logic Network . . . . .	112
11.3.1	Tensor Network Representation . . . . .	113
11.3.2	Reasoning Properties . . . . .	113
11.3.3	Expressivity . . . . .	114
11.4	Polynomial Representation . . . . .	114
11.5	Applications . . . . .	116
<b>12</b>	<b>Logical Network Inference</b>	<b>116</b>
12.1	Mean parameters of Hybrid Logic Networks . . . . .	116
12.1.1	Vertices by hard logic networks . . . . .	117
12.1.2	Faces of larger rank . . . . .	117
12.1.3	Mean parameters in the interior . . . . .	117
12.1.4	Mean parameters outside the interior . . . . .	118
12.1.5	Expressivity of Hybrid Logic Networks . . . . .	119
12.1.6	Case of tree computation networks . . . . .	121

12.1.7	Examples . . . . .	121
12.2	Entropic Motivation of unconstrained Parameter Estimation . . . . .	121
12.2.1	Maximum Likelihood in Hybrid Logic Networks . . . . .	121
12.2.2	Maximum Entropy in Hybrid Logic Networks . . . . .	122
12.3	Alternating Algorithms to Approximate the Backward Map . . . . .	122
12.4	Forward and backward mappings in closed form . . . . .	124
12.4.1	Maxterms and Minterms . . . . .	124
12.4.2	Atomic formulas . . . . .	125
12.5	Constrained parameter estimation in the minterm family . . . . .	126
12.5.1	Parameter Estimation . . . . .	127
12.5.2	Structure Learning . . . . .	127
12.6	Greedy Structure Learning . . . . .	127
12.6.1	Greedy formula inclusions . . . . .	128
12.6.2	Gain Heuristic . . . . .	128
12.6.3	Gradient heuristic and the proposal distribution . . . . .	130
12.6.4	Iterations . . . . .	130
12.7	Proposal distribution . . . . .	131
12.7.1	Mean parameter polytope . . . . .	131
12.8	Discussion . . . . .	131
<b>13</b>	<b>Probabilistic Guarantees</b>	<b>132</b>
13.1	Fluctuations of random data . . . . .	132
13.1.1	Fluctuation of the empirical distribution . . . . .	132
13.1.2	Mean parameter of the empirical distribution . . . . .	133
13.1.3	Noise tensor and its width . . . . .	133
13.2	Error bounds based on the noise width . . . . .	134
13.2.1	Parameter Estimation . . . . .	134
13.2.2	Structure Learning . . . . .	136
13.2.3	Mode recovery . . . . .	136
13.3	Fluctuations in Logic Networks . . . . .	137
13.3.1	Energy tensor in proposal distributions . . . . .	138
13.3.2	Minterm Exponential Family . . . . .	138
13.3.3	Guarantees for Mode of the Proposal Distribution . . . . .	138
13.3.4	Guarantees for Unconstrained Parameter Estimation . . . . .	139
13.4	Width bounds for the noise tensor . . . . .	139
13.4.1	Basis Vectors . . . . .	139
13.4.2	Sphere . . . . .	141
13.5	Discussion . . . . .	142
<b>14</b>	<b>First Order Logic</b>	<b>142</b>

14.1	World Tensors . . . . .	142
14.1.1	Case of Propositional Logics . . . . .	143
14.1.2	One-hot encoding of worlds . . . . .	143
14.1.3	Probability distributions . . . . .	143
14.1.4	Semantics of formulas . . . . .	144
14.1.5	Two levels of tensor representation . . . . .	144
14.2	Formulas in a fixed first-order logic world . . . . .	144
14.2.1	Grounding tensors . . . . .	144
14.2.2	Atomic Formulas . . . . .	145
14.2.3	Formula synthesis by connectives . . . . .	145
14.2.4	Quantifiers . . . . .	146
14.2.5	Storage in basis CP decomposition . . . . .	147
14.2.6	Queries . . . . .	147
14.3	Representation of Knowledge Graphs . . . . .	148
14.3.1	Representation as unary and binary predicates . . . . .	148
14.3.2	Representation as ternary predicate . . . . .	148
14.3.3	SPARQL Queries . . . . .	149
14.4	Probabilistic Relational Models . . . . .	151
14.4.1	Hybrid First-Order Logic Networks . . . . .	151
14.4.2	Base measures by importance formulas . . . . .	151
14.4.3	Decomposition of the log likelihood . . . . .	152
14.4.4	Decomposition of the Partition function . . . . .	153
14.5	Sample extraction from first-order logic worlds . . . . .	155
14.5.1	Representation by Tensor Networks . . . . .	155
14.5.2	Basis CP Decomposition of extracted data . . . . .	156
14.6	Generation of first-order logic worlds . . . . .	158
14.6.1	Samples by single objects . . . . .	159
14.6.2	Samples by pairs of objects . . . . .	160
14.7	Discussion . . . . .	161
<b>III</b>	<b>Contraction Calculus</b>	<b>161</b>
<b>15</b>	<b>Introduction into Part III</b>	<b>161</b>
<b>16</b>	<b>Coordinate Calculus</b>	<b>161</b>
16.1	One-hot encodings as basis . . . . .	161
16.2	Coordinatewise Transforms . . . . .	163
16.3	Directed Tensors . . . . .	164
16.3.1	normalization . . . . .	165
16.3.2	normalization Equations . . . . .	165

16.3.3	Contraction of Directed Tensors	166
16.4	Proof of Hammersley-Clifford Theorem	166
16.5	Differentiation of Contraction	168
16.6	Discussion	170
<b>17</b>	<b>Basis Calculus</b>	<b>170</b>
17.1	Basis Encoding of Subsets	170
17.1.1	Binary Relations	171
17.1.2	Higher order relations	171
17.2	Basis Encoding of Functions	172
17.2.1	Basis encoding of Functions	172
17.2.2	Function Evaluation	173
17.3	Calculus of basis encodings	173
17.3.1	Composition of function	173
17.3.2	Compositions with real functions	174
17.3.3	Decomposition in case of structured images	175
17.4	Selection Encodings	176
17.4.1	Basis representations of linear maps	176
17.4.2	Selection encodings as basis representations	177
17.5	Indicator features to functions	179
17.5.1	Connections with computable families	179
17.5.2	Composition of functions	181
17.5.3	Effective Representation of Partition Statistics	182
17.6	Hybrid Basis and Coordinate Calculus	183
17.7	Applications in Machine Learning	184
<b>18</b>	<b>Sparse Calculus</b>	<b>185</b>
18.1	CP Decomposition	185
18.1.1	Directed Leg Cores	186
18.1.2	Basis CP decompositions and the $\ell_0$ -norm	187
18.1.3	Basis+ CP decompositions and polynomials	188
18.2	Constructive Bounds on CP Ranks	190
18.2.1	Cascade of ranks	190
18.2.2	Operations on CP decompositions	191
18.3	Sparse Encoding of Functions	193
18.4	Optimization of sparse tensors	195
18.4.1	Unconstrained Binary Optimization	195
18.4.2	Integer Linear Programming	196
<b>19</b>	<b>Tensor Approximation</b>	<b>198</b>
19.1	Selection tensor networks for CP decompositions	198



19.1.1 Applications . . . . .	200
19.2 Approximation of Energy tensors . . . . .	200
19.3 Transformation of Maximum Search to Risk Minimization . . . . .	201
19.3.1 Weighted Squares Loss Trick . . . . .	201
19.3.2 Problem of the trivial tensor . . . . .	201
19.4 Alternating Solution of Least Squares Problems . . . . .	202
19.4.1 Choice of Representation Format . . . . .	202
19.5 Regularization and Compressed Sensing . . . . .	202
19.6 Discussion and Outlook . . . . .	203
<b>20 Message Passing</b>	<b>203</b>
20.1 Commutation of Contractions . . . . .	203
20.2 Exact Contractions . . . . .	204
20.2.1 Construction of Cluster Graphs . . . . .	204
20.2.2 Message Passing to calculate contractions . . . . .	205
20.2.3 Variable Elimination Cluster Graphs . . . . .	206
20.2.4 Bethe Cluster Graphs . . . . .	206
20.2.5 Computational Complexity . . . . .	207
20.3 Boolean Message Passing . . . . .	207
20.3.1 Monotonocity of tensor contraction . . . . .	207
20.3.2 Invariance of adding subcontractions . . . . .	208
20.3.3 Basis Calculus as message passing scheme . . . . .	208
20.3.4 Application . . . . .	208
20.4 Discussion . . . . .	209
<b>A Implementation in the tnreason package</b>	<b>209</b>
A.1 Architecture . . . . .	209
A.2 Implementation of basic notation . . . . .	210
A.2.1 Categorical Variables and Representations . . . . .	210
A.2.2 Tensors . . . . .	210
A.2.3 Contractions . . . . .	211
A.2.4 Function encoding schemes . . . . .	211
A.3 Subpackage engine . . . . .	211
A.3.1 Basis+ CP Decompositions storing values . . . . .	211
A.3.2 Contractions . . . . .	212
A.4 Subpackage representation . . . . .	213
A.4.1 Computation Activation Networks . . . . .	214
A.5 Subpackage reasoning . . . . .	215
A.5.1 Sampling . . . . .	215
A.5.2 Variational Inference . . . . .	215

A.5.3	Optimization . . . . .	215
A.6	Subpackage application . . . . .	216
A.6.1	Representation of formulas . . . . .	216
A.6.2	Script Language . . . . .	216
A.6.3	Distributions . . . . .	218
A.6.4	Inference . . . . .	219
A.6.5	Learning . . . . .	220
<b>B</b>	<b>Glossary</b>	<b>220</b>
B.1	Tensors . . . . .	220
B.2	Variables . . . . .	221
B.3	Maps . . . . .	222
B.4	Contraction equations . . . . .	222
<b>C</b>	<b>*</b>	<b>222</b>

## 2 Introduction

Artificial intelligence is a long-standing dream, which has in recent years received enormous attention, especially driven by breakthroughs in large language models. Among the key priorities towards an economic and trustworthy usage is the improvement of efficiency and explainability of models.

Instead of post-hoc explainability of a models given inference on specific data, this work aims at the intrinsic human understandability of a model. We are motivated by the theory of logic, whose formalization of human thoughts serves as an interface between mechanized reasoning on a machine and human understandability. This advanced form of explainability enables novel forms of human interactions with a model based on verbalizations, manipulations and guarantees on the models inference output.

The need for efficiency stems more from economic concerns on the resource demand of training and inferring a model. Tensors naturally represent states of systems with multiple variables, both in logical and probabilistic approaches towards artificial intelligence. However, even for a moderate numbers of variables, the curse of dimensionality prevents a typical machine's memory to store a generic representation. Carefully designing representation formats is therefore crucial to prevent exponential storage growth while balancing expressivity and efficiency.

In this work, we utilize the formalism of tensor networks in the creation of efficient representation schemes. The chosen tensor network formats are motivated as explainable model architectures and provide a synergy between the aims of efficiency and explainability. More precisely, tensor networks appear as the numerical structures behind probabilistic graphical models and logical knowledge bases. Understanding these foundations of tensor networks reveals their vast application potential in neuro-symbolic artificial intelligence.

### 2.1 Background

Before presenting an overview over the contents, we further motivate this work based on the broach approaches towards artificial intelligence and more recent developments.

#### 2.1.1 Logic and Explainability in AI

The logical tradition of artificial intelligence is motivated by the resembling of human thought in logics McCarthy (1959). Historic approaches towards artificial intelligence have focused on models by vast knowledge bases and inference by logical reasoning. The main problems hindering the success of this approach is the inability of classical first-order logic to handle uncertainty of information, as present in realistic scenarios.

Integrating observed data into a learning process has been framed Inductive Logic Programming Muggleton and De Raedt (1994). Along that line the Amie method Galárraga et al. (2013) has been developed to learn Horn clauses using a refinement operator. Class Expression Learning Lehmann et al. (2011) is a more recent approach

to assist in the design of reasoning capabilities in Knowledge Graphs. However, this approach is limited by the expressivity of description logics and the exponentially large hypothesis sets for the choice of formulas. Efficient search methods in these exponentially large hypothesis sets have been provided based on reinforcement learning Demir and Ngonga Ngomo (2021) and neural networks Kouagou et al. (2022, 2023).

Logical approaches are still dominant in the description of data. Here the semantic web initiative developed data storage formats based on Knowledge Graphs Antoniou et al. (2012); Hogan et al. (2021), which describe structured data based on description logic.

Towards extending the practical usage of logics, the field of Statistical Relational AI Nickel et al. (2016); Getoor and Taskar (2019) studies statistical models of logical relations. This directly treats uncertainty and therefore unifies logics with statistical approaches. This aims have more recently reframed as neuro-symbolic AI Hochreiter (2022); Sarker et al. (2022), with close relations to statistical relational AI Marra et al. (2024). Neuro-symbolic AI focuses on the unification of the neural and the symbolic paradigm Garcez et al. (2019), where early approaches are Towell and Shavlik (1994); Avila Garcez and Zaverucha (1999). While the symbolic paradigm is roughly understood as human understandable reasoning in formal logics, the neural paradigm is the computational benefit of decomposing a model into layerwise computation. These decompositions provide both expressive and efficient to train and infer model architectures. While modern black-box AI focuses on large neural networks, which size prevents human understanding of the inference process, neuro-symbolic AI aims at a re-implementation of the symbolic paradigm into such architectures.

### 2.1.2 Tensor Networks in AI

Decomposition schemes of tensors have been developed in numerics to efficiently operate in high-dimensional tensor spaces Hackbusch (2012) and to avoid the curse of dimensionality Bellman (1961). Each decomposition schemes has a graphical depiction, as we will introduce in Chapter 3, and decompositions are therefore referred to as networks. The first decomposition schemes by Tucker, originally introduced in Hitchcock (1927), suffered from exponential increases of the degrees of freedom with the tensor order. The CP format (see Chapter 18) can in principle establish storage in linear with the order. Sets of tensors with fixed rank with respect to this format are however not closed Beylkin and Mohlenkamp (2005) and approximation problems are often ill posed de Silva and Lim (2008). The Tensor Train decomposition Oseledets and Tyrtshnikov (2009), which appears in quantum mechanics as matrix product states Perez-Garcia et al. (2007), overcomes these numerical problems Holtz et al. (2012). Hierarchical Tucker decompositions Hackbusch and Kühn (2009) are generalizations of tensor train decompositions, which have useful properties for tensor approximations Grasedyck (2010); Falco and Hackbusch (2012).

Tensors are used in the processing of big data Cichocki (2014) and in many-body physics Orús (2019). Besides that, there have been pioneering approaches to exploit them in the data-driven identification of governing equations Gelß et al. (2019); Goeßmann et al. (2020), more general supervised learning Stoudenmire and Schwab (2016) and the simulation of noisy quantum mechanics Sander et al. (2025). The duality between tensor networks and graphical models has been first discussed in Robeva and Seigal (2019) and motivated further expressivity studies such as Glasser et al. (2019). Tensor Networks have further been applied for batch logical inference Sakama et al. (2017); Sato (2017); Tsilionis et al. (2024). Whereas these are conceptually interesting approaches, they have so far been limited to matrix multiplication, whereas obvious expressivity benefits would come from more general contraction schemes. Similar ideas have been led to TensorLog Cohen et al. (2020), Real Logic Serafini and d'Avila Garcez (2016) and based on that Logical Tensor Networks Badreddine et al. (2022).

Further, sparse representation of knowledge graphs by tensor networks has motivated several embedding schemes for objects in the knowledge graph. The sparse decomposition of the adjacency tensor capturing the ternary relations between objects provides embeddings schemes encoding relations between the objects in a latent space. The specific approaches distinguish between the format used, where Nickel et al. (2011) and Balazevic et al. (2019) used Tucker decompositions, Yang et al. (2015) the CP decomposition and Trouillon and Nickel (2017) complex extensions. Beyond embeddings, tensor based storage of knowledge graphs has recently shown tremendous improvements in querying knowledge graphs Bigerl et al. (2020). Here, queries on the knowledge graph are performed as contractions of the tensors efficiently representing knowledge graphs.

Tensors further serve as a central object in large-scale machine learning libraries such as TensorFlow Abadi et al. (2016) and PyTorch Paszke et al. (2019). Layerwise execution of neural network inference amounts then to tensor network contractions of tensors storing the activation of previous layers and weights. Beyond providing a central framework for the software design, also the design of AI-dedicated hardware orients on tensor contractions, with a current focus on Tensor Processing Units (TPU) Nikolić et al. (2022); Jouppi et al. (2023). Both the dedicated software and hardware design exploits the parallelization potential rooted in the contraction formalism of tensor networks. Besides these developments there exist several experimental libraries dedicated to the tensor-train tensor format ????

### 2.1.3 Representation Schemes of Systems

We start with ontological commitments in the description of a system and follow the book Russell and Norvig (2021) distinguishing atomic, factored and structured representations. While in atomic representation, the states of a systems are enumerated and represented in a single variable, factored representations describe a systems state based on a collection of variables. In the tensor formalism, each state of a system corresponds with a coordinate of a representing tensor. The order of the tensor coincides therefore with the number of variables in a system. In an atomic representation, where there is a single coordinate, each state corresponds with a coordinate of the representing vector being a tensor of order one. Having a factored representation with two variables requires order two tensors or matrices, where a coordinate is specified by a row and a column index. Given larger numbers of coordinates now extends this representation picture to tensors of larger orders, which have more abstract axes besides rows and columns. The generalization of the atomic representation to a factored system thus corresponds with the generalization of vectors towards matrices and tensors of larger orders. Along this line, we can always transform a factored representation of a system to an atomic one, just by enumerating the states of the factored system and interpreting them by a single variable. This amounts to the flattening of a representing tensor to a vector. However, by doing so, we would loose much of the structure of the representation, which we would like to exploit in reasoning processes.

A more generic representation of systems are structured representation. Structured representations involve objects of differing numbers and relations between them. As a consequence the numbers of variables can differ depending on the state of a system. This poses a challenge to the tensor representation, since a fixed number of variables is required to motivate a tensor space of representations. There are approaches to circumvent these difficulty by the development of template models such as Markov Logic Networks Richardson and Domingos (2006), which are instantiated on systems with differing number of objects. We will discuss those in Chapter 14.

In this work we treat discrete systems, where the number of states is finite. One can understand them as a discretization of continuous variables and many results will generalize by the converse limit to the situation of continuous variables.

Besides ontological commitments in the choice of a representation scheme, modelling a system also requires epistemologic commitments, by defining what properties are to be reasoned about. In logical approaches the properties of states are boolean values representing whether a state is consistent with known constraints. Probabilistic approaches assign to the coordinates of the tensors numbers in  $[0, 1]$  encoding the probability of a state. Compared with logical approached to reasoning, probabilistic approaches thus bear a more expressive modelling.

## 2.2 Structure of the work

The chapters are structured into three parts, and two focuses, see Figure 1.

### 2.2.1 Part I: Classical Approaches

The probabilistic and logical approaches towards artificial intelligence are reviewed in the tensor network formalism. We in this part restrict the discussion to atomic and factored system representation. In probability theory (see Chapter 5 and Chapter 6), tensors appear as generalized truth tables, storing the joint probability of each possible state of a system in factored representation. Tensors describing such distributions are of non-negative coordinates and are normed, which we will formalize by directed edges of hypergraphs. Applying the formalism, we introduce marginalization and conditioning operations based on contractions, and show how assumptions such as conditional independence lead to network decompositions. We then study the formalism of exponential families of probability distributions, which generalizes probabilistic graphical models. For generic exponential families we provide in Chapter 5 a tensor network representation, which structure is exploited for inference in Chapter 6. In logics (see Chapter 7), we motivate boolean tensors as a natural representation of propositional semantics. Logical entailment is then in Chapter 8 decided based contractions of these tensors, which we will further relate with marginal distributions in probabilistic inference. The syntax of propositional logics thereby hints at efficient decompositions schemes of these semantic representing tensors. We exploit the syntax to find efficient tensor network decompositions of the tensors in Chapter 7 and use them for efficient logical inference algorithms in and Chapter 8.

### 2.2.2 Part II: Neuro-Symbolic Approaches

Motivated by the classical approaches we apply the tensor network formalism towards learning and inferring neuro-symbolic models. We understand the decomposition of tensors into networks as an implementation of the neural paradigm of AI. Further, the symbolic paradigm is eminent in the interpretation of tensor networks using logical syntax, and enables the human-interbretable verbalization of learned models. Motivated by this central thoughts, we present vast classes of interpretable models in Chapter 11, which are unifying the logical and probabilistic approaches

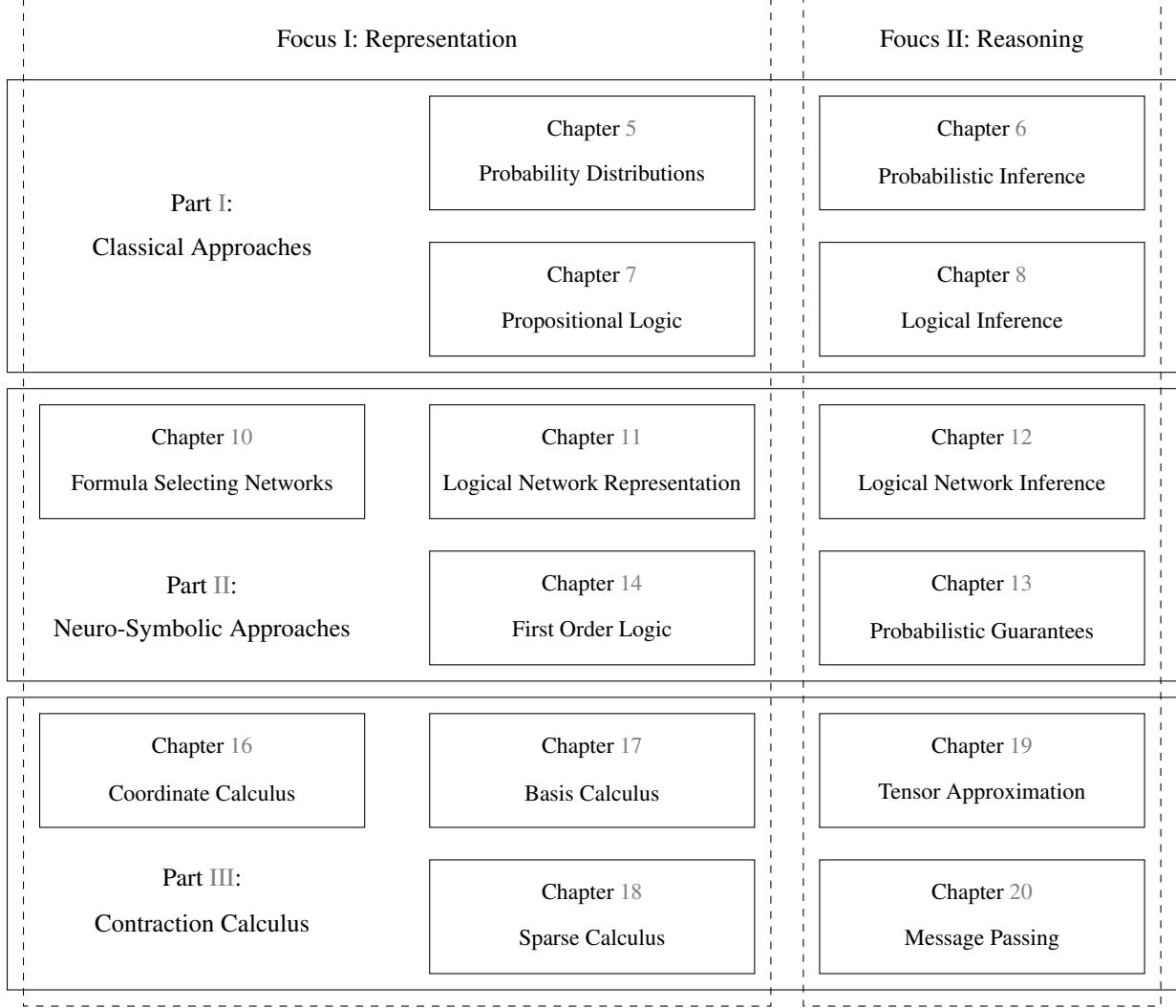


Figure 1: Sketch of the structure of this work. We assign the chapters to three parts and two focuses. The parts distinguish the coarse topics of this work into classical, neuro-symbolic approaches and the applied contraction calculus. The assigned focuses indicate whether the chapter orients more onto a representation format of the respective concepts or onto its exploitation in reasoning.

studied in Part I. The central idea here is to leverage the formalism of exponential families by choosing base measures and statistics based on logical formulas. We then turn in Chapter 12 towards inductive learning scenarios in this formalism, where new features are to be learned from data and parameters are calibrated. Here we apply the parametrization schemes developed in Chapter 10 to represent hypothesis classed for new features. While these approaches rely on propositional logics, in Chapter 14 we extend towards more expressive first-order logics. With knowledge graphs serving as examples we therein provide a tensor-network formalism to capture queries and motivate our learning schemes in propositional logics based on queries on random first-order worlds. In Chapter 13 we further derive statistical guarantees on these learning methods given random data, based on probabilistic bounds on uniform concentration events.

### 2.2.3 Part III: Contraction Calculus

In Part III the applied schemes of calculus using tensor network contractions are investigated in more detail. In particular, we distinguish between the schemes of coordinate, basis and sparse calculus. Coordinate calculus will be discussed in Chapter 16 using one-hot encodings as orthogonal basis elements. We will further properties related to directed tensors and a generic version of the Hammersley-Clifford decomposition theorem, which have been applied in

the probabilistic approached in Part I. Basis calculus in Chapter 17 introduces generic encodings of subsets, relations and functions by boolean tensors used in previous parts. We show, that these encoding schemes translate function compositions into tensor network contractions and are therefore a central technique to execute batchwise function evaluation by efficient tensor network contractions. In Chapter 18 we provide sparse schemes oriented on the CP format for the storage of tensors. We further investigate the origins of sparsity based on encodings of functions, and provide rank bounds for summations and contractions of these tensors. Then we formalize optimization problems as maximal coordinate searched among tensors and relate the investigated CP formats with standard optimization frameworks. We continue with studies of tensor approximation in Chapter 19, where we adapt formula selecting networks of Chapter 10 to select sparse CP tensors. In Chapter 20 we then investigate schemes of efficient contraction calculus based on local contractions, which are passed through the network as messages. These schemes can be regarded as generic numerical tools underlying message passing schemes such as belief propagation in probability theory and constraint propagation in logics.

#### 2.2.4 Focus I: Representation

In this focus, we motivate and investigate the efficient representation of tensors based on tensor network decompositions, where formats are captured by hypergraph as we introduce in Chapter 3. Besides being a necessity to overcome the curse of dimensionality, we show in Part I multiple motivations of tensor network decompositions originating from principles of artificial intelligence. As such, decompositions originate from conditional independence assumptions on probability distributions (see Chapter 5) and from logical syntax (see Chapter 7). Towards neuro-symbolic AI, we provide in Chapter 10 a generic representation scheme for batches of logical formulas. This scheme introduces additional axes to a tensor, which are assigned with selection variables and which slices select specific tensors. We exploit this scheme in Chapter 11 for efficient representation of exponential families, which statistics are sets of logical formulas. In Part III we investigate the applied representation scheme from a more theoretical viewpoint. More precisely, we distinguish between the schemes of coordinate calculus (Chapter 16) and basis calculus (Chapter 17). These schemes differ in the exploitation of the real coordinates of a tensor or of sums over chosen basis elements, in the encoding of information. In Chapter 18 we define restricted CP decompositions of tensors for sparse representations of  $d$ -ary relations, which appear in sparse representation of relational databases.

#### 2.2.5 Focus II: Reasoning

We develop schemes to efficiently perform inductive and deductive reasoning based on information stored in decomposed tensor. Contractions of tensor networks representing models in artificial intelligence are the central scheme to retrieve information. While in probability theory contractions compute marginal distribution (see Chapter 6), contraction of logical formulas are model counts central to the formalism of logical entailment (see Chapter 8). We will further exploit them to calculate queries in first-order logic such as on knowledge graphs (see Chapter 14). The statistical foundation on the success of contraction-based learning, which lies in the phenomenon of uniform concentration of contractions with empirical random tensors, will be investigated in Chapter 13. In Part III we further study generic tools for efficient execution of contraction-based reasoning. The tensor network approximation schemes in Chapter 19 bear the potential to approximate reasoning tasks by more efficient ones. The efficient execution of contractions using message-passing algorithms in Chapter 20 have been exploited in a variety of exact and approximated reasoning schemes.

### 3 Notation and Basic Concepts

We here provide the fundamental definitions of tensors, which are essential for the content in Part I and Part II. In Part III we will further investigate the properties of tensors focusing on their contractions.

#### 3.1 Categorical Variables and Representations

We will in this work investigate systems, which are described by a set of properties, each called categorical variables. This is called an ontological commitment, since it defines what properties a system has.

**Definition 1.** *An atomic representation of a system is described by a categorical variables  $X$  taking values  $x$  in a finite set*

$$[m] := \{0, \dots, m-1\}$$

*of cardinality  $m$ .*

We will in this work always notate categorical variables by large literals and indices by small literals, possible with other letters such as  $X, L, O, J$  and corresponding values  $x, l, o, j$ .

**Definition 2.** A factored representation of a system is a set of categorical variables  $X_k$ , where  $k \in [d]$ , taking values in  $[m_k]$ .

### 3.2 Tensors

Tensors are multiway arrays and a generalization of vectors and matrices to higher orders. We will first provide a formal definition as real maps from index sets enumerating the coordinates of vectors, matrices and larger order tensors.

**Definition 3 (Tensor).** Let there be numbers  $m_k \in \mathbb{N}$  for  $k \in [d]$  and categorical variables  $X_k$  taking their values in  $[m_k]$ . A Tensors  $\tau [X_0, \dots, X_{d-1}]$  of order  $d$  and leg dimensions  $m_d$  is defined through its coordinates

$$\tau [X_0 = x_0, \dots, X_{d-1} = x_{d-1}] \in \mathbb{R}$$

for index tuples

$$x_0, \dots, x_{d-1} \in \bigtimes_{k \in [d]} [m_k].$$

Tensors  $\tau [X_0, \dots, X_{d-1}]$  are elements of the tensor space

$$\bigotimes_{k \in [d]} \mathbb{R}^{m_k}$$

which is a linear space, enriched with the operations of coordinatewise summation and scalar multiplication.

We here introduced tensors in a non-canonical way based on categorical variables assigned to its axis. While coming as syntactic sugar at this point, this will allow us to define contractions without further specification of axes, based on comparisons of shared categorical variables. Especially, this eases the implementation of tensor network contractions without the need to further specify a graph (see Appendix A).

We abbreviate lists  $X_0, \dots, X_{d-1}$  of categorical variables by  $X_{[d]}$ , that is denote  $\tau [X_0, \dots, X_{d-1}]$  by  $\tau [X_{[d]}]$ . Occasionally, when the categorical variables of a tensor are clear from the context, we will omit the notation of the variables.

**Example 1 (Trivial Tensor).** The trivial tensor

$$\mathbb{I} [X_{[d]}] \in \bigotimes_{k \in [d]} \mathbb{R}^{m_k}$$

is defined by all coordinates being 1, that is for all  $x_0, \dots, x_{d-1} \in \bigtimes_{k \in [d]} [m_k]$

$$\mathbb{I} [X_{[d]} = x_{[d]}] = 1.$$

We will often encounter situations, where the coordinates of tensors are in  $\{0, 1\} = [2]$ .

**Definition 4.** We call a tensor  $\tau [X_{[d]}]$  boolean, when  $\text{im}(\tau) \subset [2]$ , i.e. all coordinates are either 0 or 1.

We are now ready to provide the link between tensors and states of systems with factored representations. To this end, we define the one-hot encoding of a state, which is a bijection between the states and the basis elements of a tensor space.

**Definition 5 (One-hot encodings to Atomic Representations).** Given an atomic system described by the categorical variable  $X$ , we define for each  $x \in [m]$  the basis vector  $\epsilon_x [X]$  by the coordinates

$$\epsilon_x [X = \tilde{x}] = \begin{cases} 1 & \text{if } x = \tilde{x} \\ 0 & \text{else} \end{cases} \quad (1)$$

The one-hot encoding of states  $x \in [m]$  of the atomic system described by the categorical variable  $X$  is the map

$$\epsilon : [m] \rightarrow \mathbb{R}^m$$

which maps  $x \in [m]$  to the basis vectors  $\epsilon_x [X]$ .

The basis vectors  $\epsilon_x[X]$  are tensors of order 1 and leg dimension  $m$  of the structure

$$\epsilon_x[X] = [0 \quad \cdots \quad 0 \quad 1 \quad 0 \quad \cdots \quad 0]^T, \quad (2)$$

where the 1 is at the  $x$ th coordinate of the vector.

We have so far described one-hot representations of the states of a single categorical variable, which would suffice to encode the state of an atomic system. In a factored system on the other side, we are dealing with multiple categorical variables.

**Definition 6** (One-hot encodings to Factored Representations). *Let there be a factored system defined by a tuple  $(X_0, \dots, X_{d-1})$  of variables taking values in  $\times_{k \in [d]} [m_k]$ . The one-hot encoding of its states is the tensor product of the one-hot encoding to each categorical variables, that is the map*

$$\epsilon : \times_{k \in [d]} [m_k] \rightarrow \bigotimes_{k \in [d]} \mathbb{R}^{m_k}$$

defined by mapping  $x_0, \dots, x_{d-1} = x_{[d]}$  to

$$\epsilon_{x_{[d]}}[X_{[d]}] =: \bigotimes_{k \in [d]} \epsilon_{x_k}[X_k].$$

We will call one-hot representations tensor representations and depict them as

$$\begin{array}{c} \boxed{\bigotimes_{k \in [d]} \epsilon_{x_k}} \\ | \quad | \quad \cdots \quad | \\ X_0 \quad X_1 \quad \cdots \quad X_{d-1} \end{array} = \begin{array}{c} \boxed{\epsilon_{x_0}} \\ | \\ X_0 \end{array} \otimes \begin{array}{c} \boxed{\epsilon_{x_1}} \\ | \\ X_1 \end{array} \otimes \cdots \otimes \begin{array}{c} \boxed{\epsilon_{x_{d-1}}} \\ | \\ X_{d-1} \end{array}$$

In Chapter 16 we will investigate the image of  $\epsilon$  in more detail and show that it is an orthonormal basis of the tensor space  $\bigotimes_{k \in [d]} \mathbb{R}^{m_k}$ .

**Remark 1** (Flattening of Tensors). *The use the tensor product to represent states of factored systems can be motivated by the reduction to atomic systems by enumeration of the states. We have this property reflected in the state encoding of factored systems, since the tensor space  $\bigotimes_{k \in [d]} \mathbb{R}^{m_k}$  is isomorphic to the vector spaces  $\mathbb{R}^{\prod_{k \in [d]} m_k}$ . This operation is called flattening (or unfolding) of tensors with many axes to tensors of less axes.*

### 3.3 Contractions

Contractions are the central manipulation operation on sets of tensors. To introduce them, we will develop a graphical illustration of sets of tensors, which we also call tensor networks. In Part III we will further investigate the utility of contractions in representing specific calculations, which demand different encoding schemes.

#### 3.3.1 Graphical Illustrations

Sets of tensor with categorical variables assigned to each legs implicitly carry a notion of a hypergraph. This perspective is especially useful, when some categorical variables are assigned to axis of multiple tensors, as it will often be the case in the applications considered in this work. Each variable can then be labeled by a node and each tensor as a hyperedge containing the nodes to its axis variables. Let us first formally introduce hypergraphs, which are generalizations of graphs allowing edges to be arbitrary nonempty subsets of the nodes, whereas canonical graphs demand a cardinality of two.

**Definition 7.** *A hypergraph is a pair  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$  of a set of nodes  $\mathcal{V}$  and a set of edges  $\mathcal{E}$ , where each hyperedge  $e \in \mathcal{E}$  is a subset of the nodes  $\mathcal{V}$ . A directed hypergraph is a pair  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ , such that each hyperedge  $e \in \mathcal{E}$  is the tuple of two disjoint sets  $e^{\text{in}}, e^{\text{out}} \subset \mathcal{V}$ , that is*

$$e = (e^{\text{in}}, e^{\text{out}}).$$

We will use the standard visualization by factor graphs as a diagrammatic illustration of sets of tensors, where tensors are represented by block nodes and each axis assigned with by a categorical variable  $X_k$  represented by a node, see Figure 2a). Different simplifications of these factor graph depictions have been evolved in different research fields. In the tradition of graphical models, which started with the work Pearl (1988), the categorical variables are highlighted and the tensor blocks just depicted by hyperedges. To depict dependencies with causal interpretations, the edges are further decorated by directions in the depiction of Bayesian networks, see for example Pearl (2009).



In the tensor network community on the other hand, a simplification scheme highlighting the tensors as blocks and omitting the depiction of categorical variables has been evolved. The variables, or sometimes their index or dimension, are then directly assigned to the lines depicting the axes of the tensor blocks. This depiction scheme has been established in the literature as wiring diagrams (see Landsberg (2011)) and dates back at least to the work Penrose (1987).

Both depiction schemes are simplifications of factor graphs, by highlighting the categorical variables in the depiction in Figure 2b) and the tensors in the depiction in Figure 2c). We in this work will prefer the simplification of the tensor network community, depicted in Figure 2b).

In another interpretation (see Robeva and Seigal (2019)), both simplification schemes are different hypergraphs, which are dual to each other.

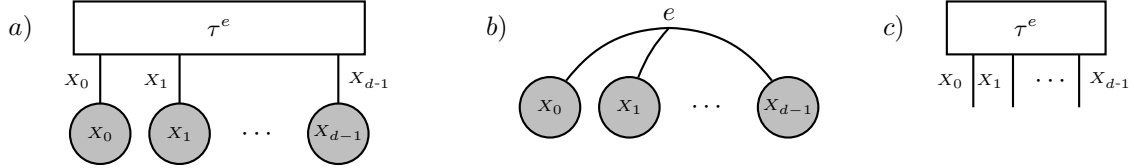


Figure 2: Depiction of Tensors a) As a factor in a factor graph, depicted by a block, and connected to categorical variables assigned to nodes. b) Highlighting only the variable dependencies by a hyperedge connecting the variables  $X_k$  to each axis  $k \in [d]$ . c) Highlighting the tensor by a blockwise notation with axes denoted by open legs represented by the variables  $X_k$ .

To depict vector calculus and its generalizations, we will apply the graphical notation (mainly version b) introduced in Chapter 3. Along this line, we represent vectors and their generalization to tensors by blocks with legs representing its indices. The basis vectors being one-hot encodings of states are in this scheme represented by



where  $\tilde{x}$  is an indexed represented by an open leg. Assigning  $x$  to this index will retrieve the  $x$ th coordinate (with value 1), whereas all other assignments will retrieve the coordinate values 0.

Drawing on the interpretation of tensors by hyperedges we can continue with the definition of tensor networks.

**Definition 8.** Let  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$  be a hypergraph with nodes decorated by categorical variables  $X_v$  with dimensions

$$m_v \in \mathbb{N}$$

and hyperedges  $e \in \mathcal{E}$  decorated by core tensors

$$\tau^e[X_e] \in \bigotimes_{v \in e} \mathbb{R}^{m_v},$$

where we denote by  $X_e$  the set of categorical variables  $X_v$  with  $v \in e$ . Then we call the set

$$\tau^{\mathcal{G}}[X_{\mathcal{V}}] = \{\tau^e[X_e] : e \in \mathcal{E}\}$$

the Tensor Network of the decorated hypergraph  $\mathcal{G}$ .

### 3.3.2 Tensor Product

Let us now exploit the developed graphical representations to define contractions of tensor networks. The simplest contraction is the tensor product, which maps a pair of two tensors with distinct variables onto a third tensor and has an interpretation by coordinatewise products. Such a contraction corresponds with a tensor network of two tensors with disjoint variables, depicted as:

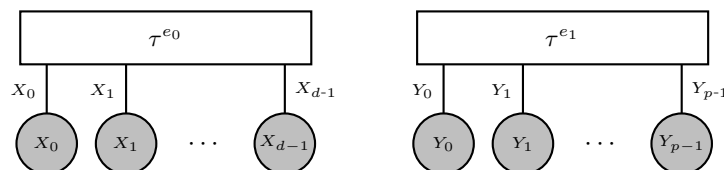




Figure 3: Example of a tensor network on a a) hypergraph with edges  $e_0 = \{X_0, X_1, X_2\}$ ,  $e_1 = \{X_1, X_2\}$  and  $e_2 = \{X_2, X_3\}$ , which is decorated by the tensor cores b), representing a contraction with leaving all variables open.

**Definition 9** (Tensor Product). *Let there be two tensor*

$$\tau^{e_0} [X_{[d]}] \in \bigotimes_{k \in [d]} \mathbb{R}^{m_k} \text{ and } \tau^{e_1} [Y_{[p]}] \in \bigotimes_{l \in [p]} \mathbb{R}^{m_l}$$

*with different categorical variables assigned to its axes. Then there tensor product is the map*

$$\langle \tau^{e_0} [X_{[d]}], \tau^{e_1} [Y_{[p]}] \rangle [X_{[d]}, Y_{[p]}] \in \left( \bigotimes_{k \in [d]} \mathbb{R}^{m_k} \right) \otimes \left( \bigotimes_{l \in [p]} \mathbb{R}^{m_l} \right)$$

*defined coordinatewise for tuples of  $x_0, \dots, x_{d-1} \in \times_{k \in [d]} [m_k]$  and  $y_0, \dots, y_{p-1} \in \times_{l \in [p]} [m_l]$  as*

$$\begin{aligned} \langle \tau, \tilde{\tau} \rangle [X_0 = x_0, \dots, X_{d-1} = x_{d-1}, Y_0 = y_0, \dots, Y_{p-1} = y_{p-1}] \\ := \tau^{e_0} [X_0 = x_0, \dots, X_{d-1} = x_{d-1}] \cdot \tau^{e_1} [Y_0 = y_0, \dots, Y_{p-1} = y_{p-1}]. \end{aligned}$$

Other popular standard notations of tensor products (see Kolda and Bader (2009); Hackbusch (2012); Cichocki et al. (2015))

$$(\tau \otimes \tilde{\tau}) = (\tau \circ \tilde{\tau}) = \langle \tau^{e_0} [X_{[d]}], \tau^{e_1} [Y_{[p]}] \rangle [X_{[d]}, Y_{[p]}].$$

We will avoid these notations in this work in favor of a consistent notation capable of depicting generic tensor network contractions.

When the tensor  $\tau^{e_1} [Y_{[p]}]$  coincides with the trivial tensor  $\mathbb{I} [Y_{[p]}]$  (see Example 1), we further make a notation convention to omit that tensor, that is

$$\langle \tau^{e_0} [X_{[d]}], \mathbb{I} [Y_{[p]}] \rangle [X_{[d]}, Y_{[p]}] = \langle \tau^{e_0} [X_{[d]}] \rangle [X_{[d]}, Y_{[p]}].$$

### 3.3.3 Generic Contractions

Contractions of Tensor Networks  $\tau^{\mathcal{G}}$  are operations to retrieve single tensors by summing products of tensors in a network over common indices. We will define contractions formally by specifying just the indices not to be summed over.

When some of the variables are not appearing as leg variables, we define the contraction as being a tensor product with the trivial tensor  $\mathbb{I}$  carrying the legs of the missing variables.

**Definition 10.** *Let  $\tau^{\mathcal{G}}$  be a tensor network on a decorated hypergraph  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ . For any subset  $\tilde{\mathcal{V}} \subset \mathcal{V}$  we define the contraction to be the tensor*

$$\langle \tau^{\mathcal{G}} \rangle [X_{\tilde{\mathcal{V}}}] \in \bigotimes_{v \in \tilde{\mathcal{V}}} \mathbb{R}^{m_v} \quad (3)$$

*defined coordinatewise by the sum*

$$\langle \tau^{\mathcal{G}} \rangle [X_{\tilde{\mathcal{V}}} = x_{\tilde{\mathcal{V}}}] = \sum_{x_{\mathcal{V}/\tilde{\mathcal{V}}} \in \times_{v \in \mathcal{V}/\tilde{\mathcal{V}}} [m_v]} \left( \prod_{e \in \mathcal{E}} \tau^e [X_e = x_e] \right). \quad (4)$$

We call  $X_{\tilde{\mathcal{V}}}$  the open variables of the contraction.

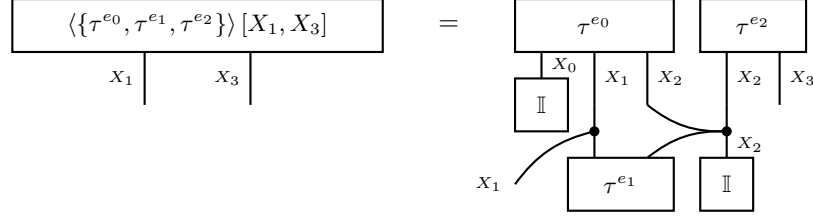


Figure 4: Example of a tensor network contraction of all but the variables  $X_1, X_3$ . Contraction of variables can always be depicted by closing the open legs with trivial tensors  $\mathbb{I}$  performing index sums.

To ease notation, we will often omit the set notation by brackets  $\{\cdot\}$  and specify the tensors to be contracted with the delimiter "," (see e.g. Example 2).

**Remark 2** (Alternative Notations). Contractions can also denoted by the Einstein summations of the indices along connected edges, understood as scalar product in each subspace. This is as in Def. 10, just omitting the sums. We found it useful in this work to do the diagrammatic representation instead, since it offers a better possibility to depict hierarchical arrangements of shared variables.

Further notations without usage of axis variables are mode products (see Kolda and Bader (2009); Hackbusch (2012); Cichocki et al. (2015)), often denoted by the operation  $\times_n$ . With our more generic variable-based notations, we can capture these more specific contractions by coloring the tensor axes, that is assignment of axis variables.

To further gain familiarity with the generic contractions, we show the connection to two more popular examples.

**Example 2. Matrix Vector Products** The matrix vector product is a special case of tensor contractions, where a matrix  $M[X_0, X_1]$  shares a categorical variable with a vector  $V[X_1]$ . When leaving the variable unique to the matrix open we get the matrix vector product as

$$\langle M[X_0, X_1], V[X_1] \rangle [X_0 = x_0] = \sum_{x_1 \in [m_1]} M[X_0 = x_0, X_1 = x_1] \cdot V[X_1 = x_1].$$

Exploiting the diagrammatic tensor network visualization we depict matrix vector by:

$$\begin{array}{c} X_0 \\ \text{---} \end{array} \boxed{M} \begin{array}{c} X_1 \\ \text{---} \end{array} \boxed{V} = \begin{array}{c} X_0 \\ \text{---} \end{array} \boxed{\langle M[X_0, X_1], V[X_1] \rangle [X_0]}$$

**Example 3.** Hadamard products of vectors  $A$  node appearing in arbitrary many hyperedges denotes a Hadamard product of the axis of the respective decorating tensors. To give an example, let  $V^k[X] \in \mathbb{R}^m$  be vectors for  $k \in [d]$ . Their hadamard product is the vector

$$\langle \{V^k[X] : k \in [d]\} \rangle [X] \in \mathbb{R}^m$$

defined by

$$\langle \{V^k[X] : k \in [d]\} \rangle [X = x] = \prod_{k \in [d]} V^k[X = x].$$

*In a contraction diagram the Hadamard product is depicted by:*

$$\boxed{\langle V^0[X], \dots, V^{d-1}[X] \rangle [X]} = \begin{array}{c} \boxed{V^0} \quad \boxed{V^1} \quad \dots \quad \boxed{V^{d-1}} \\ | \quad | \quad \quad \quad | \\ \text{---} \quad \text{---} \quad \quad \quad \text{---} \\ \quad \quad \quad \bullet \\ \quad \quad \quad | \\ \quad \quad \quad X \end{array}$$

### 3.3.4 Decompositions

Tensors can be represented by tensor network decompositions, when the contraction of the network retrieves the tensor.

**Definition 11.** A *Tensor Network Decomposition* of a tensor  $\tau [X_{\mathcal{V}}]$  is a Tensor Network  $\tau^{\mathcal{G}}$  such that

$$\tau[X_\nu] = \langle \tau^{\mathcal{G}} \rangle [X_\nu] \, .$$

We call the hypergraph  $\mathcal{G}$  the format of the decomposition.

### 3.3.5 Directed Tensors and normalizations

Directionality represents constraints on the structure of tensors, namely that the sum over outgoing trivializes the tensor.

**Definition 12.** A Tensor

$$\tau [X_{\mathcal{V}}] \in \bigotimes_{v \in \mathcal{V}} \mathbb{R}^{m_v}$$

is said to be directed with incoming variables  $\mathcal{V}^{\text{in}}$  and outgoing variables  $\mathcal{V}^{\text{out}}$ , where  $\mathcal{V} = \mathcal{V}^{\text{in}} \dot{\cup} \mathcal{V}^{\text{out}}$ , when

$$\langle \tau \rangle [X_{\mathcal{V}^{\text{out}}}] = \mathbb{I} [X_{\mathcal{V}^{\text{in}}}]$$

where  $\mathbb{I} [X_{\mathcal{V}^{\text{in}}}]$  denoted the trivial tensor in  $\bigotimes_{v \in \mathcal{V}^{\text{in}}} \mathbb{R}^{m_v}$  which coordinates are all 1.

While by default all legs are outgoing, we can change the direction by normalization.

**Definition 13.** A tensor  $\tau [X_{\mathcal{V}}]$  is said to be normable on  $\mathcal{V}^{\text{in}} \subset \mathcal{V}$ , if for any  $x_{\mathcal{V}^{\text{in}}} \in \times_{v \in \mathcal{V}^{\text{in}}} [m_v]$  we have

$$\langle \tau [X_{\mathcal{V}}], \epsilon_{x_{\mathcal{V}^{\text{in}}}} [X_{\mathcal{V}^{\text{in}}}] \rangle [\emptyset] > 0.$$

The normalization of a on  $\mathcal{V}^{\text{in}} \subset \mathcal{V}$  normable tensor is the tensor

$$\langle \tau [X_{\mathcal{V}}] \rangle [X_{\mathcal{V}^{\text{out}}} | X_{\mathcal{V}^{\text{in}}}] = \sum_{x_{\mathcal{V}^{\text{in}}} \in \times_{v \in \mathcal{V}^{\text{in}}} [m_v]} \epsilon_{x_{\mathcal{V}^{\text{in}}}} [X_{\mathcal{V}^{\text{in}}}] \otimes \frac{\langle \tau [X_{\mathcal{V}}], \epsilon_{x_{\mathcal{V}^{\text{in}}}} [X_{\mathcal{V}^{\text{in}}}] \rangle [X_{\mathcal{V}^{\text{out}}}]}{\langle \tau [X_{\mathcal{V}}], \epsilon_{x_{\mathcal{V}^{\text{in}}}} [X_{\mathcal{V}^{\text{in}}}] \rangle [\emptyset]}$$

where  $\mathcal{V}^{\text{out}} = \mathcal{V} / \mathcal{V}^{\text{in}}$ .

We will investigate the contractions of directed tensors in Part III, where we show in Theorem 100 that normalizations are directed tensors.

In our graphical tensor notation, we depict directed tensors by directed hyperedges (a), which are decorated by directed tensors (b), for example:



### 3.4 Function encoding schemes

Tensors are defined here as real-valued functions on the state set of a system described by categorical variables. We provide further schemes to represent functions in order to perform sparse calculus and to handle more generic functions.

#### 3.4.1 Basis encodings

Let us now show how we can encode maps between factored systems. The scheme is described in more generality and detail (encoding of subsets and relations) in Chapter 17, see Def. 84.

**Definition 14** (Relation encoding of maps between Factored Systems). Let  $q$  be a function

$$q : \times_{k \in [d]} [m_k] \rightarrow \times_{l \in [r]} [m_l]$$

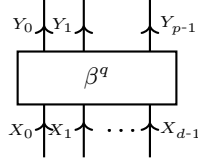
which maps the states of a factored system to variables  $X_0, \dots, X_{d-1}$  to the states of another factored system with variables  $Y_0, \dots, Y_{p-1}$ . Then the tensor representation of  $q$  is a tensor

$$\beta^f [X_0, \dots, X_{d-1}, Y_0, \dots, Y_{p-1}] \in \left( \bigotimes_{l \in [p]} \mathbb{R}^{m_l} \right) \otimes \left( \bigotimes_{k \in [d]} \mathbb{R}^{m_k} \right)$$

defined by

$$\begin{aligned} & \beta^f [Y_0, \dots, Y_{p-1}, X_0, \dots, X_{d-1}] \\ &= \sum_{x_0, \dots, x_{d-1} \in \times_{k \in [d]} [m_k]} \epsilon_{q(x_0, \dots, x_{d-1})} [Y_0, \dots, Y_{p-1}] \otimes \epsilon_{x_0, \dots, x_{d-1}} [X_0, \dots, X_{d-1}]. \end{aligned}$$

We depict basis encodings by directed tensors:



### 3.4.2 Tensor-valued functions

**Definition 15** (Selection encoding of Maps between Factored Systems). *Given a tensor space  $\bigotimes_{s \in [n]} \mathbb{R}^{p_s}$  described by categorical variables  $L_0, \dots, L_{n-1}$  and a tensor-valued function*

$$q : \bigtimes_{k \in [d]} [m_k] \rightarrow \bigotimes_{s \in [n]} \mathbb{R}^{p_s}$$

*the selection encoding of  $q$  is a tensor*

$$\sigma^q [X_{[d]}, L_{[n]}] \in \left( \bigotimes_{k \in [d]} \mathbb{R}^{m_k} \right) \otimes \left( \bigotimes_{s \in [n]} \mathbb{R}^{p_s} \right)$$

*defined by the basis decomposition*

$$\sigma^q [X_{[d]}, L_{[n]}] = \sum_{x_{[d]} \in \times_{k \in [d]} [m_k]} \epsilon_{x_0, \dots, x_{d-1}} [X_{[d]}] \otimes q(x_0, \dots, x_{d-1}) [L_{[n]}].$$

We call these tensor representation of maps selection encodings, since the coordinate of a function  $q$  to be processed is selected by another argument to  $\sigma^q$ .

We will provide more detail to the tensor representation of functions in Part III, where we distinguish between embeddings for basis and coordinate calculus.

## Part I

# Classical Approaches

## 4 Introduction into Part I

Within the introduced factored representation of systems we will in Part I present the probabilistic and logical approaches to artificial intelligence. Both the probabilistic and the logic paradigm provide a human-understandable interface to machine learning.

- **Probability:** Models describe dependencies between variables, which receive a graphical representation.
- **Logics:** Models are formulated in human interpretable logical syntax.

As we will describe in Part II, they can be combined in one formalism providing efficient reasoning. We will utilize that tensor network decompositions are in both useful tools of efficient calculus.

### 4.1 Representation of Factored Systems

**Probability** represents the uncertainty of states. The categorical variables are called random variables and their joint distribution is represented by a probability tensor. Humans interpret probabilities by Bayesian and frequentist approaches. Reasoning based on Bayes Theorem has an intuitive interpretation in terms of evidence based update of prior distributions to posterior distributions. However it is based on interpreting (large amounts) of numbers, which makes it hard for humans to assess the probabilistic reasoning process.

**Propositional Logics** explains relations between sets of worlds in a human understandable way. Categorical variables have dimension 2, where the first is interpreted as indicating a False state and the second as a True state. We mainly restrict to propositional logics, where there are finite sets of such variables called atomic formulas. Using model-theoretic semantics it defines entailment of sets by other sets, which is understandable as a consequence relation.

**Tensors** unify both approaches since they are natural numerical structures to represent properties of states in factored systems. The potential is then based in employing scalable multilinear algorithms to solve reasoning problems. Further, algorithms formulated in tensor networks have a high parallelization potential, which is why they are of central interest in the development of AI-dedicated software and hardware.

The different areas have developed separated languages to describe similar objects. Here we want to provide a rough comparison of those in a dictionary.

	Probability Theory	Propositional Logic	Tensors
<i>Atomic System</i>	Random Variable	Atomic Formula	Vector
<i>Factored System</i>	Joint Distribution	Knowledge Base	Tensor
<i>Categorical Variable</i>	Random Variable	Atomic Formula	Axis of the Tensor

While the probability theory lacks to provide an intuition about sets of events, propositional syntax has limited functionality to represent uncertainties. Tensors on the other side can build a bridge by representing both functionalities and relying on probability theory and logics for respective interpretations.

## 4.2 Mechanisms of tensor network decompositions

We investigate two mechanisms to identify tensor network decompositions of probability distributions:

- **Independence approach:** Conditional independence of random variables is a concept of probability theory. The most prominent application of this approach is the motivation of graphical models, which we introduce in Chapter 5 as tensor networks.
- **Computation approach:** When there are sufficient statistics providing probabilities, we construct tensor networks decompositions by computation of the statistics. Whenever the functions to be computed are compositions of functions of lower numbers of arguments, we utilize these representations to construct tensor network decompositions. Such decomposition schemes are provided by logical syntax as we will exploit in Chapter ???. In probability theory, we will make use of this approach in the efficient representation of sufficient statistics.

## 5 Probability Distributions

In this chapter we will establish relations between the formalism of tensor networks and basic concepts of probability theory. We will first understand distributions as tensors and connect their marginalizations and conditionings to the tensor operations of contractions and normalizations. Then we discuss independence assumptions as examples of contraction equations, which lead to tensor network decompositions known as graphical models. We then treat more generic exponential families and investigate their representation as tensor networks.

### 5.1 Classical Properties of Distributions

To start, we first relate classical properties of distributions, such as independent variables, with the tensor network formalism.

#### 5.1.1 Probability Tensors

After having discussed how to represent states of factored systems by one-hot encodings, let us now take advantage of these representation by associating properties with these states. Let there be uncertainties of the assignments  $x_k$  to the categorical variables  $X_k$  of a factored system. We then understand  $X_k$  as random variables, which have a joint distribution defined by the uncertainties of the state assignments. To capture these uncertainties we now make use of the one-hot representation of factored systems.

**Definition 16** (Probability Distribution). *Let there be for each  $k \in [d]$  a categorical variable  $X_k$  taking values in  $[m_k]$ . A joint probability distribution of these categorical variables is a tensor*

$$\mathbb{P}[X_{[d]}] \in \bigotimes_{k \in [d]} \mathbb{R}^{m_k},$$

such that

$$\langle \mathbb{P} [X_0, \dots, X_{d-1}] \rangle [\emptyset] = 1$$

and for all  $x_{[d]} \in \times_{k \in [d]} [m_k]$

$$\mathbb{P} [X_{[d]} = x_{[d]}] \in [0, 1].$$

The probability tensor can be decomposed as the sum (see Lem. 25 in Part III for more details)

$$\mathbb{P} [X_{[d]}] = \sum_{x_{[d]} \in \times_{k \in [d]} [m_k]} \mathbb{P} [X_{[d]} = x_{[d]}] \cdot \epsilon_{x_{[d]}} [X_{[d]}],$$

where we understand  $\mathbb{P} [X_{[d]} = x_{[d]}]$  as the probability of the categorical variables to take the state  $x_{[d]} \in \times_{k \in [d]} [m_k]$ .

The normalization condition  $1 = \langle \mathbb{P} [X_{[d]}] \rangle [\emptyset]$  has a more convenient equivalence by the coordinate sum

$$1 = \langle \mathbb{P} [X_{[d]}] \rangle [\emptyset] = \sum_{x_{[d]} \in \times_{k \in [d]} [m_k]} \mathbb{P} [X_{[d]} = x_{[d]}],$$

and thus ensures that all probabilities sum to 1, which is necessary for the probabilistic interpretation. While the assumptions of non-negative coordinates in Def. 16 reflects the first probability axiom of Kolmogorov, the assumption of contraction 1 implements the second axiom (see for example DeGroot (2016)). Since probability distributions contract to 1, they are directed (see Def. 12) with all distributed variables outgoing and empty incoming variables (see Figure 5).

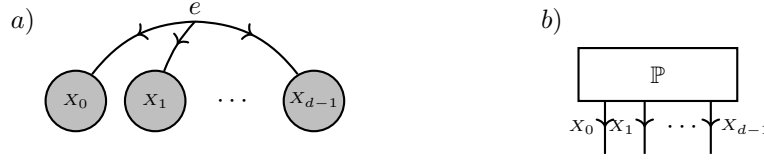


Figure 5: Probability distributions of variables  $X_0, \dots, X_{d-1}$ , sketched a) by a directed edge  $e$  with all variables outgoing, which is decorated b) by a directed tensor  $\mathbb{P} [X_{[d]}]$ .

### 5.1.2 Base measures

From a measure theoretic perspective, probabilities are measurable functions called probability densities, which integrals are 1 (see for example DeGroot (2016)). In our case of finite dimensional state spaces of factored systems, we implicitly used the trivial tensor  $\mathbb{I} [X_{[d]}]$  as a base measure, which measures subsets of states by their cardinality and is therefore referred to as state counting base measure. The distribution tensors  $\mathbb{P} [X_{[d]}]$  can then be understood as probability densities with respect to this state counting base measure. We in this work will also consider more general base measures  $\nu [X_{[d]}]$ , which we restrict to be boolean, that is  $\nu [X_{[d]} = x_{[d]}] \in \{0, 1\}$  for all states  $x_{[d]}$ . When understanding  $\mathbb{P} [X_{[d]}]$  as a probability density with respect to  $\nu [X_{[d]}]$ , any probabilistic interpretation will be through the contraction  $\langle \mathbb{P} [X_{[d]}], \nu [X_{[d]}] \rangle [X_{[d]}]$  and the normalization condition reads as

$$\langle \mathbb{P} [X_{[d]}], \nu [X_{[d]}] \rangle [\emptyset] = 1.$$

Since we restrict to boolean base measures, the contraction effectively manipulates the tensor  $\mathbb{P}$  by setting the coordinates  $\mathbb{P} [X_{[d]} = x_{[d]}]$  to zero, when  $\nu [X_{[d]} = x_{[d]}] = 0$ . Therefore, multiple tensors  $\mathbb{P}$  will have the same probabilistic interpretation, when  $\nu [X_{[d]}] \neq \mathbb{I} [X_{[d]}]$ . To avoid this ambiguity, we introduce the notation of representability with respect to a base measure  $\nu$ , by demanding that such coordinates are zero.

**Definition 17.** We say that a probability distribution  $\mathbb{P}$  is representable with respect to a boolean base measure  $\nu$ , if for all  $x_{[d]}$  with  $\nu [X_{[d]} = x_{[d]}] = 0$  we have  $\mathbb{P} [X_{[d]} = x_{[d]}] = 0$ . We denote the set of by  $\nu$  representable distributions by  $\Gamma^{\delta, \nu}$ .

When a probability distribution  $\mathbb{P}$  is representable with respect to a boolean base measure  $\nu$ , we have the invariance

$$\mathbb{P} [X_{[d]}] = \langle \mathbb{P} [X_{[d]}], \nu [X_{[d]}] \rangle [X_{[d]}]$$

and can therefore safely ignore the base measures. This enables the characterization of by  $\nu$  representable distributions by

$$\Gamma^{\delta, \nu} = \left\{ \mathbb{P}[X_{[d]}] : \forall x_{[d]} \in \bigtimes_{k \in [d]} [m_k] : \mathbb{P}[X_{[d]} = x_{[d]}] \geq 0, \langle \mathbb{P}[X_{[d]}], \nu[X_{[d]}] \rangle [X_{[d]}] = \mathbb{P}[X_{[d]}] \right\}.$$

Starting with Chapter 7 we will further investigate boolean tensors and relate them with propositional formulas. In Chapter 8 we will connect the representation and positivity with respect to boolean base measures with the formalism of entailment. The notation  $\Gamma^{\delta, \nu}$  of by  $\nu$  representable distributions will later in Chapter 11 relate to minterm exponential families introduced therein.

We now investigate, which base measures  $\nu$  can be chosen for a probability distribution  $\mathbb{P}$ , such that  $\mathbb{P}$  is representable by  $\nu$ . Here we want to find a  $\nu$ , which is in a sense to be defined minimal amount the base measures, such that  $\mathbb{P}$  is representable with respect to them. For this minimality criterion we will develop in Chapter 8 orders based on entailment and show the minimality in The. 50. Here, we just introduce the minimality criterion as positivity of a distribution with respect to a base measure.

**Definition 18.** We say that a probability distribution  $\mathbb{P}[X_{[d]}]$  is positive with respect to a boolean base measure  $\nu[X_{[d]}]$ , if the distribution is representable by  $\nu$  (i.e.  $\langle \mathbb{P}, \nu \rangle [\emptyset] = 1$ ) and for all  $x_{[d]}$  with  $\nu[X_0 = x_0, \dots, X_{d-1} = x_{d-1}] = 1$  we have  $\mathbb{P}[X_0 = x_0, \dots, X_{d-1} = x_{d-1}] > 0$ .

### 5.1.3 Marginal Distribution

Contractions of probability distributions are related to marginalizations as we introduce next.

**Definition 19** (Marginal Probability). Given a distribution  $\mathbb{P}[X_0, X_1]$  of the categorical variables  $X_0$  and  $X_1$  the marginal distribution of the categorical variable  $X_0$  is the tensor

$$\mathbb{P}[X_0] \in \mathbb{R}^{m_0}$$

defined by the contraction

$$\mathbb{P}[X_0] = \langle \mathbb{P}[X_0, X_1] \rangle [X_1].$$

To connect with a more standard defining equation of marginal distributions, let us notice that for any  $x_0 \in [m_0]$

$$\mathbb{P}[X_0 = x_0] = \langle \mathbb{P}[X_0, X_1] \rangle [X_0 = x_0] = \sum_{x_1 \in [m_1]} \mathbb{P}[X_0 = x_0, X_1 = x_1].$$

Thus, each coordinate of the marginal distribution is the sum of the joint probability of compatible states. We say that the variable  $X_1$  is marginalized out, when building the marginal distribution  $\mathbb{P}[X_0]$  of  $X_0$ . Let us now justify this terminology and show, that any marginal distribution is a probability distribution as introduced in Def. 16.

**Theorem 1.** Any marginal distribution is a probability distribution.

*Proof.* We further have that any marginal distribution is normed, since by the commutativity of contractions (see for more details The. 131 in Part III)

$$\langle \mathbb{P}[X_0] \rangle [\emptyset] = \langle \langle \mathbb{P}[X_0, X_1] \rangle [X_0] \rangle [\emptyset] = \langle \mathbb{P}[X_0, X_1] \rangle [\emptyset] = 1.$$

Further any coordinate is non-negative, since it is a sum of non-negative coordinates. It follows from Def. 16, that any marginal distribution is a probability distribution.  $\square$

In a tensor network diagram we often represent variables  $X_1$  not appearing as open variables of a contraction as contracted with the trivial tensor  $\mathbb{I}[X_1]$ . Following this notation, we depict the marginal distribution in Def. 19 by

$$\begin{array}{c} \boxed{\mathbb{P}[X_0]} \\ \downarrow x_0 \end{array} = \begin{array}{c} \boxed{\mathbb{P}[X_0, X_1]} \\ \downarrow x_0 \quad \downarrow x_1 \\ \boxed{\mathbb{I}} \end{array}$$

Since we have shown, that marginal distributions are themselves probability distributions, they inherit the outgoing directionality in tensor network diagrams.

We notice, that Def. 19 generalizes to marginalizations of arbitrary sets of variables, when having a distribution  $\mathbb{P}[X_{[d]}]$  of an arbitrary number of categorical variables. It suffices for this to interpret  $X_0$  and  $X_1$  as collections of variables, which indices take the states of the respective factored systems.



### 5.1.4 Conditional Probabilities

Normalizations of probability distributions result in conditional distributions as we define next.

**Definition 20** (Conditional Probability). *Let  $\mathbb{P}[X_0, X_1]$  be a distribution of the categorical variables  $X_0$  and  $X_1$ , such that  $\mathbb{P}[X_0, X_1]$  is normable on  $\{X_1\}$ . Then the distribution of  $X_0$  conditioned on  $X_1$  is defined by*

$$\mathbb{P}[X_0|X_1] = \langle \mathbb{P}[X_0, X_1] \rangle [X_0|X_1] .$$



Figure 6: Depiction of conditional probability distributions a) by an edge with the incoming variable  $X_1$  and the outgoing variable  $X_0$ , which is decorated by b) the directed tensor  $\mathbb{P}[X_0|X_1]$ .

Since conditional probabilities are normalizations of probability tensors, they are directed and therefore depicted by directed hyperedges (see Figure 6). For any  $x_1 \in [m_1]$  we depict the slice  $\mathbb{P}[X_0|X_1 = x_1]$  defined by a normalization operation as

$$\mathbb{P}[X_0|X_1 = x_1] = \frac{\mathbb{P}[X_0, X_1] \epsilon_{x_1}}{\langle \mathbb{P}[X_0, X_1] \rangle [X_1 = x_1]} .$$

As we have done before for marginal distribution, we relate Def. 20 with a more convenient coordinatewise definition of conditional probabilities. For any indices  $x_0 \in [m_0]$  and  $x_1 \in [m_1]$  we have

$$\mathbb{P}[X_0 = x_0|X_1 = x_1] = \frac{\mathbb{P}[X_0 = x_0, X_1 = x_1]}{\langle \mathbb{P}[X_0, X_1] \rangle [X_1 = x_1]} = \frac{\mathbb{P}[X_0 = x_0, X_1 = x_1]}{\sum_{x_0 \in [m_0]} \mathbb{P}[X_0 = x_0, X_1 = x_1]} .$$

The distribution of  $X_0$  conditioned on  $X_1$  is the normed collection of slice of the probability distribution  $\mathbb{P}[X_0, X_1]$ . Each slice of the conditioned distribution with respect to incoming variables is a probability distribution itself, as we show next.

**Theorem 2.** *For any  $x_1 \in [m_1]$  the tensor  $\mathbb{P}[X_0|X_1 = x_1]$  is a probability tensor.*

*Proof.* As a normalization of a non-negative tensor, the conditional probability  $\mathbb{P}[X_0|X_1 = x_1]$  and any of its slices is also a non-negative tensor. Further, we have for any  $x_1 \in [m_1]$

$$\begin{aligned} \langle \mathbb{P}[X_0|X_1 = x_1] \rangle [\emptyset] &= \sum_{x_0 \in [m_0]} \mathbb{P}[X_0 = x_0|X_1 = x_1] \\ &= \frac{\sum_{x_0 \in [m_0]} \mathbb{P}[X_0 = x_0, X_1 = x_1]}{\sum_{x_0 \in [m_0]} \mathbb{P}[X_0 = x_0, X_1 = x_1]} \\ &= 1 , \end{aligned}$$

and therefore each slice is normed. We can visualize this calculation exploiting our diagrammatic notation as

$$\begin{array}{c}
 \boxed{\mathbb{P}[X_0|X_1 = x_1]} \\
 \downarrow x_0 \\
 \boxed{\mathbb{I}}
 \end{array}
 =
 \begin{array}{c}
 \boxed{\mathbb{P}[X_0|X_1]} \\
 \downarrow x_0 \quad \downarrow x_1 \\
 \boxed{\mathbb{I}} \quad \boxed{\epsilon_{x_1}}
 \end{array}
 =
 \frac{
 \begin{array}{c}
 \boxed{\mathbb{P}[X_0, X_1]} \\
 \downarrow x_0 \quad \downarrow x_1 \\
 \boxed{\mathbb{I}} \quad \boxed{\epsilon_{x_1}}
 \end{array}
 }{
 \begin{array}{c}
 \boxed{\mathbb{P}[X_0, X_1]} \\
 \downarrow x_0 \quad \downarrow x_1 \\
 \boxed{\mathbb{I}} \quad \boxed{\epsilon_{x_1}}
 \end{array}
 }
 = 1.$$

Since for any  $x_1 \in [m_1]$  the slice  $\mathbb{P}[X_0|X_1 = x_1]$  is non-negative and contracts to 1, we conclude that it is a probability distribution.  $\square$

We further show, that exactly the directed tensors with non-negative coordinates are conditional probability tensors.

**Theorem 3.** *Any tensor with non-negative coordinates is a conditional distribution tensor, if and only if it is directed with the condition variables incoming and the other outgoing.*

*Proof.* " $\Rightarrow$ ": By The. 2 a conditional probability tensor  $\mathbb{P}[X_0|X_1]$  is the normalization of a tensor and by The. 100 a directed tensor. Since probability tensors have only non-negative coordinates, their contractions with one-hot encodings also have only non-negative coordinates and also their normalizations.

" $\Leftarrow$ ": Conversely, let  $\tau[X_{\mathcal{V}}]$  be a directed tensor with  $\mathcal{V}^{\text{in}}$  incoming and  $\mathcal{V}^{\text{out}}$  outgoing and non-negative coordinates. Then

$$\mathbb{P}[X_{\mathcal{V}}] = \frac{1}{\prod_{v \in \mathcal{V}^{\text{in}}} m_v} \cdot \tau[X_{\mathcal{V}}] \quad (5)$$

is a probability tensor, since

$$\sum_{x_{\mathcal{V}^{\text{in}}}} \sum_{x_{\mathcal{V}^{\text{out}}}} \mathbb{P}[X_{\mathcal{V}} = x_{\mathcal{V}}] = \sum_{x_{\mathcal{V}^{\text{in}}}} \sum_{x_{\mathcal{V}^{\text{out}}}} \frac{1}{\prod_{v \in \mathcal{V}^{\text{in}}} m_v} \cdot \tau[X_{\mathcal{V}} = x_{\mathcal{V}}] = \sum_{x_{\mathcal{V}^{\text{in}}}} \frac{1}{\prod_{v \in \mathcal{V}^{\text{in}}} m_v} = 1.$$

The conditional probability  $\mathbb{P}[X_{\mathcal{V}^{\text{out}}}|X_{\mathcal{V}^{\text{in}}}]$  coincides with  $\tau$ , since

$$\begin{aligned}
 \mathbb{P}[X_{\mathcal{V}^{\text{out}}}|X_{\mathcal{V}^{\text{in}}} = x_{\mathcal{V}^{\text{in}}}] &= \frac{\mathbb{P}[X_{\mathcal{V}^{\text{out}}}, X_{\mathcal{V}^{\text{in}}} = x_{\mathcal{V}^{\text{in}}}]}{\sum_{x_{\mathcal{V}^{\text{out}}}} \mathbb{P}[X_{\mathcal{V}^{\text{out}}} = x_{\mathcal{V}^{\text{out}}}, X_{\mathcal{V}^{\text{in}}} = x_{\mathcal{V}^{\text{in}}}] } \\
 &= \frac{\tau[X_{\mathcal{V}^{\text{out}}}, X_{\mathcal{V}^{\text{in}}} = x_{\mathcal{V}^{\text{in}}}]}{\sum_{x_{\mathcal{V}^{\text{out}}}} \tau[X_{\mathcal{V}^{\text{out}}} = x_{\mathcal{V}^{\text{out}}}, X_{\mathcal{V}^{\text{in}}} = x_{\mathcal{V}^{\text{in}}}] } = \tau[X_{\mathcal{V}^{\text{out}}}, X_{\mathcal{V}^{\text{in}}} = x_{\mathcal{V}^{\text{in}}}],
 \end{aligned}$$

where in the last equation we used that the denominator is by definition trivial since  $\tau$  is normed.  $\square$

The. 3 specifies a broad class of tensors to represent conditional probabilities. In combination with The. 106, which states that basis encodings are directed, we get that any basis encoding of a function is a conditional probability tensor.

### 5.1.5 Bayes Theorem and the Chain Rule

So far, we have connected concepts of probability theory such as marginal and conditional probabilities with contractions and normalizations of tensors. We will now proceed to show that basic theorems of probability theory translate into more general contraction equations.

**Theorem 4** (Bayes Theorem). *For any probability distribution  $\mathbb{P}[X_0, X_1]$  with positive  $\mathbb{P}[X_1]$  we have*

$$\mathbb{P}[X_0, X_1] = \langle \mathbb{P}[X_0|X_1], \mathbb{P}[X_1] \rangle [X_0, X_1].$$

*Proof.* This theorem follows from the more generic contraction equation The. 101 to be shown in Chapter 16. We note that by positivity of  $\mathbb{P}[X_1]$ , the tensor network  $\mathbb{P}$  is normable with respect to  $X_1$ . The. 101 therefore implies choosing  $\mathcal{V} = \{0, 1\}$ ,  $\mathcal{V}^{\text{in}} = \{1\}$  and  $\mathcal{V}^{\text{out}} = \{0\}$ , that For our tensor

$$\begin{aligned}
 \mathbb{P}[X_0, X_1] &= \langle \langle \mathbb{P}[X_0, X_1] \rangle [X_0|X_1], \langle \mathbb{P}[X_0, X_1] \rangle [X_1] \rangle [X_0, X_1] \\
 &= \langle \mathbb{P}[X_0|X_1], \mathbb{P}[X_1] \rangle [X_0, X_1].
 \end{aligned}$$

$\square$

Following the insight of the Bayes The. 4, probability distributions of arbitrary numbers of variables can be decomposed as a contraction of conditional probabilities, as we show in the next theorem.

**Theorem 5** (Chain Rule). *For any probability distribution  $\mathbb{P} [X_{[d]}]$  we have*

$$\mathbb{P} [X_{[d]}] = \langle \{\mathbb{P} [X_0]\} \cup \{\mathbb{P} [X_k | X_0, \dots, X_{k-1}] : k \in [d], k \geq 1\} \rangle [X_{[d]}] ,$$

*provided that all conditional probability distributions exist.*

*Proof.* The claim can be derived by an iterative application of the Bayes The. 4 theorem. We will proof this statement in more generality in Chapter 16 as The. 102, deducing it from the generalization of the Bayes The. 4 by The. 101. The claim here then follows from The. 102 using  $\mathcal{V} = [d]$  and  $\tau [X_{\mathcal{V}}] = \mathbb{P} [X_{[d]}]$ , since

$$\begin{aligned} \mathbb{P} [X_{[d]}] &= \langle \{\mathbb{P} [X_0]\} \cup \{\langle \mathbb{P} [X_{[d]}] \rangle [X_k | X_0, \dots, X_{k-1}] : k \in [d], k \geq 1\} \rangle [X_{[d]}] \\ &= \langle \{\mathbb{P} [X_0]\} \cup \{\mathbb{P} [X_k | X_0, \dots, X_{k-1}] : k \in [d], k \geq 1\} \rangle [X_{[d]}] , \end{aligned}$$

□

We observe, that the chain rule provides a generic decomposition scheme of probability distributions into conditional distributions. The conditional distribution to  $k = d - 1$ , which appears in the chain decomposition, is in the same tensor space as the decomposed distribution  $\mathbb{P} [X_{[d]}]$ . To achieve our main goal of tensor network decompositions, which is an efficient storage format of the decomposed tensor, we need to further sparsify the appearing conditional probabilities (to be more precise, we aim at basis+ CP decompositions, to be introduced in Chapter 18). These simplification require additional assumptions on the distribution, which we will introduce in the next section.

### 5.1.6 Independence

Independence leads to severe sparsifications of conditional probabilities and is therefore the key assumption to gain sparse decompositions of probability distributions. Before showing such decomposition schemes, we first provide a coordinatewise definition of independent variables.

**Definition 21** (Independence). *We say that  $X_0$  is independent of  $X_1$  with respect to a distribution  $\mathbb{P} [X_0, X_1]$ , if for any values  $x_0 \in [m_0]$  and  $x_1$  the distribution satisfies*

$$\mathbb{P} [X_0 = x_0, X_1 = x_1] = \mathbb{P} [X_0 = x_0] \cdot \mathbb{P} [X_1 = x_1] .$$

*In this case we denote  $(X_0 \perp X_1)$ .*

We state next an equivalent independence criterion based on a contraction equation of probability distributions.

**Theorem 6** (Independence Criterion as a Contraction Equation). *The variable  $X_0$  is independent from  $X_1$  with respect to a probability distribution  $\mathbb{P} [X_0, X_1]$ , if and only if*

$$\mathbb{P} [X_0, X_1] = \langle \langle \mathbb{P} [X_0, X_1] \rangle [X_0] , \langle \mathbb{P} [X_0, X_1] \rangle [X_1] \rangle [X_0, X_1] .$$

*Proof.* By The. 1 we know that marginal probabilities are equivalent to contracted probability distributions, i.e.  $\mathbb{P} [X_0] = \langle \langle \mathbb{P} \rangle \rangle [X_0]$ . By orthogonality of one-hot encodings we have that

$$\forall x_0, x_1 : \quad \mathbb{P} [X_0 = x_0, X_1 = x_1] = \mathbb{P} [X_0 = x_0] \cdot \mathbb{P} [X_1 = x_1]$$

is equivalent to

$$\sum_{x_0} \sum_{x_1} \mathbb{P} [X_0 = x_0, X_1 = x_1] \cdot \epsilon_{x_0} [X_0] \epsilon_{x_1} [X_1] = \sum_{x_0} \mathbb{P} [X_0 = x_0] \cdot \mathbb{P} [X_1 = x_1] \cdot \epsilon_{x_0} [X_0] \epsilon_{x_1} [X_1] .$$

We reorder the summations and arrive at

$$\sum_{x_0, x_1} \mathbb{P} [X_0 = x_0, X_1 = x_1] \cdot \epsilon_{x_0, x_1} [X_0, X_1] = \left( \sum_{x_0} \mathbb{P} [X_0 = x_0] \epsilon_{x_0} [X_0] \right) \cdot \left( \sum_{x_1} \mathbb{P} [X_1 = x_1] \epsilon_{x_1} [X_1] \right)$$

which is by Lem. 25 equal to the claim

$$\mathbb{P} [X_0, X_1] = \langle \langle \mathbb{P} \rangle [X_0] , \langle \mathbb{P} \rangle [X_1] \rangle [X_0, X_1] .$$

□

Two jointly distributed variables are by The. 6 independent, if and only if their joint distribution  $\mathbb{P}[X_0, X_1]$  is the tensor product of marginal probabilities. Using tensor network diagrams we depict this property by

$$\begin{array}{c} \boxed{\mathbb{P}[X_0, X_1]} \\ \downarrow x_0 \quad \downarrow x_1 \end{array} = \begin{array}{c} \boxed{\mathbb{P}[X_0, X_1]} \\ \downarrow x_0 \quad \downarrow x_1 \end{array} \otimes \begin{array}{c} \boxed{\mathbb{P}[X_0, X_1]} \\ \downarrow x_0 \quad \downarrow x_1 \end{array} = \begin{array}{c} \boxed{\mathbb{P}[X_0]} \\ \downarrow x_0 \end{array} \otimes \begin{array}{c} \boxed{\mathbb{P}[X_1]} \\ \downarrow x_1 \end{array} .$$

Let us notice, that the assumption of independence reduces the degrees of freedom from  $m_0 \cdot m_1 - 1$  to  $(m_0 - 1) + (m_1 - 1)$ . The decomposition by marginal distributions furthermore exploits this reduced freedom and provides an efficient storage. Having a joint distribution of multiple variables, which disjoint subsets are independent, we can iteratively apply the decomposition scheme. As a result, the degrees of freedom scaling exponential in the number of distributed variables would be reduced to a linear scaling, by the assumption of independence.

Independence is, as we observed, a strong assumption, which is often too restrictive. It is furthermore an undesired property, when in a supervised learning scenario a target variable has to be predicted based on known feature variables. Conditional independence instead is a less demanding assumption, which still implies efficient tensor network decompositions schemes. We introduce conditional independence as independence of variables with respect to conditional distributions.

**Definition 22** (Conditional Independence). *Given a joint distribution of variables  $X_0, X_1$  and  $X_2$ , such that  $\mathbb{P}[X_2]$  is positive. We say that  $X_0$  is independent of  $X_1$  conditioned on  $X_2$  if for any states  $x_0 \in [m_0], x_1 \in [m_1]$  and  $x_2 \in [m_2]$*

$$\mathbb{P}[X_0 = x_0, X_1 = x_1 | X_2 = x_2] = \mathbb{P}[X_0 = x_0 | X_2 = x_2] \cdot \mathbb{P}[X_1 = x_1 | X_2 = x_2] .$$

In this case we denote  $(X_0 \perp X_1) | X_2$ .

Conditional independence stated in Def. 22 has a close connection with independence stated in Def. 21. To be more precise,  $X_0$  is independent of  $X_1$  conditioned on  $X_2$ , if and only if  $X_0$  is independent of  $X_1$  with respect to any slice  $\mathbb{P}[X_0, X_1 | X_2 = x_2]$  of the conditional distribution  $\mathbb{P}[X_0, X_1 | X_2]$ . Analogously to The. 6 for independence, we further find a decomposition criterion for conditional independence. Since conditional independence can be regarded as a property of conditional probabilities, this decomposition criterion also involves conditional probabilities.

**Theorem 7** (Conditional Independence as a Contraction Equation). *Given a distribution  $\mathbb{P}$  of variables  $X_0, X_1$  and  $X_2$ , the variable  $X_0$  is independent of  $X_1$  conditioned on  $X_2$ , if and only if the equation*

$$\mathbb{P}[X_0, X_1 | X_2] = \langle \mathbb{P}[X_0 | X_2], \mathbb{P}[X_1 | X_2] \rangle [X_0, X_1, X_2]$$

holds.

*Proof.* With the same argumentation as in the proof of The. 6, we notice that the contraction equation holds, if and only if for any  $x_0 \in [m_0], x_1 \in [m_1]$  and  $x_2 \in [m_2]$

$$\mathbb{P}[X_0 = x_0, X_1 = x_1 | X_2 = x_2] = \mathbb{P}[X_0 = x_0 | X_2 = x_2] \cdot \mathbb{P}[X_1 = x_1 | X_2 = x_2] .$$

This is equivalent to conditional independence by Def. 22. □

We can further exploit conditional independence to find tensor network decompositions of probabilities, as we show as the next corollary.

$$\begin{array}{c} \boxed{\mathbb{P}[X_0, X_1, X_2]} \\ \downarrow x_0 \quad \downarrow x_1 \quad \downarrow x_2 \end{array} = \begin{array}{c} \boxed{\mathbb{P}[X_0 | X_2]} \\ \downarrow x_0 \quad \downarrow x_2 \end{array} \quad \begin{array}{c} \boxed{\mathbb{P}[X_2]} \\ \downarrow x_2 \end{array} \quad \begin{array}{c} \boxed{\mathbb{P}[X_1 | X_2]} \\ \downarrow x_2 \quad \downarrow x_1 \end{array}$$

Figure 7: Diagrammatic visualization of the contraction equation in Cor. 1. Conditional independence of  $X_0$  and  $X_1$  given  $X_2$  holds if the contraction on the right side is equal to the probability tensor on the left side.

**Corollary 1.** *If and only if  $X_0$  is independent of  $X_1$  conditioned on  $X_2$  the probability distribution  $\mathbb{P}$  satisfies (see Figure 7)*

$$\mathbb{P}[X_0, X_1, X_2] = \langle \mathbb{P}[X_0 | X_2], \mathbb{P}[X_1 | X_2], \mathbb{P}[X_2] \rangle [X_0, X_1, X_2] .$$

*Proof.* With the Bayes The. 4 it holds that

$$\mathbb{P}[X_0, X_1, X_2] = \langle \mathbb{P}[X_0, X_1|X_2], \mathbb{P}[X_2] \rangle [X_0, X_1, X_2] .$$

Decomposing the first tensor in the contraction, The. 7 implies, that  $X_0$  is independent of  $X_1$  conditioned on  $X_2$ , if and only if

$$\mathbb{P}[X_0, X_1, X_2] = \langle \mathbb{P}[X_0|X_2], \mathbb{P}[X_1|X_2], \mathbb{P}[X_2] \rangle [X_0, X_1, X_2] .$$

□

Let us now recall our motivation of the study of conditional independence, namely to find sparsifications of conditional probabilities as those appearing in chain decompositions The. 5. As we state as the next theorem, such sparsifications follow from conditional independence.

**Theorem 8.** *Whenever  $X_0$  is independent of  $X_1$  given  $X_2$ , we have for any  $x_1 \in [m_1]$*

$$\mathbb{P}[X_0|X_1 = x_1, X_2] = \mathbb{P}[X_0|X_2] .$$

*Proof.* By the Bayes The. 4 we have for any indices to the variables

$$\mathbb{P}[X_0 = x_0|X_1 = x_1, X_2 = x_2] = \frac{\mathbb{P}[X_0 = x_0, X_1 = x_1|X_2 = x_2]}{\langle \mathbb{P}[X_0, X_1 = x_1|X_2 = x_2] \rangle [\emptyset]}$$

If  $X_0$  is independent of  $X_1$  given  $X_2$  it follows that

$$\begin{aligned} \mathbb{P}[X_0 = x_0|X_1 = x_1, X_2 = x_2] &= \frac{\mathbb{P}[X_0 = x_0|X_2 = x_2] \cdot \mathbb{P}[X_1 = x_1|X_2 = x_2]}{\langle \mathbb{P}[X_0, X_1 = x_1|X_2 = x_2] \rangle [\emptyset]} \\ &= \mathbb{P}[X_0 = x_0|X_2 = x_2] . \end{aligned}$$

□

Following our motivation of sparse decompositions, we now combine this result with the generic chain rule, to show Markov Chain decompositions.

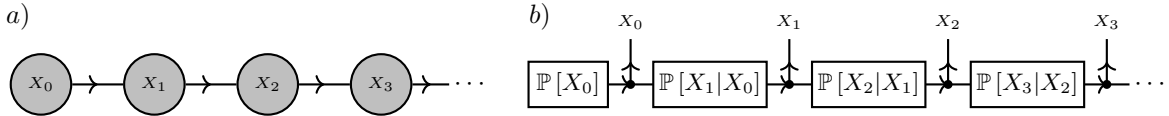


Figure 8: Depiction of a Markov Chain Decomposition by a a) hypergraph with the nodes  $\mathcal{V} = [d]$  and edges  $\mathcal{E} = \{\{0\} \cup \{\{k, k+1\} : k \in [d], k > 1\}\}$  and b) a decorating Tensor Network representing the sparsified conditional probabilities.

**Theorem 9 (Markov Chain).** *Let there be a set of variables  $X_k$  where  $k \in [d]$ , and let us denote for  $k \in [d]$  by  $X_{[k]}$  the collection of variables  $X_0, \dots, X_{k-1}$ . Let us assume, that for any  $k \in [d]$  with  $k \geq 2$  the variable  $X_k$  is independent of  $X_{[k-1]}$  conditioned on  $X_{k-1}$ , then*

$$\mathbb{P}[X_{[d]}] = \langle \{\mathbb{P}[X_0]\} \cup \{\mathbb{P}[X_k|X_{k-1}] : k \in [d], k \geq 1\} \rangle [X_{[d]}]$$

*We depict this decomposition in Figure 8.*

*Proof.* By the chain rule shown in The. 5 we have

$$\mathbb{P}[X_{[d]}] = \langle \{\mathbb{P}[X_k|X_{[k]}] : k \in [d]\} \rangle [X_{[d]}]$$

Using that  $X_k$  is conditional independent of  $X_{[k-1]}$  conditioned on  $X_{k-1}$  we further have by The. 8

$$\mathbb{P}[X_k|X_{[k]}] = \mathbb{P}[X_k|X_{k-1}] \otimes \mathbb{I}[X_{[k-1]}] .$$

Composing both equalities and omitting the trivial tensors shows the claim.

□

The assumption of  $X_k$  being independent of  $X_{[k-1]}$  conditioned on  $X_{k-1}$  is called the Markov property and the corresponding collection of random variables is called a Markov Chain. The. 9 states an efficient decomposition of the probability distribution into a concatenated product of matrices representing conditional probability distributions. Marginal distributions of Markov Chains can therefore consecutively be computed by matrix-vector products, that is for  $k \in [d]$  with  $k \geq 1$

$$\mathbb{P}[X_k] = \langle \mathbb{P}[X_k|X_{k-1}], \mathbb{P}[X_{k-1}] \rangle [X_k] .$$

The conditional probability matrices are therefore called stochastic transition matrices.

We notice, that the decomposition scheme of The. 9 hints at an efficient representation of  $\mathbb{P}[X_{[d]}]$  based on transition matrices. While  $\mathbb{P}[X_{[d]}]$  is a tensor in a space of dimension

$$\prod_{k \in [d]} m_k ,$$

the sum of the dimension of the transition matrices is

$$m_0 + \sum_{k \in [d], k \geq 1} m_k \cdot m_{k-1} .$$

We therefore observe a linear increase of the storage demand of the transition matrices in the order  $d$ , whereas a naive storage of  $\mathbb{P}[X_{[d]}]$  by its coordinates would have an exponentially demand.

The Markov Chain serves as a toy example drawing on a restrictive chain arrangement of conditional independencies. In the following section, we will investigate decomposition schemes, which relax this assumption and draw on more general collections of conditional independencies. The computation of marginal distribution by consecutive transition matrix multiplications will then be replaced by more general tensor network contractions.

## 5.2 Sufficient Statistics and Exponential Families

We have seen, that conditional independence of variables corresponds with decomposition properties of probability tensors. Another mechanism is through sufficient statistics, which leads to exponential families. When restricting to graphical models in the next section, we will see that both mechanisms are related through the Hammersley-Clifford theorem.

### 5.2.1 Sufficient Statistics

Let us consider a tuple of random variables  $X_{[d]}$ , which take values in  $\times_{k \in [d]} [m_k]$ . We now understand the probability  $\mathbb{P}[X_{[d]}]$  as another random variable taking values in  $[0, 1]$ , which has a deterministic dependence on  $X_{[d]}$ .

**Definition 23** (Sufficient Statistics). *Let  $X_{[d]}$  be a tuple of by  $\mathbb{P}[X_{[d]}]$  jointly distributed random variables and  $\mathcal{S}(X_{[d]}, L)$  be a tensor. We consider the tuple of random variables  $(X_{[d]}, \mathbb{P}[X_{[d]}], \mathcal{S}(X_{[d]}, L))$ , which takes for  $x_{[d]} \in \times_{k \in [d]} [m_k]$  with probability  $\mathbb{P}[X_{[d]} = x_{[d]}]$  the value*

$$(x_{[d]}, \mathbb{P}[X_{[d]} = x_{[d]}], \mathcal{S}(X_{[d]} = x_{[d]}, L)) .$$

We say, that  $\mathcal{S}$  is a sufficient statistic for  $\mathbb{P}$ , if this tuple obeys

$$(X_{[d]} \perp \mathbb{P}[X_{[d]}]) | \mathcal{S}(X_{[d]}, L) .$$

**Theorem 10.** *If and only if  $\mathcal{S}(X_{[d]}, L)$  is a sufficient statistic, i.e.  $X_{[d]}$  is independent of  $\mathbb{P}[X_{[d]}]$  conditioned on  $\mathcal{S}(X_{[d]}, L)$ , there is a tensor  $\alpha[Y_{[p]}]$  with*

$$\mathbb{P}[X_{[d]}] = \langle \beta^{\mathcal{S}}[Y_{[p]}, X_{[d]}], \alpha^{Y_{[p]}} \rangle [X_{[d]}] .$$

*Proof.* We exploit conditional entropies, for which definition we refer to Chapter 2 in ?. By the data processing inequality (see e.g. Theorem 2.8.1 in ?), we have

$$\mathbb{H}[\mathbb{P}|\mathcal{S}] \leq \mathbb{H}[\mathbb{P}|X_{[d]}] = 0$$

and thus  $\mathbb{H}[\mathbb{P}|\mathcal{S}] = 0$ . Moreover,  $\mathbb{H}[\mathbb{P}|\mathcal{S}] = 0$  is equivalent to a straight satisfaction of the data processing inequality, and  $\mathcal{S}$  being a sufficient statistic for  $\mathbb{P}$ .  $\mathbb{H}[\mathbb{P}|\mathcal{S}] = 0$  is further equivalent to the existence of a function  $q : \mathbb{R}^p \rightarrow [0, 1]$ , such that for each  $x_{[d]}$

$$q_{\mathcal{S}(X_{[d]}=x_{[d]}, L)} = \mathbb{P}[X_{[d]} = x_{[d]}] .$$

For an index interpretation function  $I$ , enumerating  $\times_{l \in [p]} \text{im}(S_l)$  using the variables  $Y_{[p]}$ , we define

$$\alpha := q \circ I.$$

Using basis calculus (see Chapter 17) and in particular The. 109 we have

$$\mathbb{P}[X_{[d]}] = \langle \beta^S[Y_{[p]}, X_{[d]}], \alpha^{Y_{[p]}} \rangle [X_{[d]}].$$

□

Given a statistic  $S$ , we can thus characterize the set of probability distributions, for which  $S$  is sufficient, as

$$\Lambda^{S, \mathcal{G}^{\max}} := \{ \langle \beta^S[Y_{[p]}, X_{[d]}], \alpha^{Y_{[p]}} \rangle [X_{[d]} | \emptyset] \}$$

## 5.2.2 Exponential families

The. 10 states the existence of an activation core, once a sufficient vector statistic has been identified. However, since the dimension of the activation core space is increasing exponential with the number of features (it is the product of the image cardinalities of the features), representation of generic  $\alpha$  is not feasible. We now restrict the activation cores to specific elementary tensors, which correspond with further assumptions on the dependence of  $\mathbb{P}$  and  $S$  made by exponential families.

The probability distributions, which are members of an exponential family, share the computation of the probability tensor based on a boolean base measure, marking the support of the distribution, and a statistic function containing features. They differ only by canonical parameters which weight the features at a given state to calculate the respective probability. Exponential families consist the most generic distributions investigated in this work and will also serve as a generic framework in the discussion of probabilistic reasoning in Chapter 6, as well as for neuro-symbolic models in Part II.

**Definition 24.** *Given a statistic function*

$$S : \times_{k \in [d]} [m_k] \rightarrow \mathbb{R}^p$$

*and a boolean base measure*

$$\nu : \times_{k \in [d]} [m_k] \rightarrow \{0, 1\}$$

*with  $\langle \nu \rangle [\emptyset] \neq 0$ , the set  $\Gamma^{S, \nu} = \{ \mathbb{P}^{(S, \theta, \nu)} : \theta[L] \in \mathbb{R}^p \}$  of probability distributions*

$$\mathbb{P}^{(S, \theta, \nu)} [X_{[d]}] = \langle \exp[\langle \sigma^S[X_{[d]}, L], \theta[L] \rangle [X_{[d]}], \nu[X_{[d]}] \rangle [X_{[d]} | \emptyset]$$

*is called the exponential family to  $S$ . We further define for each member with parameters  $\theta$  the associated energy tensor*

$$\phi^{(S, \theta, \nu)} [X_{[d]}] = \langle \sigma^S, \theta \rangle [X_{[d]}]$$

*and the cumulant function*

$$A^{(S, \nu)}(\theta) = \ln [\langle \nu, \exp[\langle \sigma^S, \theta \rangle [X_{[d]}]] \rangle [\emptyset]].$$

We used the selection encoding to represent the weighted summation over the statistics, that is the tensor (see Def. 15)

$$\sigma^S [X_{[d]}, L] : \times_{k \in [d]} [m_k] \times [p] \rightarrow \mathbb{R}$$

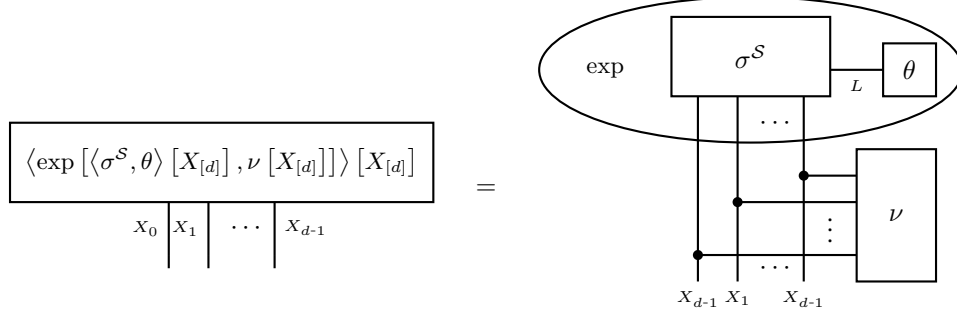
defined for  $x_{[d]} \in \times_{k \in [d]} [m_k]$  and  $l \in [p]$  as

$$\sigma^S [X_{[d]} = x_{[d]}, L = l] = S_l [X_{[d]} = x_{[d]}].$$

The selection encoding represent the weighted sum of the statistic coordinates by the canonical parameter vector  $\theta[L]$  as a contraction

$$\sum_{l \in [p]} \theta[L = l] \cdot S_l [X_{[d]}] = \langle \sigma^S [X_{[d]}, L], \theta[L] \rangle [X_{[d]}].$$

For more details on this representation scheme, we refer to The. 114 in Chapter 16. Up to normalization, we sketch the probability distribution of any member by the tensor network diagram



We here denote by an ellipsis the coordinatewise transformation by the exponential function (see Sect. 16.2). Since such coordinatewise transformation are nonlinear, they are a caveat for efficient contraction of the diagram.

Since we restrict the discussion to finite state spaces, the distribution  $\mathbb{P}^{(\mathcal{S}, \theta, \nu)}$  is well-defined for any  $\theta[L] \in \mathbb{R}^p$ . For infinite state space there are sufficient statistics and parameters, such that the partition function  $\langle \nu, \exp [\langle \sigma^S, \theta \rangle [X_{[d]}]] \rangle [\emptyset]$  diverges and the normalization  $\mathbb{P}^{(\mathcal{S}, \theta, \nu)}$  is not well-defined. In that cases, the canonical parameters need to be chosen from a subset where the partition function is finite Wainwright and Jordan (2008).

As before, we restrict to boolean base measures, which have to satisfy  $\langle \nu \rangle [\emptyset] \neq 0$  for respective distributions to exist. We notice, that by positivity of the exponential function, any distribution in an exponential family  $\Gamma^{\mathcal{S}, \nu}$  is positive with respect to  $\nu$  (see Def. 18). In Chapter 11 we will investigate distributions, where the base measures and the sufficient statistics share a common decomposition framework.

**Lemma 1.** *For any member of an exponential family  $\Gamma^{\mathcal{S}, \nu}$  we have*

$$\mathbb{P}^{(\mathcal{S}, \theta, \nu)} [X_{[d]}] = \left\langle \exp \left[ \phi^{(\mathcal{S}, \theta, \nu)} [X_{[d]}] - A^{(\mathcal{S}, \nu)}(\theta) \cdot \mathbb{I} [X_{[d]}] \right], \nu^{X_{[d]}} \right\rangle [X_{[d]}] .$$

*Proof.* By definition we have

$$\begin{aligned} \mathbb{P}^{(\mathcal{S}, \theta, \nu)} [X_{[d]}] &= \langle \exp [\langle \sigma^S, \theta \rangle [X_{[d]}]], \nu [X_{[d]}] \rangle [X_{[d]} | \emptyset] \\ &= \frac{\langle \exp [\langle \sigma^S, \theta \rangle [X_{[d]}]], \nu [X_{[d]}] \rangle [X_{[d]}]}{\langle \exp [\langle \sigma^S, \theta \rangle [X_{[d]}]], \nu [X_{[d]}] \rangle [\emptyset]} \\ &= \frac{\langle \exp [\phi^{(\mathcal{S}, \theta, \nu)} [X_{[d]}]], \nu [X_{[d]}] \rangle [X_{[d]}]}{\exp [A^{(\mathcal{S}, \nu)}(\theta)]} \\ &= \left\langle \exp \left[ \phi^{(\mathcal{S}, \theta, \nu)} [X_{[d]}] - A^{(\mathcal{S}, \nu)}(\theta) \cdot \mathbb{I} [X_{[d]}] \right], \nu^{X_{[d]}} \right\rangle [X_{[d]}] . \end{aligned}$$

□

A further useful criterion is that of minimality of an exponential family, as we define next.

**Definition 25 (Minimal).** *We say that a statistic  $\mathcal{S}$  is minimal with respect to a boolean base measure  $\nu$ , if there is no pair of a non-vanishing vector  $V[L]$  and a scalar  $\lambda \in \mathbb{R}$  with*

$$\langle \sigma^S [X_{[d]}, L], V[L], \nu [X_{[d]}] \rangle [X_{[d]}] = \lambda \cdot \nu [X_{[d]}] .$$

If a statistic is not minimal, we can omit coordinates of it without affecting the expressivity  $\Gamma^{\mathcal{S}, \nu}$ . As long as we find a non-vanishing vector  $V[L]$  and  $\lambda \in \mathbb{R}$  as in Def. 25, we can choose a coordinate  $\mathcal{S}_l$  such that  $V[L = l] \neq 0$ , conclude that the coordinate is linear dependent on the others and drop it as redundant.

### 5.2.3 Tensor Network Representation

As we have observed, the selection encoding formalism can efficiently represent the energy tensor to a member of an exponential family, but through coordinatewise transform by the exponential does not provide an efficient decomposition scheme of the probability distribution itself. We now overcome this problem with usage of the basis encoding formalism to represent members of exponential families by a single contraction without nonlinear transforms.



**Theorem 11** (Generic Representation of Exponential Families). *Given any base measure  $\nu$  and a sufficient statistic  $\mathcal{S}$  we enumerate for each coordinate  $l \in [p]$  the image  $\text{im}(\mathcal{S}_l)$  by a variable  $Y_l$  taking values in  $[\text{im}(\mathcal{S}_l)]$  (see for more details on this scheme Chapter 17), given an interpretation map*

$$I_l : [\text{im}(\mathcal{S}_l)] \rightarrow \text{im}(\mathcal{S}_l) .$$

*For any canonical parameter vector  $\theta[L] \in \mathbb{R}^p$  we build the activation cores*

$$\alpha^{l, \theta[L=l]}[Y_l = y_l] = \exp[\theta[L=l] \cdot I_l(y_l)]$$

*and have*

$$\mathbb{P}^{(\mathcal{S}, \theta, \nu)}[X_{[d]}] = \left\langle \{\nu[X_{[d]}]\} \cup \{\beta^{\mathcal{S}_l}[Y_l, X_{[d]}] : l \in [p]\} \cup \{\alpha^{\mathcal{S}_l, \theta[L=l]} : l \in [p]\} \right\rangle [X_{[d]} | \emptyset] .$$

*Proof.* We embed the image of  $\mathcal{S}$  in the cartesian product of the coordinate images

$$\text{im}(\mathcal{S}) \subset \bigtimes_{l \in [p]} \text{im}(\mathcal{S}_l)$$

and design enumerate the embedded image of  $\mathcal{S}$  by the variables  $Y_{[p]}$ . The. 109, to be shown in Chapter 17, implies

$$\exp[\langle \sigma^{\mathcal{S}}, \theta \rangle [X_{[d]}]] = \left\langle \beta^{\mathcal{S}}[Y_{[p]}, X_{[d]}], \exp[\langle \cdot, \theta[L] \rangle] |_{\bigtimes_{l \in [p]} \text{im}(\mathcal{S}_l)} \right\rangle [X_{[d]}] .$$

Here we denote by  $\langle \cdot, \theta[L] \rangle$  the dual function to  $\theta[L]$ , which assigns to vectors their contraction with  $\theta[L]$ . Its restriction onto the vectors in  $\bigtimes_{l \in [p]} \text{im}(\mathcal{S}_l)$  is the tensor satisfying

$$\exp[\langle \cdot, \theta \rangle] |_{\text{im}(\mathcal{S})} [Y_{[p]}] = \bigotimes_{l \in [p]} \exp[\cdot \theta[L=l]] |_{\text{im}(\mathcal{S}_l)} [Y_l] = \bigotimes_{l \in [p]} \alpha^{l, \theta[L=l]} [Y_l] .$$

We further have (see The. 110 in Chapter 17)

$$\beta^{\mathcal{S}}[Y_{[p]}, X_{[d]}] = \left\langle \{\beta^{\mathcal{S}_l}[Y_l, X_{[d]}] : l \in [p]\} \right\rangle [Y_{[p]}, X_{[d]}] .$$

Refining the above decomposition of  $\exp[\langle \sigma^{\mathcal{S}}, \theta \rangle [X_{[d]}]]$  by these further decompositions we arrive at the claim.  $\square$

In the proof of The. 11 we have observed, that the basis encoding  $\beta^{\mathcal{S}}[Y_{[p]}, X_{[d]}]$  of the statistics decomposed into a tensor network of basis encodings  $\beta^{\mathcal{S}_l}[Y_l, X_{[d]}]$  to the coordinate of the statistic. We can exploit further decomposition mechanisms, which will be discussed in full detail in Chapter 17, to find even sparser decompositions. This is for example the case, when the coordinates of the statistic are compositions of functions depending on small numbers of variables. When the coordinates of the statistic furthermore share similar parts in their compositions, these parts can be shared in the decomposition. We will investigate such sparsification mechanisms in more detail in Chapter 11, where the coordinates of the statistic are propositional formulas with a natural decomposition by their syntactical description.

The tensor network representation of an exponential family by The. 11 is a Markov Network consistent of two types of cores. First, we refer to the basis encodings  $\beta^{\mathcal{S}_l}$  of the coordinates of a statistic as computation cores. Our intuition is that they compute the hidden variable  $Y_l$ , based on Basis Calculus (see Chapter 17), which encode the value of the coordinate with respect to the image interpretation map  $I_l$ . We notice, that since they are directed with  $Y_l$  being the only outgoing variable, they do not influence any contraction with open variables  $X_{[d]}$ , unless further tensors sharing the variable  $Y_l$  are present in the contraction. The influence of the contraction is performed by the activation cores  $\alpha^{l, \theta[L=l]}[Y_l]$ , which exploit the computed statistic variable and provide in combination with the basis encoding a factor

$$\left\langle \beta^{\mathcal{S}_l}[Y_l, X_{[d]}], \alpha^{l, \theta[L=l]}[Y_l] \right\rangle [X_{[d]}]$$

to the Markov Network reduced to the observed variables  $X_{[d]}$ . When the canonical parameter is vanishing at a coordinate, that is  $\theta[L=l] = 0$ , then this factor is trivial, since  $\alpha^{l, 0}[Y_l] = \mathbb{I}[Y_l]$  and as a consequence of the directionality of basis encodings we have

$$\left\langle \beta^{\mathcal{S}_l}[Y_l, X_{[d]}], \alpha^{l, \theta[L=l]}[Y_l] \right\rangle [X_{[d]}] = \left\langle \beta^{\mathcal{S}_l}[Y_l, X_{[d]}], \mathbb{I}[Y_l] \right\rangle [X_{[d]}] = \mathbb{I}[X_{[d]}] .$$

In that case both the activation core and the corresponding computation core can be dropped from the network without changing its distribution.

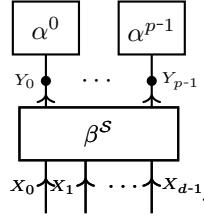
By The. 11 any member of an exponential family is represented by the normed contraction of a collection of unary activation cores contracted with the computation network  $\beta^S [Y_{[p]}, X_{[d]}]$ . We understand these activation cores as a member of a simple Markov Network distributing the head variables  $Y_{[p]}$ . This Markov Network has a graph, where the edges contain single variables, that is  $\mathcal{G}^{\text{EL}} = ([p], \{\{l\} : l \in [p]\})$ . We call this graph the elementary graph, since it also corresponds with elementary tensor network formats consistent of tensor products of vectors. A straightforward generalization of probability distributions representable by exponential families then allows for arbitrary decomposition formats for activation tensors, as we define next.

**Definition 26.** Given a statistic  $\mathcal{S} : \times_{k \in [d]} [m_k] \rightarrow \mathbb{R}^p$ , and a hypergraph  $\mathcal{G} = ([p], \mathcal{E})$  with nodes associated to the coordinates of the statistic, we define the by  $\mathcal{S}$  and  $\mathcal{G}$  computable family of distributions by

$$\Lambda^{\mathcal{S}, \mathcal{G}} = \left\{ \langle \{ \beta^{\mathcal{S}} [Y_{[p]}, X_{[d]}] \} \cup \{ \tau^e [Y_e] \} \rangle [X_{[d]} | \emptyset] : \tau^e [Y_e] \in \bigotimes_{l \in e} \mathbb{R}^{n_l}, 0 [Y_e] \prec \tau^e [Y_e] \right\}.$$

Note that we restrict to non-negative activation cores by demanding  $0 [Y_e] \prec \tau^e [Y_e]$ , a notation which will be introduced in more detail in Chapter 8 as partial order of tensors. We refer to any member  $\mathbb{P} [X_{[d]}] \in \Lambda^{\mathcal{S}, \mathcal{G}}$  as a by  $\mathcal{S}$  and  $\mathcal{G}$  computable distribution.

For unary activation cores, that is for the elementary graph  $\mathcal{G}^{\text{EL}}$ , any member of  $\Lambda^{\mathcal{S}, \mathcal{G}^{\text{EL}}}$  has up to a normalization factor a tensor network decomposition by the diagram

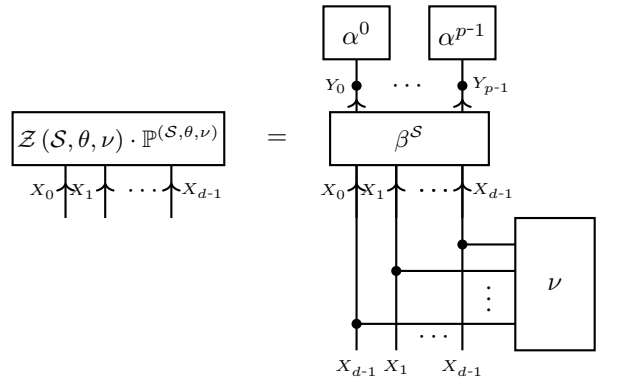


Comparing this representation scheme with The. 11, we conclude as the next corollary, that any member of an exponential family with trivial base measure can be represented by an elementary activation tensors.

**Corollary 2** (Corollary of The. 11). For any statistic  $\mathcal{S} : \times_{k \in [d]} [m_k] \rightarrow [p]$  and trivial base measure  $\nu [X_{[d]}] = \mathbb{I} [X_{[d]}]$  we have

$$\Gamma^{\mathcal{S}, \mathbb{I}} \subset \Lambda^{\mathcal{S}, \mathcal{G}^{\text{EL}}}.$$

For elements of the exponential family with general boolean base measure we have with the activation cores constructed in The. 11



where the partition function represents the normalizing contraction of the tensor network. Let us note, that when choosing activation cores with nontrivial support, we can also prepare boolean base measures and in principle extend Cor. 2 to families of nontrivial base measures. We will investigate such schemes later in Chapter 11, where we call them hybrid logic networks.

In comparison with the selection encoding representation of energy tensors, we have prepared a contraction without non-linear transforms, which represents the probability distributions being members of an exponential family.

However, relation encoding come with the expense of introducing more auxiliary variables compared with selection encodings. To be more precise, while selection encodings bundle the coordinates of the statistic in single selection variables, relation encodings create for each state  $l \in [p]$  of these selection variable an own auxiliary variable  $Y_l$ , which enumerated the image of the coordinate and can therefore be of high dimension. Thus, selection encodings offer in general a more efficient storage format coming at the expense of nonlinear operations in the computation of probabilities. We later will encounter situations, where selection encodings are feasible while relation encodings are not, when applying the formalism of formula selecting networks (see Chapter 10) in neuro-symbolic reasoning (see Chapter 12).

Based on Cor. 2 a further natural question is, whether  $\Gamma^{\mathcal{S}, \mathbb{I}}$  is a proper subset of  $\Lambda^{\mathcal{S}, \mathcal{G}^{\text{EL}}}$ . This is the case for most statistics  $\mathcal{S}$ , since members of exponential families are positive with respect to their base measure, which is in the corollaries setting trivial, while in  $\Lambda^{\mathcal{S}, \mathcal{G}^{\text{EL}}}$  we allow also for activation cores with vanishing coordinates, which in general do not produce positive distributions. The only statistics where  $\Gamma^{\mathcal{S}, \mathbb{I}}$  is not a proper subset of  $\Lambda^{\mathcal{S}, \mathcal{G}}$  are along this argumentation constant, since then the activation cores are one-dimensional vectors and vanishing coordinates are prohibited by the need for normalizability. We will follow these intuitions in the discussion of logical reasoning, starting with Chapter 8, and will use the formats  $\Lambda^{\mathcal{S}, \mathcal{G}}$  as hybrid formats storing probability distributions and logical knowledge bases.

While we have restricted our discussion on the elementary decomposition of the activation tensor, further decomposition schemes have interesting interpretations as well. Given a CP decomposition of the activation tensor (see for more details Chapter 18), the corresponding distributions are weighted mixture distributions built from the elementary decompositions. In general, the expressivity increases monotonously with the introduction of additional auxiliary variables and hyperedges in the representation format of activation tensors.

### 5.3 Graphical Models

**Specific instances of Exponential families are graphical models** Wainwright and Jordan (2008); Murphy (2022). They combine both the independence approach and the computation approach to tensor network representations of probability distributions.

Graphical models provide a more generic framework to relate conditional dependency assumptions on a distribution with tensor network decompositions. Following the tensor network formalism we in this section introduce graphical models based on hypergraphs. First, we study Markov Networks in most generality and then connect with conditional probabilities in the discussion of Bayesian Networks.

#### 5.3.1 Markov Networks

We now define Markov Networks based on hypergraphs, to establish a direct connection with tensor network decorating the hypergraph. In a more canonical way, Markov Networks are instead defined by graphs, where instead of the edges the cliques are decorated by factor tensors (see for example Koller and Friedman (2009)).

**Definition 27** (Markov Network). *Let  $\tau^{\mathcal{G}}$  be a tensor network of non-negative tensors decorating a hypergraph  $\mathcal{G}$ . Then the Markov Network  $\mathbb{P}^{\mathcal{G}}$  to  $\tau^{\mathcal{G}}$  is the probability distribution of  $X_{\mathcal{V}}$  defined by the tensor*

$$\mathbb{P}^{\mathcal{G}}[X_{\mathcal{V}}] = \frac{\langle \{\tau^e : e \in \mathcal{E}\} \rangle [X_{\mathcal{V}}]}{\langle \{\tau^e : e \in \mathcal{E}\} \rangle [\emptyset]} = \langle \tau^{\mathcal{G}} \rangle [X_{\mathcal{V}} | \emptyset] .$$

We call the denominator

$$\mathcal{Z}(\tau^{\mathcal{G}}) = \langle \{\tau^e : e \in \mathcal{E}\} \rangle [\emptyset]$$

the partition function of the tensor network  $\tau^{\mathcal{G}}$ .

The marginalization of a Markov Network to  $\tau^{\mathcal{G}}$  on subsets of variables  $X_{\tilde{\mathcal{V}}}$  is

$$\mathbb{P}^{\mathcal{G}}[X_{\tilde{\mathcal{V}}}] = \langle \tau^{\mathcal{G}} \rangle [X_{\tilde{\mathcal{V}}} | \emptyset] .$$

This can be derived from The. 131, which established an equivalence of contractions with sequences of consecutive contractions.

Further, the distribution of  $X_{\tilde{\mathcal{V}}}$  conditioned on  $X_{\bar{\mathcal{V}}}$ , where  $\tilde{\mathcal{V}}, \bar{\mathcal{V}}$  are disjoint subsets of  $\mathcal{V}$ , is

$$\mathbb{P}^{\mathcal{G}}[X_{\tilde{\mathcal{V}}} | X_{\bar{\mathcal{V}}}] = \langle \tau^{\mathcal{G}} \rangle [X_{\tilde{\mathcal{V}}} | X_{\bar{\mathcal{V}}}] .$$

While we have directly defined Markov Networks as decomposed probability distributions, we now want to derive assumptions on a distribution assuring that such decompositions exist. As we will see, the sets of conditional independencies encoded by a hypergraph are captured by its separation properties, as we define next.

**Definition 28** (Separation of Hypergraph). *A path in a hypergraph is a sequence of nodes  $v_k$  for  $k \in [d]$ , such that for any  $k \in [d - 1]$  we find a hyperedge  $e \in \mathcal{E}$  such that  $(v_k, v_{k+1}) \subset e$ . Given disjoint subsets  $A, B, C$  of nodes in a hypergraph  $\mathcal{G}$  we say that  $C$  separates  $A$  and  $B$  with respect to  $\mathcal{G}$ , when any path starting at a node in  $A$  and ending in a node in  $B$  contains a node in  $C$ .*

To characterize Markov Networks in terms of conditional independencies we need to further define the property of clique-capturing. This property of clique-capturing established a correspondence of hyperedges with maximal cliques in the more canonical graph-based definition of Markov Networks Koller and Friedman (2009).

**Definition 29** (Clique-Capturing Hypergraph). *We call a hypergraph  $\mathcal{G}$  clique-capturing, when each subset  $\tilde{\mathcal{V}} \subset \mathcal{V}$  is contained in a hyperedge, if for any  $a, b \in \tilde{\mathcal{V}}$  there is a hyperedge  $e \in \mathcal{E}$  with  $a, b \in \tilde{\mathcal{V}}$ .*

Let us now show a characterization of Markov Networks in terms of conditional independencies, which is analogous to The. 13.

**Theorem 12** (Hammersley-Clifford). *Given a clique-capturing hypergraph  $\mathcal{G}$ , the set of positive Markov Networks on the hypergraph coincides with the set of positive probability distributions, such that each for each disjoint subsets of variables  $A, B, C$  we have  $X_A$  is independent of  $X_B$  conditioned on  $X_C$ , when  $C$  separates  $A$  and  $B$  in the hypergraph.*

*Proof.* " $\Rightarrow$ ": Let there be a hypergraph  $\mathcal{G}$ , a Markov Network  $\tau^{\mathcal{G}}$  on  $\mathcal{G}$  and nodes  $A, B, C \subset \mathcal{V}$ , such that  $C$  separates  $A$  from  $B$ . Let us denote by  $\mathcal{V}_0$  the nodes with paths to  $A$ , which do not contain a node in  $C$ , and by  $\mathcal{V}_1$  the nodes with paths to  $B$ , which do not contain a node in  $C$ . Further, we denote by  $\mathcal{E}_0$  the hyperedges which contain a node in  $\mathcal{V}_0$  and by  $\mathcal{E}_1$  the hyperedges which contain a node in  $\mathcal{V}_1$ . By assumption of separability, both sets  $\mathcal{E}_0$  and  $\mathcal{E}_1$  are disjoint and no node in  $A$  is in a hyperedge in  $\mathcal{E}_1$ , respectively no node in  $B$  is in a hyperedge in  $\mathcal{E}_0$ . We then have

$$\begin{aligned} \langle \{\tau^e[X_e] : e \in \mathcal{E}\} [X_A, X_B | X_C = x_C] &= \langle \{\tau^e[X_e] : e \in \mathcal{E}\} \cup \{\epsilon_{x_C}\} [X_A, X_B | \emptyset] \\ &= \langle \{\tau^e : e \in \mathcal{E}_0\} \cup \{\epsilon_{x_C}\} [X_A | \emptyset] \\ &\quad \otimes \langle \{\tau^e : e \in \mathcal{E}_1\} \cup \{\epsilon_{x_C}\} [X_B | \emptyset] \rangle. \end{aligned}$$

By The. 7, it now follows that  $X_A$  is independent of  $X_B$  conditioned on  $X_C$ .

" $\Leftarrow$ ": The converse direction, i.e. that positive distributions respecting the conditional independence assumptions are representable as Markov Networks, is known as the Hammersley Clifford Theorem (see Clifford and Hammersley (1971)), which we will proof later in Sect. 16.4 of Chapter 16.  $\square$

From the proof of The. 12 Markov Networks with zero coordinates still satisfy the conditional independence assumption. However, the reverse is not true, that is there are distributions with vanishing coordinates, which satisfy the conditional independence assumptions, but cannot be represented as a Markov Network (see Example 4.4 in Koller and Friedman (2009)).

### 5.3.2 Bayesian Networks

Compared to Markov Networks, Bayesian Networks impose further conditions on tensor networks representing a distribution. They assume a directed hypergraph and each tensor decorating the edges to be normed according to the direction. We will observe, that if the hypergraph is in addition acyclic, then each tensor core coincides with the conditional distribution of the underlying Markov Network. To introduce Bayesian Networks, we extend Def. 7 by introducing the property of acyclicity for hypergraphs.

**Definition 30.** *A directed path is a sequence  $v_0, \dots, v_r$  such that for any  $l \in [r]$  there is an hyperedge  $e = (e^{\text{in}}, e^{\text{out}}) \in \mathcal{E}$  such that  $v_l \in e^{\text{in}}$  and  $v_{l+1} \in e^{\text{out}}$ . We call the hypergraph  $\mathcal{G}$  acyclic, if there is no path with  $r > 0$  such that  $v_0 = v_r$ . Given a directed hypergraph  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$  we define for any node  $v \in \mathcal{V}$  its parents by*

$$\text{Pa}(v) = \{\tilde{v} : \exists e = (e^{\text{in}}, e^{\text{out}}) \in \mathcal{E} : \tilde{v} \in e^{\text{in}}, v \in e^{\text{out}}\}$$

*and its non-descendants  $\text{NonDes}(v)$  as the set of nodes  $\tilde{v}$ , such that there is no directed path from  $v$  to  $\tilde{v}$ .*

Based on these additional graphical properties, we now define Bayesian Networks.

**Definition 31** (Bayesian Network). *Let  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$  be a directed acyclic hypergraph with edges of the form*

$$\mathcal{E} = \{(\text{Pa}(v), \{v\}) : v \in \mathcal{V}\}.$$

*A Bayesian Network is a decoration of each edge  $(\text{Pa}(v), \{v\})$  by a conditional probability distribution*

$$\mathbb{P}[X_v | X_{\text{Pa}(v)}]$$

which represents the probability distribution

$$\mathbb{P}[X_{\mathcal{V}}] = \langle \{ \mathbb{P}[X_v | X_{\text{Pa}(v)}] : v \in \mathcal{V} \} \rangle [X_{\mathcal{V}}] .$$

By definition each tensor decorating a hyperedge is directed with  $X_{\text{Pa}(v)}$  incoming and  $X_v$  outgoing. Thus, the directionality of the hypergraph is reflected in each tensor decorating a directed hyperedge. This allows us to verify with The. 103 that their contraction defines a probability distribution.

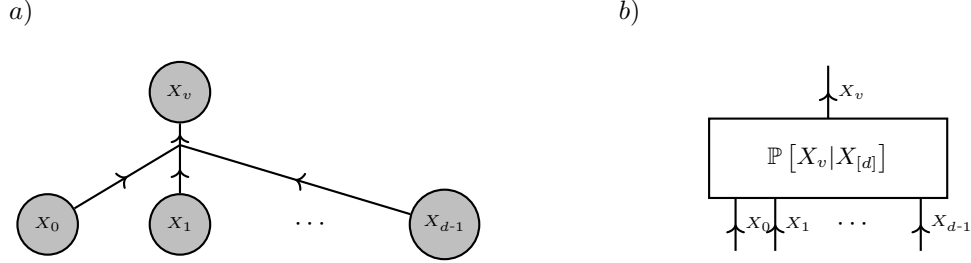


Figure 9: Example of a Factor of a Bayesian Network to the node  $X_v$  with parents  $X_0, \dots, X_{d-1}$ , as an directed edge a) which is decorated by a directed tensor b).

Marginalization of a Bayesian Network are still Bayesian Networks on a graph where the edges directing to variables, which are not marginalized over, are replaced by directed edges to the children. Conditioned Bayesian Network do not have a simple Bayesian Network representation, which is why we will treat them as Markov Networks to be introduced next.

**Theorem 13** (Independence Characterization of Bayesian Networks). *A probability distribution  $\mathbb{P}[X_{\mathcal{V}}]$  has a representation by a Bayesian Network on a directed acyclic graph  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ , if and only if for any  $v \in \mathcal{V}$  the variables  $X_v$  are independent on  $\text{NonDes}(v)$  conditioned on  $\text{Pa}(v)$ .*

*Proof.* We choose a topological order  $\prec$  on the nodes of  $\mathcal{G}$ , which exists since  $\mathcal{G}$  is acyclic.

" $\Rightarrow$ ": Let us assume, that the conditional independencies are satisfied and apply the chain rule with respect to that ordering to get

$$\mathbb{P}[X_{\mathcal{V}}] = \langle \mathbb{P}[X_v | X_{\tilde{v}} : \tilde{v} \prec v] \rangle [X_{\mathcal{V}}] .$$

Since  $\prec$  is a topological ordering we have

$$\text{Pa}(v) \subset \{ \tilde{v} : \tilde{v} \prec v \}$$

We apply the assumed conditional independence with The. 8 and get

$$\mathbb{P}[X_{\mathcal{V}}] = \langle \mathbb{P}[X_v | X_{\text{Pa}(v)}] \rangle [X_{\mathcal{V}}] .$$

" $\Leftarrow$ ": To show the converse direction, let there be a Bayesian Network  $\mathbb{P}[X_{\mathcal{V}}]$  on  $\mathcal{G}$ . To show for any node  $v$ , that  $X_v$  is independent of  $\text{NonDes}(v)$  conditioned on  $\text{Pa}(v)$ , we reorder the tensors in the contraction

$$\begin{aligned} & \mathbb{P}[X_v, X_{\text{NonDes}(v)} | X_{\text{Pa}(v)} = x_{\text{Pa}(v)}] \\ &= \langle \{ \mathbb{P}[X_{\tilde{v}} | X_{\text{Pa}(\tilde{v})}] : \tilde{v} \in \mathcal{V} \} \rangle [X_v, X_{\text{NonDes}(v)} | X_{\text{Pa}(v)} = x_{\text{Pa}(v)}] \\ &= \langle \{ \mathbb{P}[X_{\tilde{v}} | X_{\text{Pa}(\tilde{v})}] : \tilde{v} \in \mathcal{V} \cup \{ \epsilon_{x_{\text{Pa}(v)}} \} \} \rangle [X_v, X_{\text{NonDes}(v)} | \emptyset] \\ &= \langle \{ \mathbb{P}[X_{\tilde{v}} | X_{\text{Pa}(\tilde{v})}] : \tilde{v} \in \text{NonDes}(v) \} \cup \{ \epsilon_{x_{\text{Pa}(v)}}, \mathbb{P}[X_v | X_{\text{Pa}(v)}] \} \rangle [X_v, X_{\text{NonDes}(v)} | \emptyset] \\ &= \langle \{ \mathbb{P}[X_{\tilde{v}} | X_{\text{Pa}(\tilde{v})}] : \tilde{v} \in \text{NonDes}(v) \} \cup \{ \epsilon_{x_{\text{Pa}(v)}} \} \rangle [X_{\text{NonDes}(v)} | \emptyset] \\ &\quad \cdot \langle \{ \mathbb{P}[X_v | X_{\text{Pa}(v)}], \epsilon_{x_{\text{Pa}(v)}} \} \rangle [X_v | \emptyset] \\ &= \langle \{ \mathbb{P}[X_{\text{NonDes}(v)} | X_{\text{Pa}(v)} = x_{\text{Pa}(v)}], \mathbb{P}[X_v | X_{\text{Pa}(v)} = x_{\text{Pa}(v)}] \} \rangle [X_v, X_{\text{NonDes}(v)}] \end{aligned}$$

Here we have dropped in the third equation all tensors to the descendants, since their marginalization is trivial (which can be shown by a leaf-stripping argument). In the fourth equation we made use of the fact, that any directed path between the non-descendants and the node is through the parents of the node. By The. 7, it now follows that  $X_v$  is independent of  $\text{NonDes}(v)$  conditioned on  $\text{Pa}(v)$ .  $\square$

### 5.3.3 Bayesian Networks as Markov Networks

Markov Networks are more flexible compared with Bayesian Networks, since any Bayesian Network is a Markov Network by ignoring the directionality of the hypergraph and understanding the conditional distributions as generic tensor cores. In the next theorem we provide the conditions for the interpretation of a Markov Network as a Bayesian Network.

**Theorem 14.** *Let  $\tau^{\mathcal{G}}$  be a tensor network on a directed acyclic hypergraph, such that the edges are of the structure*

$$\mathcal{E} = \{(\text{Pa}(v), \{v\}) : v \in \mathcal{V}\}$$

*and each tensor  $\tau^e$  respects the directionality of the graph, that is each  $\tau^{(\text{Pa}(v), \{v\})}$  is directed with the variables to  $\text{Pa}(v)$  incoming and  $v$  outgoing. Then  $\mathcal{Z}(\tau^{\mathcal{G}}) = 1$  and for each  $v \in \mathcal{V}$  we have*

$$\tau^{(\text{Pa}(v), \{v\})} = \langle \tau^{\mathcal{G}} \rangle [X_v | X_{\text{Pa}(v)}] .$$

*In particular,  $\tau^{\mathcal{G}}$  is a Bayesian Network.*

*Proof.* We show the claim by induction over the cardinality of  $\mathcal{V}$ .

$|\mathcal{V}| = 1$ : In this case we find a unique node  $v \in \mathcal{V}$  and have  $\mathcal{E} = \{(\emptyset, \{v\})\}$ . The tensor  $\tau^{(\emptyset, \{v\})}$  is then normed with no incoming variables and we thus have

$$\mathcal{Z}(\tau^{\mathcal{G}}) = \langle \tau^{\mathcal{G}} \rangle [\emptyset] = \langle \tau^{(\emptyset, \{v\})} \rangle [\emptyset] = 1$$

and

$$\langle \tau^{\mathcal{G}} \rangle [X_v | \emptyset] = \tau^{(\emptyset, \{v\})} .$$

$|\mathcal{V}| - 1 \rightarrow |\mathcal{V}|$ : Let there now be a directed hypergraph  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$  and let us now assume, that the theorem holds for any tensor networks with node cardinality  $|\mathcal{V}| - 1$ . Since the hypergraph is acyclic, we find a root  $v \in \mathcal{V}$  such that  $v \notin \text{Pa}(\tilde{v})$  for  $\tilde{v} \in \mathcal{V}$ . We denote  $\tau^{\tilde{\mathcal{G}}}$  the tensor network on the hypergraph  $\tilde{\mathcal{G}} = \{\mathcal{V}/\{v\}, \mathcal{E}/\{(\text{Pa}(v), \{v\})\}\}$  with decorations inherited from  $\tau^{\mathcal{G}}$ . With Theorem 131, the directionality of  $\tau^{(\text{Pa}(v), \{v\})}$  and the induction assumption on  $\tau^{\tilde{\mathcal{G}}}$  we have

$$\langle \tau^{\tilde{\mathcal{G}}} \cup \{ \tau^{(\text{Pa}(v), \{v\})} \} \rangle [\emptyset] = \langle \tau^{\tilde{\mathcal{G}}} \cup \{ \langle \tau^{(\text{Pa}(v), \{v\})} \rangle [X_{\text{Pa}(v)}] \} \rangle [\emptyset] = \langle \tau^{\tilde{\mathcal{G}}} \cup \{ \mathbb{I}[X_{\text{Pa}(v)}] \} \rangle [\emptyset] = 1$$

and thus a trivial partition function. Since  $v$  does not appear in  $\tilde{\mathcal{G}}$ , we have for any index  $x_{\text{Pa}(v)}$

$$\langle \tau^{\mathcal{G}} \rangle [X_v, X_{\text{Pa}(v)} = x_{\text{Pa}(v)}] = \langle \tau^{(\text{Pa}(v), \{v\})} \rangle [X_v, X_{\text{Pa}(v)} = x_{\text{Pa}(v)}] \cdot \langle \tau^{\tilde{\mathcal{G}}} \rangle [X_{\text{Pa}(v)} = x_{\text{Pa}(v)}]$$

and thus, since  $\tau^{(\text{Pa}(v), \{v\})}$  is directed, that

$$\langle \tau^{\mathcal{G}} \rangle [X_v | X_{\text{Pa}(v)}] = \tau^{(\text{Pa}(v), \{v\})} .$$

□

Theorem 14 states that Bayesian Networks are a subset of Markov Networks. While Markov Network allow generic tensor cores, Bayesian Networks impose a local directionality condition on each tensor core by demanding it to be a conditional probability tensor. In our diagrammatic notation, the local normalization of Bayesian Networks is highlighted by the directionality of the hypergraph. Generic Markov Networks are on undirected hypergraphs, where in general no local directionality condition is assumed. As a consequence, tasks such as the determination of the partition functions or calculation of conditional distributions involve global contractions.

### 5.3.4 Hidden Markov Models

Hidden Markov Models are examples of Bayesian Networks, constructed as follows. Let us recall Markov Chains as investigated in The. 9 and extend them by observation variables  $E_k$  for  $k \in [d]$ , representing limited observations of the state variables  $X_k$ . To be more precise, we assume the following conditional independencies:

- As for Markov Chains, we assume that for  $k \in [d]$  with  $k \geq 1$  the variable  $X_k$  is independent of  $X_{[k-1]}$  and  $E_{[k-1]}$  conditioned on  $X_{k-1}$

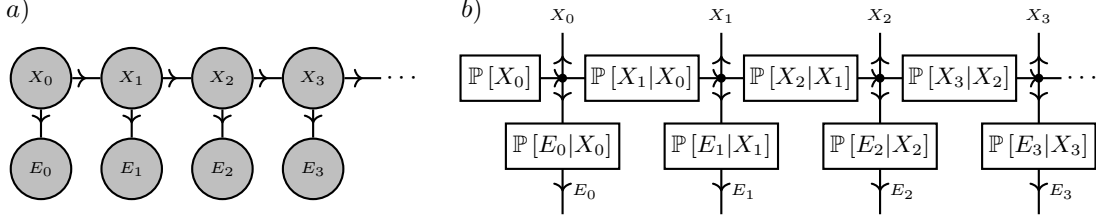


Figure 10: Decomposition of a probability distribution in the Hidden Markov Model, consistent of state variables  $X_k$  and observation variables  $E_k$ . Given the models conditional independence assumptions, the distribution is a Bayesian Network on the directed hypergraph a). The hypergraph is decorated by the network of conditional probability tensors b), which are interpreted as stochastic transition matrices  $\mathbb{P}[X_k|X_{k-1}]$  and stochastic observation matrices  $\mathbb{P}[E_k|X_k]$ .

- In addition, for we assume that for  $k \in [d]$  the observation variable  $E_k$  is independent of  $X_{[k]}$  and  $E_{[k]}$  conditioned on  $X_k$

From this conditional independence assumption, we apply the Chain Rule The. 5 given the order of variables

$$X_0, E_0, X_1, E_1, \dots, X_{d-1}, E_{d-1}$$

and get

$$\begin{aligned} \mathbb{P}[X_{[d]}, E_{[d]}] &= \langle \{\mathbb{P}[X_0], \mathbb{P}[E_0|X_0]\} \\ &\cup \{\mathbb{P}[X_k|X_{[k]}, E_{[d]}] : k \in [d]\} \\ &\cup \{\mathbb{P}[E_k|X_{[k+1]}, E_{[d]}] : k \in [d]\} \rangle [X_{[d]}, E_{[d]}]. \end{aligned}$$

We now apply the conditional independence assumptions to sparsify the appearing conditional distributions by application of The. 8. This results in the decomposition (see Figure 10b)

$$\begin{aligned} \mathbb{P}[X_{[d]}, E_{[d]}] &= \langle \{\mathbb{P}[X_0], \mathbb{P}[E_0|X_0]\} \\ &\cup \{\mathbb{P}[X_k|X_{k-1}] : k \in [d]\} \\ &\cup \{\mathbb{P}[E_k|X_k] : k \in [d]\} \rangle [X_{[d]}, E_{[d]}]. \end{aligned}$$

In addition to the stochastic transition matrices  $\mathbb{P}[X_k|X_{k-1}]$  appearing in Markov Chains, we further have stochastic observation matrices  $\mathbb{P}[E_k|X_k]$  for  $k \in [d]$ . Their contraction with marginal distribution of the respective state variables delivers the marginal distribution of the observation matrix by

$$\mathbb{P}[E_k] = \langle \mathbb{P}[E_k|X_k], \mathbb{P}[X_k] \rangle [E_k]$$

We notice, that this is a Bayesian Network on a directed acyclic hypergraph  $\mathcal{G}$  (see Figure 10a) consistent in nodes  $\{X_{[d]}\} \cup \{E_{[d]}\}$  to each state and observation variables, and the directed hyperedges by

- $(\emptyset, \{X_0\})$ , decorated by the initial marginal distribution of  $X_0$
- $(\{X_{k-1}\}, \{X_k\})$  for  $k \in [d]$  with  $k \geq 1$ , decorated by stochastic transition matrices
- $(\{X_k\}, \{E_k\})$  for  $k \in [d]$ , decorated stochastic observation matrices

While we have derived this directed graph structure directly based on the chain rule decomposition with sparsified conditional distributions, it also follows from the more generic hypergraph characterization of Bayesian Networks through separability by The. 13.

### 5.3.5 Markov Networks as Exponential Families

As we have claimed before, exponential families can be regarded as a generalization of graphical models. We here show this claim by a construction of exponential families representing Markov Networks on constant hypergraphs.

**Theorem 15** (Exponential Representation of Markov Networks). *Let there be a hypergraph  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$  with a coloring of the nodes by dimensions  $m_v$ , we define a sufficient statistics*

$$\mathcal{S}^{\mathcal{G}} : \times_{v \in \mathcal{V}} [m_v] \rightarrow \times_{e \in \mathcal{E}} \left( \times_{v \in e} [m_v] \right)$$

by a cartesian product  $\mathcal{S}^{\mathcal{G}} = \times_{e \in \mathcal{E}} \mathcal{S}_e$  of statistics

$$\mathcal{S}_e : \times_{v \in \mathcal{V}} [m_v] \rightarrow \times_{v \in e} [m_v]$$

defined by the restriction of indices to the respective edge, that is for  $x_{\mathcal{V}} \in \times_{v \in \mathcal{V}} [m_v]$

$$\mathcal{S}_e(x_{\mathcal{V}}) = x_e.$$

Given any Markov Network with positive tensors  $\{\tau^e : e \in \mathcal{E}\}$  decorating the hyperedges of  $\mathcal{G}$  we define

$$\theta[L] = \times_{e \in \mathcal{E}} \theta_e[X_e]$$

where

$$\theta_e[X_e] = \ln[\tau^e[X_e]]$$

and  $L$  enumerates a concatenation of the states of  $X_e$ . Then, the Markov Network is the member with canonical parameter  $\theta$  of the exponential family with trivial base measure, statistic  $\mathcal{S}^{\mathcal{G}}$ , which we denote by  $\Gamma^{\mathcal{S}^{\mathcal{G}}, \mathbb{I}}$ .

*Proof.* We have for any  $x_{\mathcal{V}}$

$$\begin{aligned} \prod_{e \in \mathcal{E}} \tau^e[X_e = x_e] &= \exp \left[ \sum_{e \in \mathcal{E}} \ln[\tau^e[X_e = x_e]] \right] \\ &= \exp \left[ \sum_{e \in \mathcal{E}} \theta_e[X_e = x_e] \right] \\ &= \exp \left[ \sum_{e \in \mathcal{E}} \langle \theta_e[X_e], \mathcal{S}_e(x_{\mathcal{V}}) \rangle [\emptyset] \right]. \end{aligned}$$

By contraction, we further have

$$\langle \mathcal{S}(X_{\mathcal{V}}, L), \theta[L] \rangle [X_{\mathcal{V}}] = \sum_{e \in \mathcal{E}} \langle \mathcal{S}_e[X_{\mathcal{V}}, X_e], \theta_e[X_e] \rangle [X_{\mathcal{V}}]$$

we thus get with the above

$$\langle \{\tau^e : e \in \mathcal{E}\} \rangle [X_{\mathcal{V}}] = \exp[\langle \mathcal{S}(X_{\mathcal{V}}, L), \theta[L] \rangle [X_{\mathcal{V}}]]. \quad (6)$$

This implies, that the contraction of the tensors in the Markov Network coincides with the exponential of the energy tensor of the constructed member of the exponential family. It follows for the normalization, that

$$\langle \{\tau^e : e \in \mathcal{E}\} \rangle [X_{\mathcal{V}} | \emptyset] = \langle \exp[\langle \theta, \mathcal{S} \rangle [X_{\mathcal{V}}]] \rangle [X_{\mathcal{V}} | \emptyset]. \quad (7)$$

We thus conclude, that the Markov Network coincides with the constructed member of the exponential family.  $\square$

The mean parameter of the Markov Network exponential family is the cartesian product of the marginals  $\mu_e[X_e]$ . They are often referred to as beliefs in the literature, as introduced by Pearl (1988). For Markov Networks on tree hypergraphs, and their embedding into junction tree formats, the corresponding mean parameter polytope can be characterized by local consistency constraints. More precisely, it can be shown, that for the statistic constructed in The. 15, in case of tree hypergraphs, the

$$\begin{aligned} \mathcal{M}_{\mathcal{S}^{\mathcal{G}}} = \left\{ \mu[L] = (\mu_e[X_e])_{e \in \mathcal{E}} : \forall e, \tilde{e} \in \mathcal{E} : \langle \mu_e[X_e] \rangle [X_{e \cap \tilde{e}}] = \langle \mu_{\tilde{e}}[X_{\tilde{e}}] \rangle [X_{e \cap \tilde{e}}], \right. \\ \left. \forall e : 0[X_e] \prec \mu_e[X_e] \wedge \langle \mu_e[X_e] \rangle [\emptyset] = 1 \right\}. \end{aligned}$$

That is, the polytope of realizable mean parameters consists of those non-negative and normed beliefs, which are coinciding on the contraction of shared variables. Capturing these constraints by Lagrange parameters and performing optimization of certain objectives then results in message-passing schemes, as we will discuss in Chapter 20. If the hypergraph is not minimally connected, this constructed polytope is only an outer bound of the true mean parameter polytope, but still serves as a motivation of loopy belief propagation schemes (see Chapter 4 in Wainwright and Jordan (2008)).



### 5.3.6 Representation of generic distributions

We now present a universal exponential family, which contains all positive with respect to a base measure distributions.

The formalism of exponential families can capture any probability distribution, when applying statistic functions of large expressivity. Taking for the statistic the identity function  $\delta [X_{[d]}, L_{[d]}]$  defined as

$$\delta [X_{[d]} = x_{[d]}, L_{[d]} = l_{[d]}] = \begin{cases} 1 & \text{if } x_{[d]} = l_{[d]} \\ 0 & \text{else} \end{cases},$$

we can represent any positive probability distribution  $\mathbb{P} [X_{[d]}]$  as a member of the exponential family  $\Gamma^{\delta, \mathbb{I}}$ . To see this, it is enough to choose

$$\theta [L_{[d]}] = \langle \ln [\mathbb{P} [X_{[d]}]] , \delta [X_{[d]} = x_{[d]}, L_{[d]} = l_{[d]}] \rangle [L_{[d]}],$$

where the contraction with  $\delta$  copies the variables  $X_{[d]}$  to  $L_{[d]}$ . The energy tensor of this member of  $\Gamma^{\delta, \mathbb{I}}$  is then

$$\langle \theta [L_{[d]}] , \delta [X_{[d]} = x_{[d]}, L_{[d]} = l_{[d]}] \rangle [X_{[d]}] = \ln [\mathbb{P} [X_{[d]}]]$$

and thus

$$\mathbb{P}^{\delta, \theta, \mathbb{I}} [X_{[d]}] = \langle \exp [\ln [\mathbb{P} [X_{[d]}]] \rangle [X_{[d]} | \emptyset] = \mathbb{P} [X_{[d]}].$$

We further note, that the mean parameter of this constructed element of  $\Gamma^{\delta, \mathbb{I}}$  is

$$\mu [L_{[d]}] = \langle \mathbb{P}^{\delta, \theta, \mathbb{I}} [X_{[d]}] , \delta [X_{[d]} = x_{[d]}, L_{[d]} = l_{[d]}] \rangle [L_{[d]}] = \langle \mathbb{P} [X_{[d]}] , \delta [X_{[d]} = x_{[d]}, L_{[d]} = l_{[d]}] \rangle [L_{[d]}],$$

and thus coincides with the distribution itself, after a relabelling of the distributed variables. Let us notice, that this family also correspond with the Markov Network on the maximal hypergraph  $\mathcal{G}^{\max} = (\mathcal{V}, \{\mathcal{V}\})$ . We will further revisit this family in Chapter 11, where we will refer to it by the minterm family in order to connect with terminology developed for logical reasoning in Chapter 8.

## 5.4 Polytopes of mean parameters

We in this section investigate properties of probability distributions based on their mean parameters. Given a statistic  $\mathcal{S}$ , we first define a mean parameter to any distribution by the expectation of the statistic.

**Definition 32.** Let there be a statistic  $\mathcal{S}$  and a boolean base measure  $\nu [X_{[d]}]$ . We call the tensor

$$\mu [L] = \langle \mathbb{P} [X_{[d]}] , \sigma^{\mathcal{S}} [X_{[d]}, L] \rangle [L]$$

the mean parameter tensor to a distribution  $\mathbb{P} [X_{[d]}]$ . The set

$$\mathcal{M}_{\mathcal{S}, \nu} = \{ \langle \mathbb{P} , \sigma^{\mathcal{S}} , \nu \rangle [L] : \mathbb{P} [X_{[d]}] \in \Gamma^{\delta, \nu} \},$$

is called the polytope of realizable mean parameters. Here we denote by  $\Gamma^{\delta, \nu}$  the set of all probability distributions representable with respect to  $\nu$  (see Def. 17).

While introduced here as a property of a distribution, the mean parameters will be central to probabilistic inference in Chapter 6. We in the reminder of this section prepare for this application and derive tensor network representations for distributions having sufficient statistics  $\mathcal{S}$ , depending on their corresponding mean parameter in the polytope.

### 5.4.1 Representation by convex hulls

First of all, we provide a simple characterization of the sets of mean parameters as the convex hull of the slices to the selection encoding of the statistic (see Figure 11). Convex hulls of finite vectors are called  $\mathcal{V}$ -polytopes (see Lecture 1 in Ziegler (2013)).

**Theorem 16.** For any statistic  $\mathcal{S}$  the polytope of mean parameters is the convex hull of the slices of  $\sigma^{\mathcal{S}}$  with fixed indices to  $X_{[d]}$ , that is

$$\mathcal{M}_{\mathcal{S}, \nu} = \text{conv} \left( \sigma^{\mathcal{S}} [X_{[d]} = x_{[d]}, L] : x_{[d]} \in \bigtimes_{k \in [d]} [m_k], \nu [X_{[d]} = x_{[d]}] = 1 \right).$$

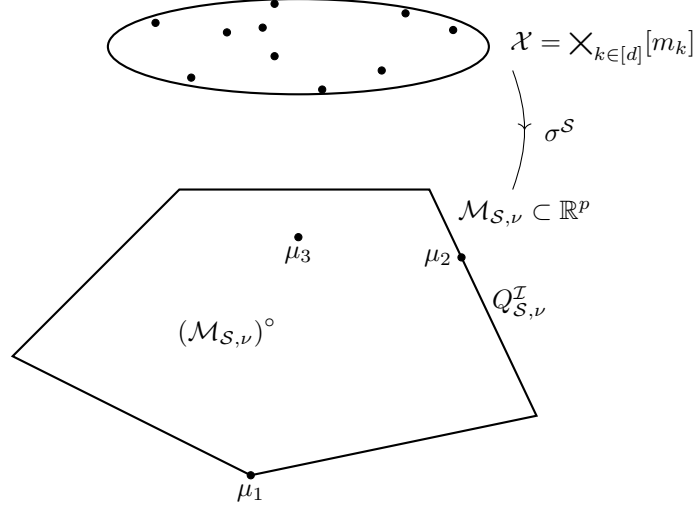


Figure 11: Sketch of the mean polytope  $\mathcal{M}_{S, \nu}$  to a statistic  $\mathcal{S}$ , which is minimal with respect to  $\nu$ . The mean polytope is a bounded subset of  $\mathbb{R}^p$  (here sketched as a 2-dimensional projection). For any vertex  $\mu_1$  we find  $x_{[d]}$  such that  $\mu_1[L] = \sigma^S[X_{[d]} = x_{[d]}, L]$ . Generic mean parameters  $\mu_2$  outside the interior are on a face  $Q_{S, \nu}^I$ . Interior points  $\mu_3 \in \mathcal{M}_{S, \nu}^\circ$  are exactly those reproducible by positive distributions with respect to  $\nu$ .

*Proof.* First we realize that the characterization of by  $\nu$  representable distributions is a standard simplex extended by trivial coordinates, that is

$$\Gamma^{\delta, \nu} = \text{conv}(\epsilon_{x_{[d]}}[X_{[d]}] : \nu[X_{[d]} = x_{[d]}] = 1) .$$

This follows from the fact, that the support of any by  $\nu$  representable distribution is contained in the support of  $\nu$ . Further, each representable distribution is contained in the convex hull of the one-hot encoded support elements, since any distribution is normed.

The polytope of mean parameters is a linear transform of the elements in  $\Gamma^{\delta, \nu}$ , since the contraction with  $\sigma^S$  is linear. It follows that

$$\begin{aligned} \mathcal{M}_{S, \nu} &= \text{conv}(\langle \sigma^S[X_{[d]}, L], \epsilon_{x_{[d]}}[X_{[d]}] \rangle [L] : \nu[X_{[d]} = x_{[d]}] = 1) \\ &= \text{conv}(\sigma^S[X_{[d]} = x_{[d]}, L] : \nu[X_{[d]} = x_{[d]}] = 1) . \end{aligned}$$

□

We thus understand the map

$$\sigma^S : \times_{k \in [d]} [m_k] \rightarrow \mathcal{M}_{S, \nu} \quad , \quad x \rightarrow \sigma^S[X_{[d]} = x_{[d]}, L]$$

as an encoding of states with respect to a statistic  $\mathcal{S}$ . When the statistic is the universal statistic  $\mathcal{S} = \delta$ , this statistic encoding coincides with the one-hot encoding and the mean polytope is the simplex of dimension  $\langle \nu \rangle[\emptyset] - 1$

$$\mathcal{M}_{\delta, \nu} = \text{conv}(\epsilon_{x_{[d]}}[X_{[d]}] : \nu[X_{[d]} = x_{[d]}] = 1) .$$

A generic mean polytope  $\mathcal{M}_{S, \nu}$  is then the linear transform of the simplex by the contraction with  $\sigma^S[X_{[d]}, L]$ , where  $L$  is left open.

#### 5.4.2 Representation as intersecting half-spaces

For any vector  $a[L] \in \mathbb{R}^p$  and a scalar  $b \in \mathbb{R}$ , we call the set

$$\{\mu[L] : \langle \mu[L], a[L] \rangle [\emptyset] \leq b\} \subset \mathbb{R}^p$$

a half-space of  $\mathbb{R}^p$ . Bounded intersections of finitely many half-spaces are called  $\mathcal{H}$ -polytopes Ziegler (2013). We state next, that the polytope  $\mathcal{M}_{S, \nu}$  of mean parameters is a  $\mathcal{H}$ -polytopes.

**Theorem 17.** *The set  $\mathcal{M}_{\mathcal{S},\nu}$  is for any statistic  $\mathcal{S}$  and base measure  $\nu$  a  $\mathcal{H}$ -polytope, i.e. there exists a finite collection*

$$((a_i[L], b_i) : i \in [n])$$

where  $a_i[L]$  a vector and  $b_i \in \mathbb{R}$  for all  $i \in [n]$  such that

$$\mathcal{M}_{\mathcal{S},\nu} = \{ \mu[L] : \forall_{i \in [n]} \langle \mu[L], a_i[L] \rangle [\emptyset] \leq b_i \} .$$

*Proof.* By The. 16, the set  $\mathcal{M}_{\mathcal{S},\nu}$  is the convex hull of a finite set of vectors and is therefore a  $\mathcal{V}$ -polytope. We therefore apply the main theorem for polytopes ?, which states the equivalence of  $\mathcal{V}$ -polytopes and  $\mathcal{H}$ -polytope, for which a proof can be found as Theorem 1.1 in Ziegler (2013). Therefore,  $\mathcal{M}_{\mathcal{S},\nu}$  is also a  $\mathcal{H}$ -polytope and has is thus the intersection of finitely many half-spaces.  $\square$

The determination of the half-space parametrizing  $((a_i[L], b_i) : i \in [n])$  is, however, in general difficult and the main reason for the intractability of probabilistic inference (see e.g. Wainwright and Jordan (2008)).

### 5.4.3 Characterization of the interior

The interior of the mean polytope consists of the mean parameters to positive distributions as we show next.

**Theorem 18.** *For any minimal statistics  $\mathcal{S}$  with respect to a boolean base measure  $\nu$  (see Def. 25) and a with respect to  $\nu$  positive distribution  $\mathbb{P}[X_{[d]}]$  we have*

$$\langle \mathbb{P}[X_{[d]}], \sigma^{\mathcal{S}}[X_{[d]}, L] \rangle [L] \in (\mathcal{M}_{\mathcal{S},\nu})^{\circ} .$$

*Proof.* Since by assumption the statistics is minimal, the convex set  $\mathcal{M}_{\mathcal{S},\nu}$  is full dimensional (see e.g. Appendix B in Wainwright and Jordan (2008)). We thus use a well-known property for full-dimensional convex sets (see Rockafellar (1997); Hiriart-Urruty and Lemarechal (1993)), that  $\mu \in \mathcal{M}_{\mathcal{S},\nu}^{\circ}$  if for any non-vanishing vector  $V[L]$  there is a there is a  $\tilde{\mu}[L]$  with

$$\langle V[L], \mu[L] \rangle [\emptyset] < \langle V[L], \tilde{\mu}[L] \rangle [\emptyset] .$$

It thus suffices to show for an arbitrary non-vanishing vector  $V[L]$  the existence of a distribution  $\tilde{\mathbb{P}}$ , such that

$$\langle V[L], \mu[L] \rangle [\emptyset] < \langle V[L], \sigma^{\mathcal{S}}[X_{[d]}, L], \tilde{\mathbb{P}}[X_{[d]}] \rangle [\emptyset] .$$

We define for  $\epsilon \in \mathbb{R}$

$$\mathbb{P}^{\epsilon}[X_{[d]}] = \langle \mathbb{P}[X_{[d]}], \exp[\epsilon \cdot \langle \sigma^{\mathcal{S}}[X_{[d]}, L], V[L] \rangle [X_{[d]}]] \rangle [X_{[d]} | \emptyset]$$

The derivation of this map at  $\epsilon = 0$  is

$$\frac{\partial}{\partial \epsilon} \mathbb{P}^{\epsilon}[X_{[d]}] |_{\epsilon=0} = \langle \mathbb{P}[X_{[d]}], \sigma^{\mathcal{S}}[X_{[d]}, L], V[L] \rangle [X_{[d]}] - \langle \mathbb{P}[X_{[d]}], \sigma^{\mathcal{S}}[X_{[d]}, L], V[L] \rangle [\emptyset] \cdot \mathbb{P}[X_{[d]}]$$

and thus

$$\begin{aligned} \frac{\partial}{\partial \epsilon} \langle \mathbb{P}^{\epsilon}[X_{[d]}], \sigma^{\mathcal{S}}[X_{[d]}, L], V[L] \rangle [\emptyset] |_{\epsilon=0} &= \left\langle \mathbb{P}[X_{[d]}], (\langle \sigma^{\mathcal{S}}[X_{[d]}, L], V[L] \rangle [\emptyset])^2 \right\rangle [X_{[d]}] \\ &\quad - (\langle \mathbb{P}[X_{[d]}], \sigma^{\mathcal{S}}[X_{[d]}, L], V[L] \rangle [X_{[d]}])^2 . \end{aligned}$$

We can interpret this quantity as the variance of the random variable  $\langle \sigma^{\mathcal{S}}[X_{[d]}, L], V[L] \rangle [X_{[d]} = x_{[d]}]$ , where  $x_{[d]}$  is drawn from  $\mathbb{P}[X_{[d]}]$ . The variance is greater than zero, if this random variable is not constant. But from the minimality of  $\mathcal{S}$  with respect to  $\nu$  it follows, that this variable is not constant and we therefore have

$$0 < \frac{\partial}{\partial \epsilon} \langle \mathbb{P}^{\epsilon}[X_{[d]}], \sigma^{\mathcal{S}}[X_{[d]}, L], V[L] \rangle [\emptyset] |_{\epsilon=0} .$$

Thus, there is a  $\epsilon > 0$  with

$$\langle V[L], \mu[L] \rangle [\emptyset] < \langle V[L], \sigma^{\mathcal{S}}[X_{[d]}, L], \mathbb{P}^{\epsilon}[X_{[d]}] \rangle [\emptyset] .$$

$\square$

While The. 18 only states that the mean parameter of each positive distribution is in the interior, we will construct for each interior point a positive distribution in Chapter 6 by a member of the corresponding exponential family.

#### 5.4.4 Characterization of the boundary by faces

Let us now continue with the investigation of the faces of the mean parameter polytope.

**Definition 33.** Given a mean parameter polytope  $\mathcal{M}_{\mathcal{S},\nu}$  in the half space representation of The. 17, and any subset  $\mathcal{I} \subset [n]$  we say that the set

$$Q_{\mathcal{S},\nu}^{\mathcal{I}} = \{\mu[L] \in \mathcal{M}_{\mathcal{S},\nu} : \forall_{i \in \mathcal{I}} \langle \mu[L], a_i[L] \rangle [\emptyset] = b_i\}$$

is the face to the constraints  $\mathcal{I}$ .

While all inequalities in a half-space representation are satisfied for any element of the polytope, we defined faces by the additional sharp satisfaction of a subset of the half-space inequalities. In this way, the faces build the boundary of  $\mathcal{M}_{\mathcal{S},\nu}$ . This can be easily verified, since for any vector  $\mu[L] \in \mathcal{M}_{\mathcal{S},\nu}$ , for which no halfspace inequalities hold sharply, also a neighborhood satisfies the halfspace inequalities. If any halfspace inequality holds sharply, in the other case, the vector is a member of the corresponding face.

If  $\mathcal{S}$  is not minimal with respect to  $\nu$ , we find a non-vanishing vector  $V[L]$  and a scalar  $\lambda \in \mathbb{R}$  such that

$$\langle \sigma^{\mathcal{S}}[X_{[d]}, L], V[L], \nu[X_{[d]}] \rangle [X_{[d]}] = \lambda \cdot \nu[X_{[d]}] .$$

This implies, that any probability distribution  $\mathbb{P}[X_{[d]}]$  representable with  $\nu$  satisfies

$$\langle \mathbb{P}[X_{[d]}], \sigma^{\mathcal{S}}[X_{[d]}, L], V[L], \nu[X_{[d]}] \rangle [\emptyset] = \lambda \cdot \langle \mathbb{P}[X_{[d]}], \nu[X_{[d]}] \rangle [\emptyset] = \lambda .$$

Any  $\mu[L] \in \mathcal{M}_{\mathcal{S},\nu}$  then satisfies

$$\langle \mu[L], V[L] \rangle [\emptyset] = \lambda .$$

Thus, the polytope  $\mathcal{M}_{\mathcal{S},\nu}$  is contained in an affine linear subspace and has vanishing interior. We can further understand this equation as two half-space inequalities

$$\langle \mu[L], V[L] \rangle [\emptyset] \leq \lambda \quad \text{and} \quad \langle \mu[L], V[L] \rangle [\emptyset] \geq \lambda ,$$

which can be integrated into any half-space representation. We conclude, that in the case of non-minimal statistics, the whole polytope  $\mathcal{M}_{\mathcal{S},\nu}$  is a face itself, since it satisfies these half-space inequalities sharply.

**Lemma 2.** For each face  $Q_{\mathcal{S},\nu}^{\mathcal{I}}$  we have

$$Q_{\mathcal{S},\nu}^{\mathcal{I}} = \text{conv}(\sigma^{\mathcal{S}}[X_{[d]} = x_{[d]}, L] : x_{[d]} \in (\sigma^{\mathcal{S}})^{-1}(Q_{\mathcal{S},\nu}^{\mathcal{I}}), \nu[X_{[d]} = x_{[d]}] = 1) .$$

*Proof.* This holds, since each face is the convex hull of the contained vertices (see Proposition 2.2 and 2.3 in Ziegler (2013)). Since the vertices are contained in the image of the statistic encoding  $\sigma^{\mathcal{S}}$ , the vertices contained in  $Q_{\mathcal{S},\nu}^{\mathcal{I}}$  are contained in the set

$$\sigma^{\mathcal{S}}[X_{[d]} = x_{[d]}, L] : x_{[d]} \in (\sigma^{\mathcal{S}})^{-1}(Q_{\mathcal{S},\nu}^{\mathcal{I}}) .$$

□

Lem. 2 implies in particular, that faces are mean parameter polytopes with respect to refined base measures. For reference in later chapters, we define these refined base measures next as face measures.

**Definition 34.** The base measure to the face  $Q_{\mathcal{S},\nu}^{\mathcal{I}}$  of  $\mathcal{M}_{\mathcal{S},\nu}$  is the boolean tensor

$$\nu^{\mathcal{S},\mathcal{I}}[X_{[d]}] = \mathbb{I}_{\gamma^{x_{[d]}} \in Q_{\mathcal{S},\nu}^{\mathcal{I}}}[X_{[d]}] .$$

**Theorem 19.** For any face  $Q_{\mathcal{S},\nu}^{\mathcal{I}}$  of  $\mathcal{M}_{\mathcal{S},\nu}$ , we have with the refined base measure

$$\tilde{\nu}[X_{[d]}] = \langle \nu[X_{[d]}], \nu^{\mathcal{S},\mathcal{I}}[X_{[d]}] \rangle [X_{[d]}]$$

that

$$Q_{\mathcal{S},\nu}^{\mathcal{I}} = \mathcal{M}_{\mathcal{S},\tilde{\nu}} .$$

*Proof.* We notice that for any  $x_{[d]} \in \times_{k \in [d]} [m_k]$ ,  $x_{[d]} \in (\sigma^S)^{-1}(Q_{S,\nu}^{\mathcal{I}})$  is equal to  $\nu^{S,\mathcal{I}} [X_{[d]} = x_{[d]}] = 1$  and thus

$$\{x_{[d]} : x_{[d]} \in (\sigma^S)^{-1}(Q_{S,\nu}^{\mathcal{I}}), \nu [X_{[d]} = x_{[d]}] = 1\} = \{x_{[d]} : \tilde{\nu} [X_{[d]} = x_{[d]}] = 1\}.$$

In combination with Lem. 2 we then get

$$\begin{aligned} Q_{S,\nu}^{\mathcal{I}} &= \text{conv} (\sigma^S [X_{[d]} = x_{[d]}, L] : x_{[d]} \in (\sigma^S)^{-1}(Q_{S,\nu}^{\mathcal{I}}), \nu [X_{[d]} = x_{[d]}] = 1) \\ &= \text{conv} (\sigma^S [X_{[d]} = x_{[d]}, L] : x_{[d]} : \tilde{\nu} [X_{[d]} = x_{[d]}] = 1) = \mathcal{M}_{S,\tilde{\nu}}. \end{aligned}$$

□

Positivity of a distribution with respect to face measures is an equivalent condition for the mean parameter of a distribution to be on a face, as we show next.

**Theorem 20.** *If and only if for a distribution  $\mathbb{P} [X_{[d]}]$  and a face  $\mathcal{I}$  we have*

$$\langle \mathbb{P} [X_{[d]}], \sigma^S [X_{[d]}, L] \rangle [L] \in Q_{S,\nu}^{\mathcal{I}},$$

*then  $\mathbb{P} [X_{[d]}]$  is representable with respect to the face measure  $\nu^{S,\mathcal{I}}$ .*

*Proof.* We have

$$\mu [L] = \sum_{x_{[d]}} \mathbb{P} [X_{[d]} = x_{[d]}] \cdot \gamma^S [X_{[d]} = x_{[d]}, L].$$

Now, the  $x_{[d]}$  with  $\nu^{S,\mathcal{I}} [X_{[d]} = x_{[d]}] = 1$  are exactly those, for which the conditions  $\mathcal{I}$  hold straight. If and only if for a  $x_{[d]}$  with  $\nu^{S,\mathcal{I}} [X_{[d]} = x_{[d]}] = 0$  we have  $\mathbb{P} [X_{[d]} = x_{[d]}] > 0$ , one of the conditions  $\mathcal{I}$  would not hold straight. Thus, if and only if  $\mathbb{P} [X_{[d]}]$  is representable with respect to  $\nu^{S,\mathcal{I}} [X_{[d]}]$ , we have  $\mu [L] \in Q_{S,\nu}^{\mathcal{I}}$ . □

Let us now investigate tensor network representations of face measures, based on the basis encoding  $\beta^S$  of a statistic. Vertices of  $\mathcal{M}_{S,\nu}$  are faces with single elements, that is  $\{\mu [L]\}$ . By Lem. 2 there must be  $\mu$  must lie in the image of  $\sigma^S$ , since otherwise  $\mathcal{M}_{S,\nu}$  would be empty. We denote  $y_{[p]}^{\mu}$  as the index such that  $I_S(y_{[p]}^{\mu}) = \mu$ . The vertex measure is then

$$\nu^{S,\mathcal{I}} [X_{[d]}] = \left\langle \beta^S [Y_{[p]}, X_{[d]}], \epsilon_{y_{[p]}^{\mu}} [Y_{[p]}] \right\rangle [X_{[d]}]$$

**Theorem 21** (Face measure representation). *For any face  $Q_{S,\nu}^{\mathcal{I}}$  of  $\mathcal{M}$  we have*

$$\nu^{S,\mathcal{I}} [X_{[d]}] = \left\langle \beta^S [Y_{[p]}, X_{[d]}], \alpha [Y_{[p]}] \right\rangle [X_{[d]}]$$

where

$$\alpha [Y_{[p]}] = \sum_{\mu \in Q_{S,\nu}^{\mathcal{I}} \cup \text{im}(\sigma^S)} \epsilon_{y_{[p]}^{\mu}} [Y_{[p]}].$$

*Proof.* For any  $\mu \in Q_{S,\nu}^{\mathcal{I}} \cup \text{im}(\sigma^S)$  the tensor

$$\tau^{\mu} [X_{[d]}] = \left\langle \beta^S [Y_{[p]}, X_{[d]}], \epsilon_{y_{[p]}^{\mu}} [Y_{[p]}] \right\rangle [X_{[d]}]$$

is the indicator of the preimage of  $\mu$  under  $\sigma^S$ . Since preimages of different  $\mu$  are disjoint, the support of  $\tau^{\mu} [X_{[d]}]$  is disjoint and their sum

$$\sum_{\mu \in Q_{S,\nu}^{\mathcal{I}} \cup \text{im}(\sigma^S)} \tau^{\mu} [X_{[d]}]$$

is the indicator of the preimage of  $Q_{S,\nu}^{\mathcal{I}}$  under  $\sigma^S$ , which is the face measure  $\nu^{S,\mathcal{I}} [X_{[d]}]$ . Exploiting linearity of contraction we have

$$\begin{aligned} \nu^{S,\mathcal{I}} [X_{[d]}] &= \sum_{\mu \in Q_{S,\nu}^{\mathcal{I}} \cup \text{im}(\sigma^S)} \tau^{\mu} [X_{[d]}] \\ &= \left\langle \beta^S [Y_{[p]}, X_{[d]}], \sum_{\mu \in Q_{S,\nu}^{\mathcal{I}} \cup \text{im}(\sigma^S)} \epsilon_{y_{[p]}^{\mu}} [Y_{[p]}] \right\rangle [X_{[d]}] \\ &= \left\langle \beta^S [Y_{[p]}, X_{[d]}], \alpha [Y_{[p]}] \right\rangle [X_{[d]}]. \end{aligned}$$

□

Let us notice, that  $\alpha [Y_{[p]}]$  in The. 21 is a sparse tensor with basis CP rank  $|Q_{S,\nu}^{\mathcal{I}} \cup \text{im}(\sigma^S)|$  (see Chapter 18).

## 5.5 Discussion and Outlook

This chapter has established a foundational treatment of probability distributions by tensors, and motivated tensor network decompositions along classical approaches towards graphical models. To show this correspondence, we defined both tensor networks and graphical models based on the same hypergraph. This then enabled us to define Markov Networks simply as the normalizations of tensor networks with non-negative coordinates. In the literature, tensor networks are, however, often treated as being dual to the graphs defining graphical models (see e.g. Robeva and Seigal (2019)). The duality becomes clear, when one interpretes the tensors as nodes and their common variables as edges, as might be natural given the applied notation of wiring diagrams to represent tensor networks. We in this work avoid the discussion of this ambiguity, and treat tensors as decorations of hyperedges.

In the literature, the tensors decorating hyperedges are often referred to as "factors" and their coordinatewise logarithm as "features" Koller and Friedman (2009). With the scope of this work, we avoided such further terminology.

Further, graphical models follow a tradition of definition on graphs, instead of hypergraphs. Tensors, or "factors", are then assigned to maximal cliques. We observed that the notion of maximality is an important assumption, for example in the proof of the Hammersley-Clifford theorem, and therefore introduced the property of clique capturing hypergraphs (see Def. 29) to connect with this graph-based formalism.

While we here restricted our discussion to finite state spaces to each variable, probability distributions can in general be defined for arbitrary measurable spaces. Joint distributions of these more generic variables still have a tensor structure. The discussion of them, however, needs to be more careful, since integrals might diverge and tensors therefore not be normable. By restriction in this work to finite state spaces of factored systems, we were able to exclude such situations.

## 6 Probabilistic Inference

After having investigated sparse decomposition schemes of probability distributions into tensor networks, we now exploit these schemes to derive efficient reasoning schemes. We first introduce by queries a generic scheme to retrieve information by contractions, and introduce the method of maximum likelihood estimation related to entropy optimization. Then we focus on inference tasks in exponential families, which have been introduced as a generalization of graphical models in the previous chapter.

### 6.1 Queries

The efficient retrieval of information stored in probability distributions has to exploit the available decomposition schemes. To avoid the instantiation of a distribution based on its decomposition, we directly define deductive reasoning schemes by contractions, which can be executed using the available decomposition.

#### 6.1.1 Querying by functions

We now formalize queries by retrieving expectations of functions given a distribution specified by probability tensors. We exploit basis calculus in defining categorical variables  $Y_q$  to tensors  $q$ , which are enumerating the set  $\text{im}(q)$ . More details on this scheme are provided in Chapter 17, see Def. 84 therein.

**Definition 35.** *The marginal query of a probability distribution  $\mathbb{P}[X_{[d]}]$  by a query statistic*

$$q : \bigtimes_{k \in [d]} [m_k] \rightarrow \mathcal{U}^q$$

*is the vector  $\mathbb{P}[Y_q]$ , where  $Y_q$  is an image enumerating variable (see Chapter 17 for more details), defined as the contraction*

$$\mathbb{P}[Y_q] = \langle \mathbb{P}[X_{[d]}], \beta^q[Y_q, X_{[d]}] \rangle [Y_q] .$$

*If  $\mathcal{U}^q \subset \mathbb{R}$ , and the statistic  $q$  is therefore a tensor, we further define the expectation query of  $\mathbb{P}[X_{[d]}]$  by  $q$  as*

$$\mathbb{E}[q] = \langle q(X_{[d]}), \mathbb{P}[X_{[d]}] \rangle [\emptyset] .$$

*Given another query statistic  $g : \bigtimes_{k \in [d]} [m_k] \rightarrow \mathcal{U}^g$ , which image is enumerated by a variable  $Y_g$ , the conditional query of the probability distribution  $\mathbb{P}[X_{[d]}]$  by the statistic  $q$  conditioned on  $g$  is the matrix  $\mathbb{P}[Y_q|Y_g] \in \mathbb{R}^{|\text{im}(q)|} \otimes \mathbb{R}^{|\text{im}(g)|}$  defined as the normalization*

$$\mathbb{P}[Y_q|Y_g] = \langle \mathbb{P}[X_{[d]}], \beta^q[Y_q, X_{[d]}], \beta^g[Y_g, X_{[d]}] \rangle [Y_q|Y_g] .$$

We notice, that marginal and conditional queries generalize the schemes of marginalization and conditioning on the distributed variables. As such, marginal distributions are marginal queries with respect to a identity query statistic acting on the respective variables. Conditional distributions can be similarly retrieved by two identity statistics, representing the incoming and outgoing variables.

Expectation queries return the expectation of a real-valued feature  $q : \times_{k \in [d]} [m_k] \rightarrow \mathbb{R}$ . When understanding  $q$  as a random variable given the probability measure  $\mathbb{P} [X_{[d]}]$ , the expectation query returns the expectation of that random variable. Expectation queries are further contractions of marginal queries with the identity function restricted on the image of  $q$ , since

$$\mathbb{E} [q] = \langle \mathbb{P} [Y_q], \text{Id}_{|\text{im}(q)} Y_q \rangle [\emptyset] .$$

This contraction equations follows from the more general Cor. 13, which will be shown in Chapter 17. **Expectation queries on statistics build the mean parameter of a distribution.**

### 6.1.2 Mode Queries

A different kind of queries are mode queries, which we formalize by the searches of the indices to maximal coordinates in a tensor.

**Definition 36.** *Given a tensor  $\tau [X_{[d]}]$  the mode query is the problem*

$$\text{argmax}_{x_{[d]} \in \times_{k \in [d]} [m_k]} \tau [X_{[d]} = x_{[d]}] . \quad (8)$$

We can pose mode queries as convex optimization problems. To this end we recall, that the set of all probability distributions is a convex hull of the one-hot encoded states, that is

$$\Gamma^{\delta, \mathbb{I}} = \text{conv} \left( \epsilon_{x_{[d]}} [X_{[d]}] : x_{[d]} \in \times_{k \in [d]} [m_k] \right) .$$

With this we have

$$\begin{aligned} \max_{x_{[d]}} \tau [X_{[d]} = x_{[d]}] &= \max_{x_{[d]}} \langle \tau [X_{[d]}], \epsilon_{x_{[d]}} [X_{[d]}] \rangle [\emptyset] \\ &= \max_{\mu [X_{[d]}] \in \Gamma^{\delta, \mathbb{I}}} \langle \tau [X_{[d]}], \mu [X_{[d]}] \rangle [\emptyset] . \end{aligned}$$

We note that the maximization over  $\Gamma^{\delta, \mathbb{I}}$  is a convex optimization problem and the maxima are taken at

$$\text{argmax}_{\mu [X_{[d]}] \in \Gamma^{\delta, \mathbb{I}}} \langle \tau [X_{[d]}], \mu [X_{[d]}] \rangle [\emptyset] = \text{conv} \left( \epsilon_{x_{[d]}} [X_{[d]}] : x_{[d]} \in \text{argmax}_{x_{[d]}} \tau [X_{[d]} = x_{[d]}] \right) .$$

We will further apply this generic trick to approach mode queries when studying inference problems for exponential families in the following sections.

Let us note, that we can also approach modified mode queries, where we restrict to  $\mathcal{U} \subset \times_{k \in [d]} [m_k]$ , namely

$$\text{argmax}_{x_{[d]} \in \mathcal{U}} \tau [X_{[d]} = x_{[d]}] .$$

We can choose the base measure by the subset encoding (see Chapter 17) of  $\mathcal{U}$ , namely

$$\nu [X_{[d]}] = \sum_{x_{[d]} \in \mathcal{U}} \epsilon_{x_{[d]}} [X_{[d]}] ,$$

and conclude that

$$\text{argmax}_{\mu [X_{[d]}] \in \Gamma^{\delta, \nu}} \langle \tau [X_{[d]}], \mu [X_{[d]}] \rangle [\emptyset] = \text{conv} \left( \epsilon_{x_{[d]}} [X_{[d]}] : x_{[d]} \in \text{argmax}_{x_{[d]} \in \mathcal{U}} \tau [X_{[d]} = x_{[d]}] \right) .$$

### 6.1.3 Energy representations

For exponential families (see Sect. 5.2) we have observed, that often energy tensors have feasible tensor network representations, whereas the corresponding probability distributions can get infeasible. We therefore investigate here schemes to answer queries based on the energy tensor instead of the distribution.

**Theorem 22.** Let  $\phi[X_{[d]}]$  be an energy tensor and  $\mathbb{P}[X_{[d]}] = \langle \exp[\phi[X_{[d]}]] \rangle [X_{[d]}|\emptyset]$  the corresponding distribution. For disjoint subsets  $A, B \subset [d]$  with  $A \cup B = [d]$  and any  $x_B$  we have

$$\mathbb{P}[X_A|X_B = x_B] = \langle \exp[\phi[X_A, X_B = x_B]] \rangle [X_A|\emptyset] .$$

*Proof.* To show the theorem, we use a generic simplification property of coordinatewise transforms, which we will show as Lem. 26 in Chapter 16 and get

$$\langle \exp[\phi[X_A, X_B]] \rangle [X_A, X_B = x_B] = \langle \exp[\phi[X_A, X_B = x_B]] \rangle [X_A]$$

Based on this we get

$$\begin{aligned} \mathbb{P}[X_A|X_B = x_B] &= \langle \exp[\phi[X_A, X_B]] \rangle [X_A|X_B = x_B] \\ &= \frac{\langle \exp[\phi[X_A, X_B]] \rangle [X_A, X_B = x_B]}{\langle \exp[\phi[X_A, X_B]] \rangle [X_B = x_B]} \\ &= \frac{\langle \exp[\phi[X_A, X_B = x_B]] \rangle [X_A]}{\langle \exp[\phi[X_A, X_B = x_B]] \rangle [\emptyset]} \\ &= \langle \exp[\phi[X_A, X_B = x_B]] \rangle [X_A|\emptyset] . \end{aligned}$$

□

Importantly, The. 22 does not generalize to situations, where  $A \cup B \neq [d]$ , since summation over the indices of the variables  $[d]/A \cup B$  and contraction do not commute. In this more generic situation, we would need to sum over exponentiated coordinates, that is

$$\mathbb{P}[X_A|X_B = x_B] = \left\langle \sum_{x_{[d]/A \cup B} \in [m_{[d]/A \cup B}]} \exp[\phi[X_A, X_B = x_B, X_{[d]/A \cup B} = x_{[d]/A \cup B}]] \right\rangle [X_A|\emptyset] .$$

Mode queries on probability distributions in an energy representation can always be reduced to mode queries on the energy tensor. This is due to the monotonicity of the exponential function, which implies

$$\begin{aligned} \operatorname{argmax}_{x_{[d]} \in \times_{k \in [d]} [m_k]} \mathbb{P}[X_{[d]} = x_{[d]}] &= \operatorname{argmax}_{x_{[d]} \in \times_{k \in [d]} [m_k]} \langle \exp[\phi[X_{[d]}]] \rangle [X_{[d]} = x_{[d]}|\emptyset] \\ &= \operatorname{argmax}_{x_{[d]} \in \times_{k \in [d]} [m_k]} \exp[\phi[X_{[d]} = x_{[d]}]] \\ &= \operatorname{argmax}_{x_{[d]} \in \times_{k \in [d]} [m_k]} \phi[X_{[d]} = x_{[d]}] . \end{aligned}$$

Since we are only interested in identifying the index of the maximum coordinate, and not its value, we have further dropped the normalization term by partition functions. When one instead need to get the value of the maximal, the partition function cannot be ignored.

## 6.2 Sampling

Let us here investigate how to draw samples from a probability distribution, based on queries on it. Naive methods, such as drawing a random number in  $[0, 1]$ , adding iteratively the coordinates and stopping when the sum exceeds the random variables, are infeasible when having large tensor orders causing exponential increases of the coordinate number. We recall, that the number of coordinates of  $\mathbb{P}[X_{[d]}]$  is  $\prod_{k \in [d]} m_k$ , which increases exponentially in the number  $d$  of the variables. Efficient methods instead have to exploit tensor network decompositions of the decompositions.

### 6.2.1 Exact Methods

The first insight to derive efficient sampling algorithms is to sample a single variable in each step. Forward sampling (see Algorithm 1) exploits to this end the generic chain decomposition (see The. 5) of a probability distribution, namely

$$\mathbb{P}[X_{[d]}] = \langle \{\mathbb{P}[X_0]\} \cup \{\mathbb{P}[X_k|X_0, \dots, X_{k-1}] : k \in [d], k \geq 1\} \rangle [X_{[d]}] ,$$

It then samples iteratively a state  $x_k$  for the variable  $X_k$  conditioned on the previously sampled states, that is from the conditional distribution

$$\mathbb{P}[X_k|X_{[k]} = x_{[k]}] .$$



**Algorithm 1** Forward Sampling**Require:** Probability distribution  $\mathbb{P}$ **Ensure:** Exact sample  $x_0, \dots, x_{d-1}$  of  $\mathbb{P}$ **for**  $k \in [d]$  **do**    Draw  $x_k \in [m_k]$  from the conditional distribution

$$\mathbb{P} [X_k | X_{[k]} = x_{[k]}]$$

**end for**

The generic chain decomposition thereby ensures that probability of getting a state  $x_{[d]}$  by this procedure coincides with  $\mathbb{P} [X_{[d]} = x_{[d]}]$ .

Forward Sampling is especially efficient, when sampling from a Bayesian Network respecting the topological order of its nodes. More technically, when the parents  $\text{Pa}(k)$  of a node  $k$  are contained in the preceding variables  $[k]$ , we apply the conditional independence assumption (more precisely The. 8 in combination with The. 13) to get

$$\mathbb{P} [X_k | X_{[k]} = x_{[k]}] = \mathbb{P} [X_k | X_{\text{Pa}(k)} = x_{\text{Pa}(k)}] .$$

Since this conditional probability coincides with a local tensor in the Bayesian Network, we can avoid to contract the network for preparing the conditional distribution. Different to more general Markov Networks, forward sampling from Bayesian Network can therefore be done efficiently by reduction to conditional queries answerable using local tensors. We note that it is important to sample in the topological order induced by the underlying directed hypergraph, since the computation of generic conditional distributions is also for Bayesian Networks NP-hard (see Chapter 13 in Koller and Friedman (2009)). Sampling along the topological variable order requires only tractable to answer conditional queries on the Bayesian Network.

**6.2.2 Gibbs Sampling**

While we have seen that forward sampling can be performed efficiently on Bayesian Networks, Gibbs sampling can be also performed efficiently for Markov Networks. Gibbs sampling Algorithm 2 overcomes the intractability problems of sampling steps during forward sampling at the expense of repetitions of the sampling step. When performing finite repetitions, Gibbs sampling in general samples from an approximate distribution to the one desired. It can be shown, that these approximate distribution tend to one desired in the asymptotic limit of infinite repetitions of the sampling step (see Chapter 12 in Koller and Friedman (2009)).

**Algorithm 2** Gibbs Sampling**Require:** Probability distribution  $\mathbb{P}$ **Ensure:** Approximative sample  $x_0, \dots, x_{d-1}$  of  $\mathbb{P}$ **for**  $k \in [d]$  **do**    Draw  $x_k$  from an initialization distribution.**end for****while** Stopping criterion is not met **do**    **for**  $k \in [d]$  **do**        Draw  $x_k$  from the conditional distribution

$$\mathbb{P} [X_k | X_{[d]/\{k\}} = x_{[d]/\{k\}}]$$

**end for****end while return**  $x_0, \dots, x_{d-1}$ 

The central problem of forward sampling on Markov Networks has been the need for global contractions to answer the required conditional queries, which originates from large numbers of variables to be marginalized out. When avoiding the marginalization of variables, and conditioning on them instead, global contractions can be avoided. To be more precise, for any tensor network  $\tau^G$  on  $\mathcal{G} = ([d], \mathcal{E})$  and any  $k \in [d]$  we have

$$\langle \tau^G \rangle [X_k, X_{[d]/\{k\}}] = \langle \{ \tau [e] X_k, X_{e/\{k\}} = x_{e/\{k\}} : e \in \mathcal{E}, k \in e \} \rangle [X_k] \cdot \prod_{e \in \mathcal{E}, k \notin e} \tau [e] X_e = x_e .$$

As a consequence, we get for the Markov Network  $\mathbb{P} = \mathbb{P}^{\mathcal{G}}$  to  $\tau^{\mathcal{G}}$ , that

$$\begin{aligned}\mathbb{P}[X_k | X_{[d]/\{k\}} = x_{[d]/\{k\}}] &= \langle \tau^{\mathcal{G}} \rangle [X_k, X_{[d]/\{k\}} | \emptyset] \\ &= \langle \{ \tau[e] X_k, X_{e/\{k\}} = x_{e/\{k\}} : e \in \mathcal{E}, k \in e \} \rangle [X_k | \emptyset] .\end{aligned}$$

The conditional queries on a Markov Network asked in Gibbs Sampling can therefore be answered by contractions only of those tensors containing the variable  $X_k$ . To find further locally answerable conditional queries, we need to condition only on the neighbored variables, referred to as Markov blanket, such that the other variables are conditionally independent. This follows from the characterization of the conditional independences eminent in Markov Networks, which has been shown in The. 12, and can be used to design further tractable sampling schemes for Markov Networks.

We can further answer these conditional queries efficiently, when we perform Gibbs sampling on a probability distribution in an energy representation, that is  $\mathbb{P}[X_{[d]}] = \langle \exp[\phi[X_{[d]}]] \rangle [X_{[d]} | \emptyset]$ . Using The. 22, we have

$$\mathbb{P}[X_k | X_{[d]/\{k\}} = x_{[d]/\{k\}}] = \langle \exp[\phi[X_k, X_{[d]/\{k\}} = x_{[d]/\{k\}}]] \rangle [X_k | \emptyset]$$

We note, that the main property of the conditional query exploited here, is that all variables but the one sampled appear as a condition and none is marginalized out. In the scheme of forward sampling, where most of the variables are marginalized out in many queries, we cannot apply this trick and would have to perform sums over exponentiated coordinates to the variables marginalized out.

### 6.2.3 Simulated Annealing

Simulated annealing is an adapted sampling scheme that targets mode queries rather than generating representative samples from a distribution. It employs an annealing process that gradually transforms the probability distribution by increasingly favoring high-likelihood configurations, thereby improving the chances of sampling a solution to a mode query. To be more precise, let there be a distribution in energy representation, that is  $\mathbb{P}[X_{[d]}] = \langle \exp[\phi[X_{[d]}]] \rangle [X_{[d]} | \emptyset]$ . We introduce a parameter  $\beta \in \mathbb{R}$ , which we understand as the inverse temperature, and anneal the distribution through scaling its energy by this parameter. In the limit of  $\beta \rightarrow \infty$ , for each state  $x_{[d]} \in \times_{k \in [d]} [m_k]$  the annealed distribution behaves as

$$\langle \exp[\beta \cdot \phi[X_{[d]}]] \rangle [X_{[d]} = x_{[d]} | \emptyset] \rightarrow \left\langle \mathbb{I}_{\text{argmax}_{x_{[d]} \in \times_{k \in [d]} [m_k]} \phi[X_{[d]} = x_{[d]}]} \right\rangle [X_{[d]} = x_{[d]} | \emptyset] .$$

In this limit, the annealed distribution tends to the uniform distribution of the maximal coordinates, that is the uniform distribution of the set

$$\text{argmax}_{x_{[d]} \in \times_{k \in [d]} [m_k]} \phi[X_{[d]} = x_{[d]}] = \text{argmax}_{x_{[d]} \in \times_{k \in [d]} [m_k]} \mathbb{P}[X_{[d]} = x_{[d]}] .$$

i To integrate annealing into Gibbs sampling, one chooses a parameter  $\beta$  for each repetition of a sampling step and sample from the conditioned annealed distribution  $\langle \exp[\beta \cdot \phi[X_{[d]}]] \rangle [X_{[d]} | \emptyset]$ . Through increasing  $\beta$  during the algorithm, the samples are drawn towards states with larger coordinates in  $\mathbb{P}[X_{[d]}]$ . However, when  $\beta$  is large, the sampling procedure can get stuck in local maxima, whereas small  $\beta$  are in favor of overcoming such. The inverse temperature is thus understood as a tradeoff parameter between the exploration of new regions of the state space and increasing the coordinate of the sample by local coordinate optimization. It is therefore typically chosen low in the beginning of the sampling algorithm and then sequentially increased to find maximal coordinates. Due to this typical increase of the inverse temperature strategy, the algorithm is referred to as simulated annealing.

## 6.3 Maximum Likelihood Estimation

So far we have been concerned with deductive reasoning task, that is retrieve information from a given distribution or drawing a sample. We now turn to inductive reasoning tasks, where a probability distribution is estimated given data. To present the generic framework of maximum likelihood estimation in the tensor network contraction formalism, we introduce the likelihood loss exploiting the structure of empirical distribution, and then provide interpretations in terms of entropies.

### 6.3.1 Empirical Distributions

To prepare for reasoning on data, we now derive tensor network representation for empirical distributions, which are defined based on observed states  $((x_0^j, \dots, x_{d-1}^j) : j \in [m])$  of a factored system.

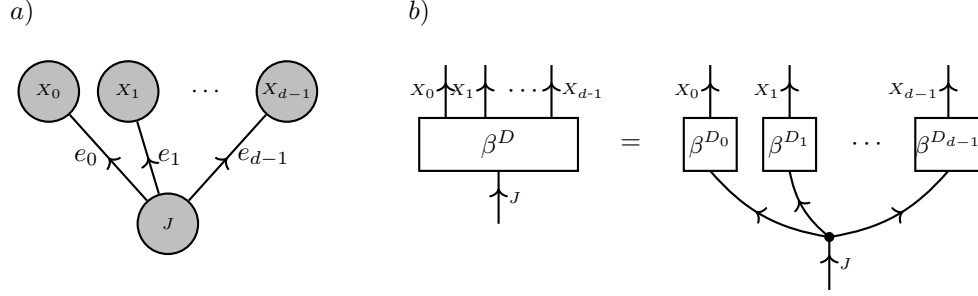


Figure 12: Decomposition of the basis encoding of a sample selector map to a dataset  $((x_0^j, \dots, x_{d-1}^j) : j \in [m])$ . a) Interpretation as a sample selection variable  $J$  selecting states for the variables  $X_{[d]}$  according to the enumerated dataset. b) Corresponding decomposition of the basis encoding  $\beta^D$  into a tensor network in the basis CP Format (see Sect. 18.1.2), where  $\tau^{e_k} = \beta^{D_k}$ .

**Definition 37.** Given a dataset  $((x_0^j, \dots, x_{d-1}^j) : j \in [m])$  of samples of the factored system we define the sample selector map

$$D : [m] \rightarrow \bigtimes_{k \in [d]} [m_k]$$

elementwise by

$$D(j) = (x_0^j, \dots, x_{d-1}^j).$$

The empirical distribution to the sample selector map  $D$  is the probability distribution

$$\mathbb{P}^D [X_{[d]}] := \langle \beta^D [X_{[d]}, J] \rangle [X_{[d]} | \emptyset],$$

where we introduced as single incoming for the basis encoding of the sample selector map the sample selecting variable  $J$  taking values in  $[m]$ .

The basis encoding of the sample selector map has a decomposition by

$$\beta^D [X_{[d]}, J] = \sum_{j \in [m]} \epsilon_{x_0^j, \dots, x_{d-1}^j} [X_{[d]}] \otimes \epsilon_j [J].$$

Each coordinate  $x_{[d]}$  of the empirical distribution can thus be calculated by

$$\begin{aligned} \mathbb{P}^D [X_{[d]} = x_{[d]}] &= \frac{1}{\langle \beta^D \rangle [\emptyset]} \left( \sum_{j \in [m]} \epsilon_{x_0^j, \dots, x_{d-1}^j} [X_{[d]} = x_{[d]}] \right) \\ &= \frac{\left| \{ j \in [m] : (x_0^j, \dots, x_{d-1}^j) = (x_0, \dots, x_{d-1}) \} \right|}{|j \in [m]|}. \end{aligned}$$

We can therefore interpret each coordinate of the empirical distribution as the relative frequency of the corresponding state in the observed data.

The basis encoding of the sample selector map is a sum of one-hot encodings of the data indices and the corresponding sample states. Such sums of basis tensors will be further investigated in Sect. 18.1.2 as basis CP decompositions. We now exploit this structure to find efficient tensor network decompositions (see Figure 12) based on matrices encoding its variables.

**Theorem 23.** Given a data map  $D : [m] \rightarrow \bigtimes_{k \in [d]} [m_k]$  we define for  $k \in [d]$  its coordinate maps

$$D_k : [m] \rightarrow [m_k]$$

by

$$D_k(j) = x_k^j.$$

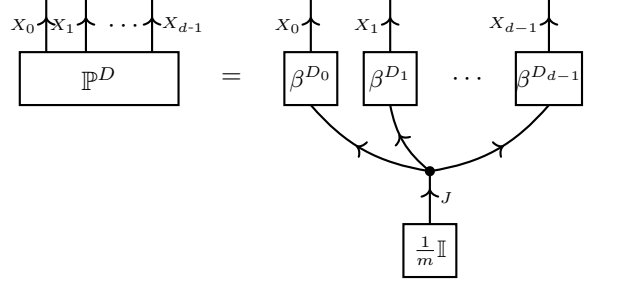
We then have

$$\beta^D [X_{[d]}, J] = \left\langle \{\beta^{D^k} [X_k, J] : k \in [d]\} \right\rangle [X_{[d]}, J]$$

and

$$\mathbb{P}^D [X_{[d]}] = \left\langle \beta^D [X_{[d]}, J], \frac{1}{m} \mathbb{I} [J] \right\rangle [X_{[d]}] = \left\langle \beta^{D_0} [X_0, J], \dots, \beta^{D_{d-1}} [X_{d-1}, J], \frac{1}{m} \mathbb{I} [J] \right\rangle [X_{[d]}] .$$

In a contraction diagram this decomposition is represented as



*Proof.* The first claim is a special case of Theorem 126, to be shown in Chapter 17. To show the second claim we notice

$$\langle \beta^D \rangle [\emptyset] = \sum_{j \in [m]} \langle \beta^D [X_{[d]}, J = j] \rangle [\emptyset] = m .$$

With the first claim it now follows that

$$\mathbb{P}^D [X_{[d]}] = \langle \beta^D \rangle [X_{[d]} | \emptyset] = \frac{\langle \beta^D \rangle [X_{[d]}]}{\langle \beta^D \rangle [\emptyset]} = \left\langle \{\beta^{D^k} [X_k, J] : k \in [d]\} \cup \left\{ \frac{1}{m} \mathbb{I} [J] \right\} \right\rangle [X_{[d]}] .$$

□

The cores  $\beta^{D^k}$  are matrices storing the value of the categorical variable  $X_k$  in the sample world indexed by  $j$ .

From the proof of Theorem 23 we notice that the scalar  $\frac{1}{m}$  could be assigned with any core in a representation of  $\mathbb{P}^D$ , and the core  $\mathbb{I} [J]$  is thus redundant in the contraction representation. However, creating the core  $\frac{1}{m} \mathbb{I} [J]$  provides us with a simple interpretation of the empirical distribution. We can understand  $\frac{1}{m} \mathbb{I} [J]$  as the uniform probability distribution over the samples, which is by the map  $D$  forwarded to a distribution over  $\times_{k \in [d]} [m_k]$ . The one-hot encoding of each sample is itself a probability distribution, which is understood as conditioned on the respective state of the sample selection variable  $J$ . The conditional distribution  $\beta^D$  therefore forwards the uniform distribution of the samples to a distribution of the variables  $X_{[d]}$ . In the perspective of a Bayesian Network (see Figure 12), the variable  $J$  serves as single parent for each categorical variable  $X_k$ .

### 6.3.2 Likelihood Loss

The likelihood of a probability distribution  $\mathbb{P} [X_{[d]}]$  to produce an observed sample is

$$\mathbb{P} [X_{[d]} = D(j)] .$$

We further introduce the likelihood of  $\mathbb{P} [X_{[d]}]$  with respect to a dataset as

$$\mathbb{P} \left[ \left( (x_0^j, \dots, x_{d-1}^j) : j \in [m] \right) \right] := \prod_{j \in [m]} \mathbb{P} [X_{[d]} = D(j)] .$$

The likelihood draws on the assumption, that each datapoint in the dataset has been drawn independently from the same distribution. When this generating distribution coincides with  $\mathbb{P} [X_{[d]}]$ , then the probability of generating a dataset by this scheme is the likelihood. In inductive reasoning, the true distribution  $\mathbb{P} [X_{[d]}]$  is unknown and needs to be approximatively estimated based on data and a learning hypothesis. We will therefore compute and compare the likelihood of distributions, which in general do not coincide with the distribution generating the data. It is thus

important to not understand the likelihood as a probability, which is only true for the generating distribution, as pointed out in Chapter 2 in MacKay (2003).

Let us now transform the likelihood to find an representation involving the empirical distribution (see Def. 37), for which efficient tensor network decompositions have been derived in The. 23. Applying a scaled logarithm we get

$$\begin{aligned} \frac{1}{m} \cdot \ln \left[ \mathbb{P} \left[ \left( (x_0^j, \dots, x_{d-1}^j) : j \in [m] \right) \right] \right] &= \frac{1}{m} \cdot \ln \left[ \prod_{j \in [m]} \mathbb{P} [X_{[d]} = D(j)] \right] \\ &= \frac{1}{m} \sum_{j \in [m]} \ln [\mathbb{P} [X_{[d]} = D(j)]] \\ &= \langle \ln [\mathbb{P} [X_{[d]}]] , \mathbb{P}^D [X_{[d]}] \rangle [ ] . \end{aligned}$$

Let us notice, that this transform of the likelihood is monotoneous and therefore does not influence the position of the maximum, when optimizing the likelihood. Motivated by this property, we use the transformed form to define the loss-likelihood loss and maximum likelihood estimation.

**Definition 38.** *The log-likelihood loss of a distribution  $\mathbb{P}$  given a dataset  $\left( (x_0^j, \dots, x_{d-1}^j) : j \in [m] \right)$  is the functional*

$$\mathcal{L}_D (\mathbb{P}) = - \langle \ln [\mathbb{P} [X_{[d]}]] , \mathbb{P}^D [X_{[d]}] \rangle [ ] .$$

Having a hypothesis  $\Gamma \subset \Gamma^{\delta, \mathbb{I}}$ , that is a set of probability distributions, the maximum likelihood estimation is the problem

$$\operatorname{argmin}_{\mathbb{P} \in \Gamma} \mathcal{L}_D (\mathbb{P}) . \quad (\mathbf{P}_{\Gamma, \mathbb{P}^D}^M)$$

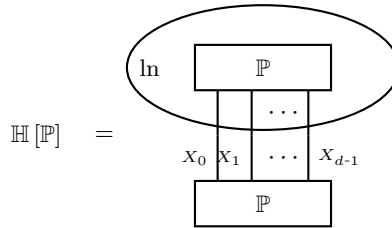
### 6.3.3 Entropic Interpretation

The Shannon entropy, which has been introduced in the seminal paper Shannon (1948), is a foundational concept in various research fields beyond statistical learning, such as information theory or statistical physics. While a detailed discussion is out of the scope of this work, we here only provide computation schemes of the entropy based on contractions of distributions.

**Definition 39** (Shannon entropy). *The Shannon entropy of a distribution  $\mathbb{P} [X_{[d]}]$  is the quantity*

$$\mathbb{H} [\mathbb{P}] = \sum_{x_{[d]} \in \times_{k \in [d]} [m_k]} \mathbb{P} [X_{[d]} = x_{[d]}] \cdot (-\ln [\mathbb{P} [X_{[d]} = x_{[d]}]]) = \langle \mathbb{P}, -\ln [\mathbb{P}] \rangle [\emptyset] .$$

We represent the Shannon entropy by the tensor network diagram



where we denote a coordinatewise transform by the logarithm as an ellipsis (see Sect. 16.2).

We here make the convention  $\ln [0] = -\infty$  and  $0 \cdot \ln [0] = 0$ , to have the Shannon entropy well-defined for distributions with non-trivial support.

Among the distributions in the same tensor space, the uniform distribution maximizes the Shannon entropy

$$\mathbb{H} [\langle \mathbb{I} \rangle [X_{[d]} | \emptyset]] = \sum_{k \in [d]} \ln [m_k]$$

and the one-hot encodings to states  $x_{[d]} \in \times_{k \in [d]} [m_k]$  minimize the Shannon entropy

$$\mathbb{H} [\epsilon_{x_{[d]}} [X_{[d]}]] = 0 .$$

The Shannon entropy measures the information content of a distribution and is therefore a central tool for regularization in statistical learning (see for an introduction Chapter 2 in MacKay (2003)). We therefore exploit this information content as a regularizer to identify a distribution among those coinciding in the answer to a collection of expectation queries. To be more precise, let there be for  $l \in [p]$  query tensors  $q^l$  (see Def. 35) and  $\mathbb{P}^D$  an empirical distribution. The problem of maximal entropy with respect to coinciding expectation queries with  $\mathbb{P}^D$  is then posed as

$$\operatorname{argmax}_{\mathbb{P} \in \Gamma^{\delta, \mathbb{I}}} \mathbb{H}[\mathbb{P}] \quad \text{subject to} \quad \forall l \in [p] : \langle \mathbb{P}, q^l \rangle [\emptyset] = \langle \mathbb{P}^D, q^l \rangle [\emptyset]$$

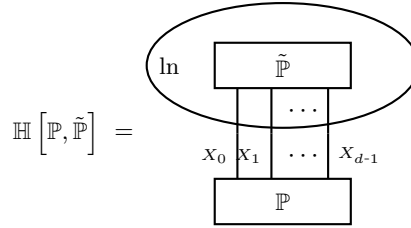
where  $\Gamma^{\delta, \mathbb{I}}$  is the set of probability distributions given a factored system. We study instances of this maximal entropy problem later in Sect. 6.5.1, where we show that its solution is a member of the exponential family, which statistic is build by the query tensors. We will further provide connections between the problems of maximal entropy and maximum likelihood estimation.

While the Shannon entropy is a property of a single distribution, the cross-entropy is a straight forward generalization towards pairs of distributions. We first introduce this quantity and then interpret the log-likelihood loss based on the cross-entropy.

**Definition 40** (Cross entropy and Kullback Leibler divergence). *The cross-entropy between two distributions  $\mathbb{P}[X_{[d]}]$  and  $\tilde{\mathbb{P}}[X_{[d]}]$  defined with respect to the same factored system is the quantity*

$$\mathbb{H}[\mathbb{P}, \tilde{\mathbb{P}}] = \sum_{x_0, \dots, x_{d-1}} \mathbb{P}[X_0 = x_0, \dots, X_{d-1} = x_{d-1}] \cdot \left( -\ln \left[ \tilde{\mathbb{P}}[X_{[d]} = x_{[d]}] \right] \right) = \left\langle \mathbb{P}, -\ln \left[ \tilde{\mathbb{P}} \right] \right\rangle [\emptyset] .$$

The cross-entropy is captured by the tensor network diagram



The Kullback-Leiber divergence between  $\mathbb{P}[X_{[d]}]$  and  $\tilde{\mathbb{P}}[X_{[d]}]$  is the quantity

$$D_{\text{KL}}[\mathbb{P} || \tilde{\mathbb{P}}] = \mathbb{H}[\mathbb{P}, \tilde{\mathbb{P}}] - \mathbb{H}[\mathbb{P}] .$$

Let us notice, that we have  $\mathbb{H}[\mathbb{P}, \tilde{\mathbb{P}}] = \infty$  if and only if there is a state  $x_{[d]} \in \times_{k \in [d]} [m_k]$  such that  $\mathbb{P}[X_{[d]} = x_{[d]}] > 0$  and  $\tilde{\mathbb{P}}[X_{[d]} = x_{[d]}] = 0$ .

The Gibbs inequality (for a proof see for example Chapter 2 in ?) states that for any distributions  $\mathbb{P}[X_{[d]}]$  and  $\tilde{\mathbb{P}}[X_{[d]}]$  we have

$$\mathbb{H}[\mathbb{P}, \tilde{\mathbb{P}}] \geq \mathbb{H}[\mathbb{P}] ,$$

where equality holds if and only if  $\mathbb{P} = \tilde{\mathbb{P}}$ . This ensures, that the Kullback-Leiber Divergence between any distributions is positive and vanishes if and only if both distributions coincide.

We in the next lemma provide an entropic interpretation of maximum likelihood estimation as defined in Def. 38.

**Lemma 3.** *The maximum likelihood estimation Problem  $\mathbb{P}_{\Gamma, \mathbb{P}^D}^M$  is equivalent to the minimization of cross-entropy and Kullback-Leibler divergence, that is*

$$\operatorname{argmin}_{\mathbb{P} \in \Gamma} \mathcal{L}_D(\mathbb{P}) = \operatorname{argmin}_{\mathbb{P} \in \Gamma} \mathbb{H}[\mathbb{P}^D, \mathbb{P}] = \operatorname{argmin}_{\mathbb{P} \in \Gamma} D_{\text{KL}}[\mathbb{P}^D || \mathbb{P}] .$$

*Proof.* Comparing the log-likelihood loss in Def. 38 with the cross-entropy in Def. 40, we get

$$\mathcal{L}_D(\mathbb{P}) = \mathbb{H}[\mathbb{P}^D, \mathbb{P}]$$

which established the equivalence of maximum likelihood estimation and cross-entropy minimization. Further, since

$$D_{\text{KL}}[\mathbb{P}^D || \mathbb{P}] = \mathbb{H}[\mathbb{P}^D, \mathbb{P}] - \mathbb{H}[\mathbb{P}^D]$$

and  $\mathbb{H}[\mathbb{P}^D]$  is a constant offset in the objective, maximum likelihood estimation is equivalent to the minimization of the Kullback-Leibler divergence.  $\square$

More general than Maximum Likelihood Estimation, we define the moment projection of an arbitrary distribution  $\mathbb{P}^*$  onto a set  $\Gamma$  of probability distributions as the problem

$$\operatorname{argmax}_{\mathbb{P} \in \Gamma} \mathbb{H}[\mathbb{P}^*, \mathbb{P}] . \quad (\mathbb{P}_{\Gamma, \mathbb{P}^*}^M)$$

With Lem. 3 we have established, that maximum likelihood estimation is the moment projection of the empirical distribution onto a set  $\Gamma$ .

We further define the information projection an arbitrary distribution  $\mathbb{P}^*$  onto a set  $\Gamma$  of probability distributions as the problem

$$\operatorname{argmax}_{\mathbb{P} \in \Gamma} \mathbb{H}[\mathbb{P}, \mathbb{P}^*] . \quad (\mathbb{P}_{\Gamma, \mathbb{P}^*}^I)$$

The cross-entropy is not symmetric, thus the information and the moment projections do not coincide in general. The differences of both are discussed in Chapter 8 in Koller and Friedman (2009).

**Example 4** (Cross entropy with respect to exponential families). *If  $\tilde{\mathbb{P}}$  is a member of an exponential family, we have*

$$\mathbb{H}[\mathbb{P}, \mathbb{P}^{(\mathcal{S}, \theta, \nu)}] = \left\langle \mathbb{P}, \ln[\mathbb{P}^{(\mathcal{S}, \theta, \nu)}] \right\rangle [\emptyset] = \left\langle \mathbb{P}, \sigma^{\mathcal{S}}, \theta \right\rangle [\emptyset] - A^{(\mathcal{S}, \nu)}(\theta) + \left\langle \mathbb{P}, \ln[\nu] \right\rangle [\emptyset] .$$

*The last term vanishes, given the convention  $0 \cdot \ln[0] = 0$ , if and only if for any  $x_{[d]}$  with  $\nu[X_{[d]} = x_{[d]}] = 0$  we have  $\mathbb{P}[X_{[d]} = x_{[d]}] = 0$ , and is infinite instead. Therefore, the cross entropy between a distribution and a member of an exponential family is finite, if and only if the distribution is representable with respect to the base measure  $\nu$  (see Def. 17). If  $\mathbb{P}$  is representable with respect to  $\nu$ , we can abbreviate the cross-entropy to*

$$\mathbb{H}[\mathbb{P}, \mathbb{P}^{(\mathcal{S}, \theta, \nu)}] = \left\langle \mathbb{P}, \phi^{(\mathcal{S}, \theta, \nu)}[X_{[d]}] \right\rangle [\emptyset] - A^{(\mathcal{S}, \nu)}(\theta) .$$

## 6.4 Variational Inference in Exponential Families

### 6.4.1 Forward and Backward Mappings

While defined for arbitrary distributions (see Def. 32), we now consider mean parameters to members of exponential families. First of all, they provide an alternative parameterization to the canonical parameter  $\theta$ . The computation of the mean parameter to a given canonical parameter and vice versa are the central inference problems in exponential families. We first formalize these inference problems by the forward and backward mapping and then provide in this section further insights into these mappings.

**Definition 41.** *Let  $\mathcal{S}$  be a statistic and  $\nu$  a base measure and consider the exponential family  $\Gamma^{\mathcal{S}, \nu}$ . The map*

$$F^{(\mathcal{S}, \nu)} : \mathbb{R}^p \rightarrow \mathcal{M}_{\mathcal{S}, \nu} \subset \mathbb{R}^p \quad , \quad F^{(\mathcal{S}, \nu)}(\theta) = \left\langle \mathbb{P}^{(\mathcal{S}, \theta, \nu)}[X_{[d]}], \sigma^{\mathcal{S}}[X_{[d]}, L] \right\rangle [L]$$

*is called the forward map of the exponential family and inverse, that is a map*

$$B^{(\mathcal{S}, \nu)} : \operatorname{im}(F^{(\mathcal{S}, \nu)}) \subset \mathbb{R}^p \rightarrow \mathbb{R}^p$$

*with  $\mathbb{P}^{(\mathcal{S}, B^{(\mathcal{S}, \nu)}(F^{(\mathcal{S}, \nu)}(\theta))), \nu} = \mathbb{P}^{(\mathcal{S}, \theta, \nu)}$  for any  $\theta[L] \in \mathbb{R}^p$ , a backward map of the exponential family.*

We notice, that the domain of  $F^{(\mathcal{S}, \nu)}$  is always  $\mathbb{R}^p$ , since the coordinates of  $F^{(\mathcal{S}, \nu)}(\theta)$  are for any  $\theta \in \mathbb{R}^p$  summations over finitely many products.

We already know by The. 18, that distributions representable by  $\nu$  have mean parameters in the interior of  $\mathcal{M}_{\mathcal{S}, \nu}$ . We now state that the elements of the corresponding exponential family  $\Gamma^{\mathcal{S}, \nu}$ , which are by construction representable by  $\nu$ , are expressive enough to reproduce the whole interior of  $\mathcal{M}_{\mathcal{S}, \nu}$ .

**Theorem 24.** *For any exponential family  $\Gamma^{\mathcal{S}, \nu}$  the image of the forward mapping is the mean polytope except its proper faces, that is*

$$\operatorname{im}(F^{(\mathcal{S}, \nu)}) = \mathcal{M}_{\mathcal{S}, \nu} / \bigcup_{Q_{\mathcal{S}, \nu}^{\mathcal{I}} \neq \mathcal{M}_{\mathcal{S}, \nu}} Q_{\mathcal{S}, \nu}^{\mathcal{I}} .$$

*If  $\mathcal{S}$  is minimal with respect to  $\nu$ , then  $\operatorname{im}(F^{(\mathcal{S}, \nu)}) = (\mathcal{M}_{\mathcal{S}, \nu})^\circ$ , the forward mapping is a bijection and the unique backward mapping its inverse.*

*Proof.* In case of minimal statistics, we refer for the proof of this statement to Theorem 3.3 in Wainwright and Jordan (2008). If  $\mathcal{S}$  is not minimal, we find a subset  $\tilde{\mathcal{S}}$  of its features such that  $\tilde{\mathcal{S}}$  is minimal with respect to  $\nu$  and there is a matrix  $M[L, \tilde{L}]$  such that for any distribution

$$\left\langle \mathbb{P}[X_{[d]}], \sigma^{\tilde{\mathcal{S}}}[X_{[d]}, \tilde{L}] \right\rangle [L] = \left\langle \mathbb{P}[X_{[d]}], \sigma^{\mathcal{S}}[X_{[d]}, L] \right\rangle [L].$$

This subset  $\tilde{\mathcal{S}}$  can be found by iteratively identifying  $V[L]$  and  $\lambda$  such that the condition in Def. 25 is violated, and dropping a feature  $S_i$  with  $V[L = i] \neq 0$ . At each manipulation step the expressivity of the exponential family stays constant and thus  $\Gamma^{\tilde{\mathcal{S}}, \nu} = \Gamma^{\mathcal{S}, \nu}$ . The matrix  $M[L, \tilde{L}]$  can be constructed based on the linear dependencies of the dropped features on the remaining. The procedure terminates, when there is no pair  $V[L], \lambda$ , which is equal to  $\tilde{\mathcal{S}}$  being minimal with respect to  $\nu$ . We then have

$$\mathcal{M}_{\mathcal{S}, \nu} / \bigcup_{Q_{\tilde{\mathcal{S}}, \nu}^{\mathcal{I}} \neq \mathcal{M}_{\mathcal{S}, \nu}} Q_{\tilde{\mathcal{S}}, \nu}^{\mathcal{I}} = \left\{ \left\langle \mu[\tilde{L}], M[L, \tilde{L}] \right\rangle [L] : \mu[\tilde{L}] \in \left( \mathcal{M}_{\tilde{\mathcal{S}}, \nu} \right)^{\circ} \right\}.$$

and using that  $\tilde{\mathcal{S}}$  is minimal we get

$$\text{im} \left( F^{(\mathcal{S}, \nu)} \right) = \mathcal{M}_{\mathcal{S}, \nu} / \bigcup_{Q_{\tilde{\mathcal{S}}, \nu}^{\mathcal{I}} \neq \mathcal{M}_{\mathcal{S}, \nu}} Q_{\tilde{\mathcal{S}}, \nu}^{\mathcal{I}}.$$

□

Forward and backward maps in an exponential family  $\Gamma^{\mathcal{S}, \nu}$  are the central classes of inference, which transform the description of a member by a canonical parameter into mean parameters and vice versa. Forward maps calculate to a canonical parameter  $\theta[L]$  the corresponding mean parameter  $\mu[L]$ . For any  $\theta[L]$  we have a closed form representation of this expectation query by the moment matching condition

$$\mu[L] = \left\langle \sigma^{\mathcal{S}}[X_{[d]}, L], \mathbb{P}^{(\mathcal{S}, \theta, \nu)}[X_{[d]}] \right\rangle [L].$$

The forward map is thus a collection of expectation queries (see Def. 35) to compute the coordinates of the mean parameter. The query  $\mathcal{S}_i$  asked against  $\mathbb{P}^{(\mathcal{S}, \theta, \nu)}[X_{[d]}]$  computes the coordinate  $\mu[L = i]$ . The contraction by the moment matching condition can, however, be infeasible, since it requires the instantiation of the probability distribution, which can be done by basis encodings of the statistic. We in this section provide alternative characterization of the forward map and approximations of it, which can be computed based on the selection encoding instead. Following Wainwright and Jordan (2008), we can characterize the forward mapping to exponential families as a variational problem and provide an alternative characterization to this contraction.

#### 6.4.2 Variational Formulation

We now formulate the forward and backward inference problems in exponential families as convex optimization problems.

The conjugate dual of the cumulative function is

$$\left( A^{(\mathcal{S}, \nu)} \right)^*(\mu) = \max_{\theta \in \mathbb{R}^p} \langle \mu, \theta \rangle [\emptyset] - A^{(\mathcal{S}, \nu)}(\theta).$$

**Lemma 4.** Let  $\theta \in \mathbb{R}^p$  and  $\mu[L]$  the corresponding mean parameter. For the gradient of  $A^{(\mathcal{S}, \nu)}$  evaluated at  $\theta$  we have

$$\nabla_{\tilde{\theta}[L]} |_{\theta} A^{(\mathcal{S}, \nu)}(\tilde{\theta}) = \mu[L]$$

If the statistic  $\mathcal{S}$  is minimal with respect to  $\nu$ , then also

$$\nabla_{\tilde{\mu}[L]} |_{\mu} \left( A^{(\mathcal{S}, \nu)} \right)^*(\tilde{\mu}) = \theta[L].$$

*Proof.* We have

$$\begin{aligned} \nabla_{\tilde{\theta}} |_{\tilde{\theta}} A^{(\mathcal{S}, \nu)}(\tilde{\theta}) &= \nabla_{\tilde{\theta}} |_{\tilde{\theta}} \ln \left[ \left\langle \exp \left[ \left\langle \sigma^{\mathcal{S}}[X_{[d]}, L], \theta[L] \right\rangle [X_{[d]}] \right], \nu \right\rangle [\emptyset] \right] \\ &= \frac{1}{\left\langle \exp \left[ \left\langle \sigma^{\mathcal{S}}[X_{[d]}, L], \theta[L] \right\rangle [X_{[d]}] \right], \nu \right\rangle [\emptyset]} \cdot \left\langle \sigma^{\mathcal{S}}[X_{[d]}, L], \exp \left[ \left\langle \sigma^{\mathcal{S}}[X_{[d]}, L], \theta[L] \right\rangle [X_{[d]}] \right], \nu \right\rangle [L] \\ &= \left\langle \mathbb{P}^{(\mathcal{S}, \theta, \nu)}[X_{[d]}], \sigma^{\mathcal{S}}[X_{[d]}, L] \right\rangle [L]. \end{aligned}$$

For the proof of the second claim we refer to Appendix B.2 in Wainwright and Jordan (2008). □



**Theorem 25** (Variational backward map). *For any  $\mu \in \mathcal{M}^\circ$  choose*

$$\theta [L] \in \operatorname{argmax}_\theta \langle \mu, \theta \rangle [\emptyset] - A^{(\mathcal{S}, \nu)}(\theta).$$

*Then  $\mu$  is the mean parameter to  $\mathbb{P}^{(\mathcal{S}, \theta, \nu)}$ .*

*Proof.* The maximization over  $\theta$  is an unconstrained concave maximization problem and the optimum is characterized by

$$0 = \nabla_\theta|_{\hat{\theta}} \left( \langle \mu, \theta \rangle [\emptyset] - A^{(\mathcal{S}, \nu)}(\theta) \right)$$

which reads

$$\mu [L] = \nabla_\theta|_{\hat{\theta}[L]} A^{(\mathcal{S}, \nu)}(\theta).$$

With Lem. 4, the gradient vanishes if and only if  $\mu$  is the mean parameter to  $\mathbb{P}^{(\mathcal{S}, \theta, \nu)}$ .  $\square$

The backward map is closely connected with the conjugate dual of  $A^{(\mathcal{S}, \nu)}$ . In particular, while the conjugate dual is defined by the maximization problem

$$(A^{(\mathcal{S}, \nu)})^*(\mu) = \max_{\theta \in \mathbb{R}^p} \langle \mu, \theta \rangle [\emptyset] - A^{(\mathcal{S}, \nu)}(\theta)$$

the backward map returns the position of the maximum in  $\mathbb{R}^p$ .

**Lemma 5.** *For any  $\theta$  with mean parameter  $\mu$  we have*

$$(A^{(\mathcal{S}, \nu)})^*(\mu) = -\mathbb{H} \left[ \mathbb{P}^{(\mathcal{S}, \theta, \nu)} \right].$$

*Proof.* We have

$$\theta \in \operatorname{argmax}_\theta \langle \mu, \theta \rangle [\emptyset] - A^{(\mathcal{S}, \nu)}(\theta)$$

if and only if

$$\mu [L] = \left\langle \mathbb{P}^{(\mathcal{S}, \theta, \nu)} [X_{[d]}], \sigma^{\mathcal{S}} [X_{[d]}, L] \right\rangle [L].$$

Therefore

$$(A^{(\mathcal{S}, \nu)})^*(\mu) = \left\langle \mathbb{P}^{(\mathcal{S}, \theta, \nu)} [X_{[d]}], \sigma^{\mathcal{S}} [X_{[d]}, L], \theta \right\rangle [\emptyset] - A^{(\mathcal{S}, \nu)}(\theta) = \left\langle \mathbb{P}^{(\mathcal{S}, \theta, \nu)} [X_{[d]}], \ln \left[ \mathbb{P}^{(\mathcal{S}, \theta, \nu)} [X_{[d]}] \right] \right\rangle [\emptyset] = -\mathbb{H} \left[ \mathbb{P}^{(\mathcal{S}, \theta, \nu)} \right].$$

$\square$

We can use this insight to provide a variational characterization of the forward mapping.

**Theorem 26** (Variational forward mapping). *Let  $\mathcal{S}$  be a minimal statistic with respect to  $\nu$ . Given  $\theta$ , there is a unique  $\mu$  with*

$$\mu \in \operatorname{argmax}_{\mu \in \mathcal{M}} \langle \mu, \theta \rangle [\emptyset] + \mathbb{H} \left[ \mathbb{P}^{\mathcal{S}, \mu, \nu} \right]$$

*and  $\mu$  is the mean parameter to  $\mathbb{P}^{(\mathcal{S}, \theta, \nu)} [X_{[d]}]$ . Here, we denote by  $\mathbb{P}^{\mathcal{S}, \mu, \nu}$  the member of the exponential family with the mean parameter  $\mu$ .*

*Proof.* By strong duality, we have  $(A^{(\mathcal{S}, \nu)})^{**} = A^{(\mathcal{S}, \nu)}$  and

$$A^{(\mathcal{S}, \nu)}(\theta) = \max_{\mu \in \mathcal{M}} \langle \mu, \theta \rangle [\emptyset] - (A^{(\mathcal{S}, \nu)})^*(\mu).$$

The statement then follows from the first order condition, where we use the gradient Lem. 4, and the characterization of  $(A^{(\mathcal{S}, \nu)})^*$  by Lem. 5.  $\square$

## 6.5 Maximum entropy distributions

We characterize in this section to any mean parameter the reproducing distribution with maximum entropy, and investigate their modes and tensor network representation. As we will show, the distributions with maximal entropy are in exponential families with base measure determined by the smallest face, which contains the mean parameter.

The mean parameters are computed as collections of expectation queries to  $S_l$ , which are answered against distributions in  $\Gamma^{\delta, \nu}$ . For any  $l \in [p]$  we have for the mean parameter  $\mu[L]$  reproduced by a distribution  $\mathbb{P}[X_{[d]}]$

$$\mu[L = l] = \mathbb{E}[S_l] = \langle \sigma^S[X_{[d]}, L = l], \mathbb{P}[X_{[d]}] \rangle [\emptyset] .$$

### 6.5.1 Entropy maximization problem

The entropy maximization problem with respect to matching expected statistics  $\mu^* \in \mathcal{M}_{S, \nu}$  is the optimization problem

$$\operatorname{argmax}_{\mathbb{P} \in \Gamma^{\delta, \nu}} \mathbb{H}[\mathbb{P}] \quad \text{subject to} \quad \langle \mathbb{P}, \sigma^S \rangle [L] = \mu^*[L] \quad (\mathbb{P}_{S, \nu, \mu^*}^{\mathbb{H}})$$

where by  $\Gamma^{\delta, \nu}$  we denote all distributions, which are representable with respect to a base measure  $\nu$ .

By definition of the mean polytope, Problem  $\mathbb{P}_{S, \nu, \mu^*}^{\mathbb{H}}$  has a feasible distribution if and only if  $\mu^*[L] \in \mathcal{M}_{S, \nu}$ . If this condition holds, we now characterize the solution of Problem  $\mathbb{P}_{S, \nu, \mu^*}^{\mathbb{H}}$ . First of all, we show that the maximum entropy distribution is in the exponential family  $\Gamma^{S, \nu}$ , when  $\mu$  does not lie on a proper face of  $\mathcal{M}_{S, \nu}$ . We then drop this assumption and generalize the statement to exponential families with refined base measures.

**Theorem 27.** *If the only face  $Q_{S, \nu}^T$  of  $\mathcal{M}_{S, \nu}$  with  $\mu \in Q_{S, \nu}^T$  is  $\mathcal{M}_{S, \nu}$  itself, then the solution of the maximum entropy distribution is the unique distribution*

$$\mathbb{P}^{S, \mu, \nu}[X_{[d]}] \in \Gamma^{S, \nu}$$

with  $\langle \mathbb{P}^{S, \mu, \nu}[X_{[d]}], \sigma^S[X_{[d]}, L] \rangle [L] = \mu[L]$ .

*Proof.* By The. 24, since by assumption

$$\mu[L] \in \mathcal{M}_{S, \nu} / \bigcup_{Q_{S, \nu}^T \neq \mathcal{M}_{S, \nu}} Q_{S, \nu}^T ,$$

there is a canonical parameter  $\theta$  with

$$\langle \mathbb{P}^{S, \theta, \nu}[X_{[d]}], \sigma^S[X_{[d]}, L] \rangle [L] = \mu[L]$$

For any other feasible distribution  $\tilde{\mathbb{P}}[X_{[d]}]$  we also have  $\langle \tilde{\mathbb{P}}[X_{[d]}], \sigma^S[X_{[d]}, L] \rangle [L] = \mu[L]$  and thus

$$\begin{aligned} \mathbb{H}[\tilde{\mathbb{P}}, \mathbb{P}^{(S, \theta, \nu)}] &= - \left\langle \tilde{\mathbb{P}}, \ln[\mathbb{P}^{(S, \theta, \nu)}[X_{[d]}]] \right\rangle [\emptyset] \\ &= - \left\langle \tilde{\mathbb{P}}, \langle \sigma^S[X_{[d]}, L], \theta[L] \rangle [X_{[d]}] \right\rangle [\emptyset] + A^{(S, \nu)}(\theta) \\ &= - \langle \theta, \mu \rangle [\emptyset] + A^{(S, \nu)}(\theta) \\ &= \mathbb{H}[\mathbb{P}^{(S, \theta, \nu)}] . \end{aligned}$$

With the Gibbs inequality we have if  $\tilde{\mathbb{P}} \neq \mathbb{P}^{(S, \theta, \nu)}$

$$\mathbb{H}[\mathbb{P}^{(S, \hat{\theta}, \nu)}] - \mathbb{H}[\tilde{\mathbb{P}}] = \mathbb{H}[\tilde{\mathbb{P}}, \mathbb{P}^{(S, \hat{\theta}, \nu)}] - \mathbb{H}[\tilde{\mathbb{P}}] > 0 .$$

Therefore, if  $\tilde{\mathbb{P}}$  does not coincide with  $\mathbb{P}^{(S, \hat{\theta}, \nu)}$ , it is not a solution of Problem  $\mathbb{P}_{S, \nu, \mu^*}^{\mathbb{H}}$ .  $\square$

Let us highlight the fact, that in Problem  $\mathbb{P}_{S, \nu, \mu^*}^{\mathbb{H}}$  we did not restrict to distributions in an exponential family and only demanded representability with respect to the base measure. When choosing the trivial base measure, this does not pose a restriction on the distributions. The. 27 states, that when the maximum entropy problem has a solution (i.e.  $\mu^* \in \mathcal{M}_{S, \nu}$ ), then the solution is in the exponential family to the statistic  $S$ .

When  $\mu^* \notin (\mathcal{M}_{S, \nu})^\circ$ , the mean parameter is by The. 18 not reproducible by a member of the exponential family  $\Gamma^{S, \nu}$ . Instead, in combination with the base measure refinement Algorithm 3, we show that the solution is in a refined exponential family.

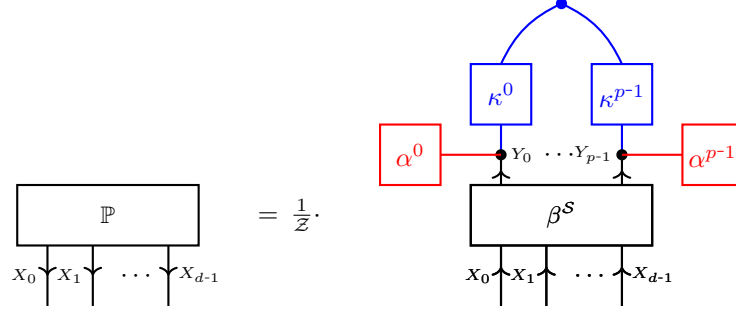


Figure 13: Tensor network decomposition of maximum entropy distributions to the constraint  $\mu[L] = \langle \mathbb{P}, \sigma^S \rangle [L]$ . Blue: Constraint activation cores  $\kappa^l$  in an CP decomposition, representing the face measure to the minimal face, such that  $\mu \in Q_{S,\nu}^{\mathcal{I}}$ . Red: Probabilistic activation cores  $\alpha^l$  in an elementary decomposition, where each leg core is a scaled exponentials evaluated on the enumerated image  $\text{im}(\mathcal{S}_l)$ .

**Theorem 28.** Let  $Q_{S,\nu}^{\mathcal{I}}$  be the minimal face of  $\mathcal{M}_{S,\nu}$  such that  $\mu[L] \in Q_{S,\nu}^{\mathcal{I}}$ . Then, the solution of the maximum entropy problem is the distribution

$$\mathbb{P}^{\mathcal{S},\mu,\tilde{\nu}}$$

where the base measure  $\tilde{\nu}$  is the refinement of  $\nu$  by the face measure  $\nu^{\mathcal{S},\mathcal{I}}$ , that is

$$\tilde{\nu}[X_{[d]}] = \langle \nu[X_{[d]}], \nu^{\mathcal{S},\mathcal{I}}[X_{[d]}] \rangle [X_{[d]}] .$$

*Proof.* By The. 20 all feasible distributions for the maximum entropy problem have to be representable by the face measure  $\nu^{\mathcal{S},\mathcal{I}}$ . Since the feasible distributions are further restricted to those representable by  $\nu$ , they are also representable by the refined base measure  $\tilde{\nu}$ . Now, by The. 19 the face itself is a polytope  $\mathcal{M}_{S,\tilde{\nu}}$  and the smallest face containing  $\mu$  is the polytope itself. We thus arrive at the claim by applying The. 27 on the polytope  $\mathcal{M}_{S,\tilde{\nu}}$ .  $\square$

The. 27 implies, that the by the face measure refined base measure  $\tilde{\nu}$  is minimal for the maximum entropy problem, in the sense that the solving distribution is positive with respect to it and all feasible distributions have to be representable by it. This highlights the fact, that the maximum entropy distribution does not vanish beyond those states, which are necessary to vanish to lie on the respective face.

### 6.5.2 Tensor Network Representation

Maximum entropy distributions with respect to constraints  $\mu[L] = \langle \mathbb{P}, \sigma^S \rangle [L]$  always have the sufficient statistic  $\mathcal{S}$ . They are represented in  $\Lambda^{\mathcal{S},\text{EL}}$ , if and only if the face measure of any face  $\mathcal{I}$  such that

$$\mu[L] \in Q_{S,\nu}^{\mathcal{I}} \mathcal{I}$$

is in  $\Lambda^{\mathcal{S},\text{EL}}$ . This is for example the case when  $\mu[L] \in (\mathcal{M}_{S,\nu})^\circ$ . In general, we find a CP decomposition as sketched in Figure 13.

### 6.5.3 Modes of maximum entropy distributions

We now show, that the face base measures coincide with the solutions of mode queries. Motivated by this fact, we define max cones of canonical parameters, which parametrized distributions coincide in their modes.

Let us now investigate, that faces of  $\mathcal{M}_{S,\nu}$  are the solutions of linear optimization problems constrained by  $\mathcal{M}_{S,\nu}$  (see Figure 14). Since such solution sets are intersections of the boundary of  $\mathcal{M}_{S,\nu}$  with half-spaces, they are faces. In the next theorem we show, how linear optimization problems are constructed to match a given face.

**Theorem 29.** For any non-empty face  $Q_{S,\nu}^{\mathcal{I}}$  to a subset  $\mathcal{I} \subset [n]$  there is a vector  $\theta[L]$ , which we call a normal of the face, such that

$$Q_{S,\nu}^{\mathcal{I}} = \text{argmax}_{\mu \in \mathcal{M}_{S,\nu}} \langle \theta[L], \mu[L] \rangle [\emptyset] .$$

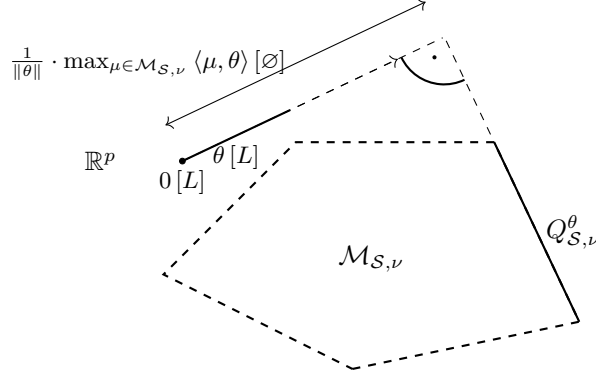


Figure 14: Sketch of a face  $Q_{S,\nu}^\theta$  to the normal  $\theta[L]$ . The face is characterized by those mean parameters in  $\mathcal{M}_{S,\nu}$ , which maximize the contraction with  $\theta[L]$ . The sketched distance of the origin  $0[\theta]$  of  $\mathbb{R}^p$  to the affine hull of the face is further related to the maximal contraction.

For any collection of positive  $\lambda_i$ , where  $i \in \mathcal{I}$ , the vector

$$\theta[L] = \sum_{i \in \mathcal{I}} \lambda_i \cdot a_i[L]$$

is a normal for  $Q_{S,\nu}^\mathcal{I}$ .

*Proof.* The first claim follows trivially from the second. To show the second claim, let there be for  $i \in \mathcal{I}$  arbitrary positive scalars  $\lambda_i$ . Since the face is non-empty, there is a  $\mu[L]$  with

$$\langle \mu[L], a_i[L] \rangle [\emptyset] = b_i$$

for all  $i \in \mathcal{I}$ . Since any  $\mu \in \mathcal{M}_{S,\nu}$  obey

$$\langle \mu[L], a_i[L] \rangle [\emptyset] \leq b_i$$

it follows that

$$\max_{\mu \in \mathcal{M}_{S,\nu}} \langle \theta[L], \mu[L] \rangle [\emptyset] = \sum_{i \in \mathcal{I}} \lambda_i \cdot b_i.$$

The maximum is attained at a  $\mu[L]$ , if and only if the equations  $\langle \mu[L], a_i[L] \rangle [\emptyset] = b_i$  are satisfied for  $i \in \mathcal{I}$ . This is equal to  $\mu[L] \in Q_{S,\nu}^\mathcal{I}$ .  $\square$

As we show next, also a converse statement holds, namely that for any vector  $\theta[L]$  we find a face  $Q_{S,\nu}^\mathcal{I}$ , such that the  $\theta[L]$  is a face normal to that face.

**Theorem 30.** For any  $\theta[L]$  we find a subset  $\mathcal{I} \subset [n]$ , such that

$$\operatorname{argmax}_{\mu \in \mathcal{M}_{S,\nu}} \langle \theta[L], \mu[L] \rangle [\emptyset] = Q_{S,\nu}^\mathcal{I}.$$

*Proof.* We first notice, that

$$\begin{aligned} & \operatorname{argmax}_{\mu \in \mathcal{M}_{S,\nu}} \langle \theta[L], \mu[L] \rangle [\emptyset] \\ &= \operatorname{conv} \left( \sigma^S [X_{[d]} = x_{[d]}, L] : x_{[d]} \in \operatorname{argmax}_{x_{[d]} \in \times_{k \in [d]} [m_k]} \langle \theta[L], \mathcal{S}(x_{[d]}) \rangle [\emptyset] \right). \end{aligned}$$

Further, since the contraction with  $\theta[L]$  is linear, the set  $\operatorname{argmax}_{\mu \in \mathcal{M}_{S,\nu}} \langle \theta[L], \mu[L] \rangle [\emptyset]$  is contained in the boundary of the polytope  $\mathcal{M}_{S,\nu}$ . We can conclude, that the set is a face, that is we find a subset  $\mathcal{I} \subset [n]$  with

$$Q_{S,\nu}^\mathcal{I} = \operatorname{argmax}_{\mu \in \mathcal{M}_{S,\nu}} \langle \theta[L], \mu[L] \rangle [\emptyset].$$

$\square$

In a slide abuse of notation, we denote in this case  $Q_{S,\nu}^\theta = Q_{S,\nu}^\mathcal{I}$ .

The. 30 provides a geometric perspective for mode queries. An arbitrary tensor  $\tau [X_{[d]}]$  can be understood to be a canonical parameter of the exponential family with statistic  $\delta$ , and base measure  $\mathbb{I} [X_{[d]}]$ . For this exponential family,  $\Gamma^{\delta,\mathbb{I}}$  coincides with the polytope of mean parameters, and is a standard simplex. Continuing our discussion in Sect. 6.1.2, we have for any  $\mathcal{U} \subset \times_{k \in [d]} [m_k]$  and  $\nu$  being the subset encoding of  $\mathcal{U}$  that

$$\max_{x_{[d]} \in \mathcal{U}} \tau [X_{[d]} = x_{[d]}] = \max_{\mu [X_{[d]}] \in \Gamma^{\delta,\nu}} \langle \tau [X_{[d]}], \mu [X_{[d]}] \rangle [\emptyset] .$$

The maximum is attained exactly for the mean parameters

$$\mu [X_{[d]}] \in \text{conv} \left( \epsilon_{x_{[d]}} [X_{[d]}] : x_{[d]} \in \text{argmax}_{x_{[d]} \in \mathcal{U}} \tau [X_{[d]} = x_{[d]}] \right) .$$

Answering the mode query is thus the characterization of the face of the standard simplex with face normal  $\tau [X_{[d]}]$ .

Let us now consider cases where the queried tensor  $\tau$  has a tensor network decomposition. This is the case for energy tensors to members of exponential families, for which we have a decomposition into selection encodings  $\sigma^S [X_{[d]}, L]$  and canonical parameters  $\theta [L]$ . In most generality we assume a decomposition of  $\tau$  by a tensor network on a hypergraph  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ , where  $[d] \subset \mathcal{V}$  as

$$\tau [X_{[d]}] = \langle \{ \tau^e [X_e] : e \in \mathcal{E} \} \rangle [X_{[d]}] .$$

Let us choose a subset  $\tilde{\mathcal{E}} \subset \mathcal{E}$  and

$$\max_{x_{[d]} \in \mathcal{U}} \tau [X_{[d]} = x_{[d]}] = \max_{x_{[d]} \in \mathcal{U}} \langle \{ \epsilon_{x_{[d]}} [X_{[d]}] \} \cup \{ \tau^e [X_e] : e \in \mathcal{E} \} \rangle [\emptyset]$$

We now split the contractions (see The. 131) to contract the cores  $\mathcal{E}/\tilde{\mathcal{E}}$  with the one-hot encoding first and keeping  $\tilde{\mathcal{V}} = \bigcup_{e \in \tilde{\mathcal{E}}} e$  open. With this we get

$$\begin{aligned} & \max_{x_{[d]} \in \mathcal{U}} \tau [X_{[d]} = x_{[d]}] \\ &= \max_{x_{[d]} \in \mathcal{U}} \left\langle \left\langle \{ \epsilon_{x_{[d]}} [X_{[d]}] \} \cup \{ \tau^e [X_e] : e \in \mathcal{E}/\tilde{\mathcal{E}} \} \right\rangle [X_{\tilde{\mathcal{V}}}], \left\langle \{ \tau^e [X_e] : e \in \tilde{\mathcal{E}} \} \right\rangle [X_{\tilde{\mathcal{V}}}] \right\rangle [\emptyset] . \end{aligned}$$

This optimization problem is the characterization of vectors, which convex hull is the face in the polytope

$$\mathcal{M} = \left\{ \left\langle \{ \epsilon_{x_{[d]}} [X_{[d]}] \} \cup \{ \tau^e [X_e] : e \in \mathcal{E}/\tilde{\mathcal{E}} \} \right\rangle [X_{\tilde{\mathcal{V}}}] : x_{[d]} \in \mathcal{U} \right\}$$

with the face normal

$$\theta [X_{\tilde{\mathcal{V}}}] = \left\langle \{ \tau^e [X_e] : e \in \tilde{\mathcal{E}} \} \right\rangle [X_{\tilde{\mathcal{V}}}] .$$

We define for each face the normal cone

$$C^\mathcal{I} = \{ \theta : \theta \in \mathbb{R}^p, \text{argmax}_{\mu \in \mathcal{M}_{S,\nu}} \langle \theta, \mu \rangle [\emptyset] = Q_{S,\nu}^\mathcal{I} \} .$$

**Lemma 6.** *For each face  $\mathcal{I}$ ,  $C^\mathcal{I}$  is a convex cone.*

*Proof.*  $C^\mathcal{I}$  is a cone, since for any  $\theta \in C^\mathcal{I}$  and  $\lambda > 0$  we have

$$Q_{S,\nu}^\mathcal{I} = \text{argmax}_{\mu \in \mathcal{M}_{S,\nu}} \langle \theta, \mu \rangle [\emptyset] = \text{argmax}_{\mu \in \mathcal{M}_{S,\nu}} \langle \lambda \cdot \theta, \mu \rangle [\emptyset] .$$

$C^\mathcal{I}$  is convex, since for  $\theta, \tilde{\theta} \in C^\mathcal{I}$  and  $\lambda \in (0, 1)$  we have

$$Q_{S,\nu}^\mathcal{I} = \text{argmax}_{\mu \in \mathcal{M}_{S,\nu}} \langle \theta, \mu \rangle [\emptyset] = \text{argmax}_{\mu \in \mathcal{M}_{S,\nu}} \langle \tilde{\theta}, \mu \rangle [\emptyset] .$$

It follows

$$Q_{S,\nu}^\mathcal{I} = \text{argmax}_{\mu \in \mathcal{M}_{S,\nu}} \langle \lambda \cdot \theta + (1 - \lambda) \cdot \tilde{\theta}, \mu \rangle [\emptyset]$$

and  $\lambda \cdot \theta + (1 - \lambda) \cdot \tilde{\theta} \in C^\mathcal{I}$ .

□

**Lemma 7.** For two non-empty faces  $\mathcal{I}^0$  and  $\mathcal{I}^1$  we have

$$Q_{\mathcal{S},\nu}^{\mathcal{I}^0} \subset Q_{\mathcal{S},\nu}^{\mathcal{I}^1}$$

if and only if

$$C_{\mathcal{S},\nu}^{\mathcal{I}^1} \subset \overline{C_{\mathcal{S},\nu}^{\mathcal{I}^0}}.$$

*Proof.* " $\Rightarrow$ " Let us assume  $Q_{\mathcal{S},\nu}^{\mathcal{I}^0} \subset Q_{\mathcal{S},\nu}^{\mathcal{I}^1}$  and let us choose arbitrary  $\theta_0 \in C_{\mathcal{S},\nu}^{\mathcal{I}^0}$ ,  $\theta_1 \in C_{\mathcal{S},\nu}^{\mathcal{I}^1}$ . It suffices to show, that for all  $\lambda \in (0, 1)$  we have  $\lambda \cdot \theta_0 + (1 - \lambda) \cdot \theta_1 \in C_{\mathcal{S},\nu}^{\mathcal{I}^0}$ , since this implies  $\theta_1 \in \overline{C_{\mathcal{S},\nu}^{\mathcal{I}^0}}$ . By assumption we have

$$\operatorname{argmax}_{\mu \in \mathcal{M}_{\mathcal{S},\nu}} \langle \theta_0, \mu \rangle [\emptyset] \subset \operatorname{argmax}_{\mu \in \mathcal{M}_{\mathcal{S},\nu}} \langle \theta_1, \mu \rangle [\emptyset]$$

and thus

$$\operatorname{argmax}_{\mu \in \mathcal{M}_{\mathcal{S},\nu}} \langle \lambda \cdot \theta_0 + (1 - \lambda) \cdot \theta_1, \mu \rangle [\emptyset] = \operatorname{argmax}_{\mu \in \mathcal{M}_{\mathcal{S},\nu}} \langle \theta_0, \mu \rangle [\emptyset]$$

which implies  $\lambda \cdot \theta_0 + (1 - \lambda) \cdot \theta_1 \in C_{\mathcal{S},\nu}^{\mathcal{I}^0}$ . " $\Leftarrow$ " Conversely, let us assume  $C_{\mathcal{S},\nu}^{\mathcal{I}^1} \subset \overline{C_{\mathcal{S},\nu}^{\mathcal{I}^0}}$ . We then find a sequence  $(\theta_n)_{n \in \mathbb{N}}$  with  $\theta_n \in C_{\mathcal{S},\nu}^{\mathcal{I}^0}$  for  $n \in \mathbb{N}$  and which limit  $\theta$  exists in  $C_{\mathcal{S},\nu}^{\mathcal{I}^1}$ . For an arbitrary  $x \in \mathcal{X}$  with  $\sigma^{\mathcal{S}}(x) \in Q_{\mathcal{S},\nu}^{\mathcal{I}^0}$  we have

$$\begin{aligned} \langle \theta, \sigma^{\mathcal{S}}(x) \rangle [\emptyset] &= \lim_{n \rightarrow \infty} \langle \theta_n, \sigma^{\mathcal{S}}(x) \rangle [\emptyset] \\ &= \lim_{n \rightarrow \infty} \max_{y \in \mathcal{X}} \langle \theta_n, \sigma^{\mathcal{S}}(y) \rangle [\emptyset] \\ &= \max_{y \in \mathcal{X}} \langle \theta, \sigma^{\mathcal{S}}(y) \rangle [\emptyset] \end{aligned}$$

and therefore  $\sigma^{\mathcal{S}}(x) \in Q_{\mathcal{S},\nu}^{\mathcal{I}^1}$ . Note, that we used in the third equation, that  $\mathcal{X}$  is finite. We use this property for all elements of the preimage of  $\mathcal{I}^0$  and get

$$\begin{aligned} Q_{\mathcal{S},\nu}^{\mathcal{I}^0} &= \operatorname{conv}(\sigma^{\mathcal{S}}(x) : x \in \operatorname{argmax}_{y \in \mathcal{X}} \langle \theta_1, \sigma^{\mathcal{S}}(y) \rangle [\emptyset]) \\ &\subset \operatorname{conv}(\sigma^{\mathcal{S}}(x) : x \in \operatorname{argmax}_{y \in \mathcal{X}} \langle \theta, \sigma^{\mathcal{S}}(y) \rangle [\emptyset]) \\ &= Q_{\mathcal{S},\nu}^{\mathcal{I}^1}. \end{aligned}$$

□

Lem. 7 suggests that the partial order of faces by inclusion is mimicked by another partial order on the max cones, sketched in Figure 15. We now consider the set of cones with this partial order, and show that this set is homomorphic to the face lattice.

**Theorem 31.** The max cone lattice of  $\mathcal{M}_{\mathcal{S},\nu}$ , partially order by the

$$C_{\mathcal{S},\nu}^{\mathcal{I}^0} \prec C_{\mathcal{S},\nu}^{\mathcal{I}^1} \quad \text{if and only if} \quad C_{\mathcal{S},\nu}^{\mathcal{I}^1} \subset \overline{C_{\mathcal{S},\nu}^{\mathcal{I}^0}}$$

is homomorphic to the face lattice of  $\mathcal{M}_{\mathcal{S},\nu}$  with the homomorphism

$$\psi(Q_{\mathcal{S},\nu}^{\mathcal{I}}) = C_{\mathcal{S},\nu}^{\mathcal{I}}.$$

*Proof.* We have to show that for all pairs of faces  $Q_{\mathcal{S},\nu}^{\mathcal{I}^0}, Q_{\mathcal{S},\nu}^{\mathcal{I}^1}$

$$\psi(Q_{\mathcal{S},\nu}^{\mathcal{I}^0}) \prec \psi(Q_{\mathcal{S},\nu}^{\mathcal{I}^1}) \Leftrightarrow Q_{\mathcal{S},\nu}^{\mathcal{I}^0} \prec Q_{\mathcal{S},\nu}^{\mathcal{I}^1}.$$

We show this first for the case, that  $Q_{\mathcal{S},\nu}^{\mathcal{I}^0} = \emptyset$  or  $Q_{\mathcal{S},\nu}^{\mathcal{I}^1} = \emptyset$ . Note, that the empty face is contained in any other face, but contains no non-empty face. Conversely, the trivial max cone  $C_{\mathcal{S},\nu}^{\emptyset}$  is for minimal statistics  $\{0[L]\}$  and contained in the boundary of any other non-empty cone. If the statistic is not minimal, the inclusion holds for the equivalence class of canonical parameters. Further, since the max cones are a disjoint partition of  $\mathbb{R}^p$ ,  $0[L]$  is not in any other max cone and thus  $C_{\mathcal{S},\nu}^{\mathcal{I}} \prec C_{\mathcal{S},\nu}^{\emptyset}$  does not hold for any non-empty  $\mathcal{I}$ .

For all pairs of non-empty faces  $Q_{\mathcal{S},\nu}^{\mathcal{I}^0}, Q_{\mathcal{S},\nu}^{\mathcal{I}^1}$ , Lem. 7 ensures that

$$\psi(Q_{\mathcal{S},\nu}^{\mathcal{I}^0}) \prec \psi(Q_{\mathcal{S},\nu}^{\mathcal{I}^1}) \Leftrightarrow Q_{\mathcal{S},\nu}^{\mathcal{I}^0} \prec Q_{\mathcal{S},\nu}^{\mathcal{I}^1}.$$

□

As a consequence of The. 31, the max cone lattice inherits all properties of the face lattice, for example those shown in Theorem 2.6 in Ziegler (2013).

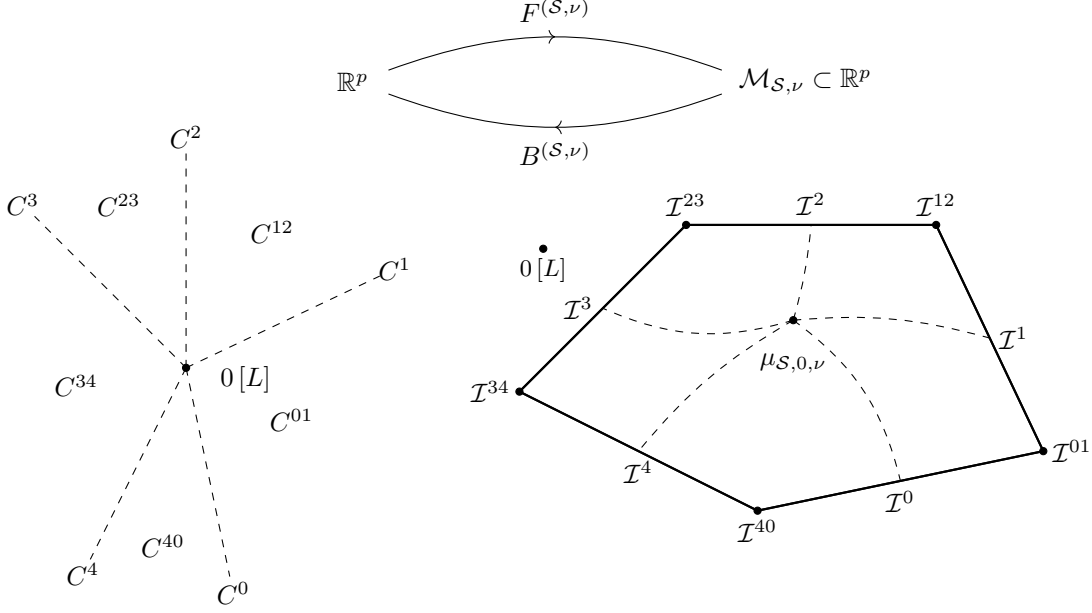


Figure 15: Partition of the canonical parameters in  $\mathbb{R}^p$  into convex and effectively open maximum cones  $C^{\mathcal{I}}$  to faces  $\mathcal{I}$  (left). The forward mapping maps the maximum cones into the mean parameter polytope (right).

#### 6.5.4 Base measure refinement

For mean parameters  $\mu[L]$  outside the interior of  $\mathcal{M}_{S,\nu}$  we know by The. 18, that any distribution with mean parameter  $\mu[L]$  is not positive with respect to  $\nu$  and is therefore not in the exponential family. We investigate this situation further and provide here a construction scheme to adapt the base measure such that there are exponential families containing these boundary distributions.

**Theorem 32.** *Let there be a statistic  $\mathcal{S}$ , which is minimal with respect to a base measure  $\nu$ , and  $\mu[L] \notin \mathcal{M}_{S,\nu}^\circ$ . Then there is a vector  $\theta[L] \in \mathbb{R}^p$  with*

$$\mu[L] \in \operatorname{argmax}_{\mu \in \mathcal{M}_{S,\nu}} \langle \theta[L], \mu[L] \rangle [\emptyset]$$

*and all distributions reproducing the mean parameter  $\mu[L]$  are representable with respect to the base measure*

$$\tilde{\nu}[X_{[d]}] = \langle \nu, \mathbb{I}_{\mathcal{U}}[X_{[d]}] \rangle [X_{[d]}],$$

*where the indicator is on the set*

$$\mathcal{U} = \operatorname{argmax}_{x_{[d]}} \langle \theta, \mathcal{S}(x_{[d]}) \rangle [\emptyset].$$

*Proof.* When  $\mu \notin \mathcal{M}_{S,\nu}^\circ$  we find a face such that  $\mu \in Q_{S,\nu}^{\mathcal{I}}$ . The existence of  $\theta[L]$  follows from The. 29, in which also a construction procedure is provided given a half-space representation (see The. 17). Now, we have

$$\mu[L] \in \operatorname{argmax}_{\mu \in \mathcal{M}_{S,\nu}} \langle \theta[L], \mu[L] \rangle [\emptyset]$$

and thus

$$\mu[L] \in \operatorname{conv} \left( \sigma^{\mathcal{S}} X_{[d]} = x_{[d]}, L : x_{[d]} \in \operatorname{argmax}_{x_{[d]} : \nu[X_{[d]}=x_{[d]}]=1} \langle \theta[L], \sigma^{\mathcal{S}} X_{[d]} = x_{[d]}, L \rangle [\emptyset] \right)$$

Thus, any distribution reproducing the mean parameter is a convex combination of the one-hot encodings of the states in  $\operatorname{argmax}_{x_{[d]}} \langle \theta[L], \sigma^{\mathcal{S}} X_{[d]} = x_{[d]}, L \rangle [\emptyset]$ , and therefore representable with respect to the base measure  $\tilde{\nu}$ .  $\square$

Each face of  $\mathcal{M}_{S,\nu}$  thus defines a refinement of a base measure, which is sufficient to reproduce the mean parameters on that face.

**Definition 42.** *The base measure to the face of  $\mathcal{M}$  with normal  $\theta$  is*

$$\nu^{S,\mathcal{I}}[X_{[d]}] = \nu^{S,\theta}[X_{[d]}]$$

The mean parameter to the normalized face base measure is

$$\mu(Q_{\mathcal{S},\nu}^\theta) = \langle \sigma^\mathcal{S} [X_{[d]}, L], \langle \nu^{\mathcal{S},\theta} \rangle [X_{[d]}|\emptyset] \rangle [L] .$$

This quantity will be interesting as the limit of annealing (see Chapter 6).

The. 32 implies that any mean parameter on a face of  $\mathcal{M}_{\mathcal{S},\nu}$  can be reproduced with a distribution representable with respect to the refined base measure

$$\tilde{\nu} [X_{[d]}] = \langle \nu, \nu^{\mathcal{S},\theta} \rangle [X_{[d]}] .$$

We now utilize these findings and provide by Algorithm 3 a procedure to refine the base measure until the reduced mean parameter is in the interior of a reduced mean parameter polytope.

---

**Algorithm 3** Base Measure Refinement

---

**Require:**  $\nu$ , statistic  $\mathcal{S}$  and mean parameter  $\mu \in \mathcal{M}_{\mathcal{S},\nu}$

**Ensure:** Refined base measure  $\tilde{\nu}$ , reduced statistic  $\tilde{\mathcal{S}}$  and reduced mean parameter  $\tilde{\mu}$  with properties described in The. 33

---

**while**  $\mu \notin (\mathcal{M}_{\mathcal{S},\nu})^\circ$  **do**

**while**  $\mathcal{S}$  not minimal with respect to  $\nu$  (see Def. 25) **do**

        Find non-vanishing vector  $V[L]$  and scalar  $\lambda \in \mathbb{R}$  such that

$$\langle \sigma^\mathcal{S} [X_{[d]}, L], V[L], \nu [X_{[d]}] \rangle [X_{[d]}] = \lambda \cdot \nu [X_{[d]}] .$$

        Choose a coordinate  $l \in [p]$  with  $V[L = l] \neq 0$  and drop it from  $\mathcal{S}$  and  $\mu$

**end while**

    Find a non-trivial face (i.e. a non-empty face, which is a proper subset of  $\mathcal{M}_{\mathcal{S},\nu}$ ) with normal  $\theta$ , such that

$$\mu \in Q_{\mathcal{S},\nu}^\theta$$

    Refine base measure

$$\nu \leftarrow \langle \nu, \nu^{\mathcal{S},\theta} \rangle [X_{[d]}]$$

**end while**

**return**  $\nu, \mathcal{S}, \mu$

---

**Theorem 33.** For arbitrary inputs  $\nu, \mathcal{S}$  and  $\mu \in \mathcal{M}_{\mathcal{S},\nu}$ , Algorithm 3 terminates in finite time and outputs a triple of base measure  $\tilde{\nu}$ , statistic  $\tilde{\mathcal{S}}$  and mean parameter  $\tilde{\mu}$  such that the following holds. Any probability distribution  $\mathbb{P}$  reproduces  $\mu$ , if and only if it reproduces  $\tilde{\mu}$ . Further, any probability distribution reproducing  $\mu$  is representable with respect to  $\tilde{\nu}$  and  $\tilde{\mu} \in (\mathcal{M}_{\tilde{\mathcal{S}},\tilde{\nu}})^\circ$ .

*Proof.* Let us first show, that Algorithm 3 always terminates. The inner while loop of Algorithm 3 always terminates, since  $\mathcal{S}$  has a finite number of coordinates, and in each iteration one of the coordinates is dropped. To show that the outer while loop also terminates, it suffices to show, that the non-vanishing coordinates of the refined base measure are a proper subset of the base measure before refinement. But if this would not be the case, we would have

$$\nu [X_{[d]}] = \langle \nu, \nu^{\mathcal{S},\theta} \rangle [X_{[d]}]$$

and thus  $Q_{\mathcal{S},\nu}^\theta = \mathcal{M}_{\mathcal{S},\nu}$ , which is a contradiction with the assumption of a non-trivial face.

The second and third claim follow from an iterative application of The. 32 and the fact, that a probability distribution reproduces  $\mu$  in a non-minimal representation, if and only if it reproduces the corresponding reduced  $\tilde{\mu}$  with respect to the reduced statistics.  $\square$

**Example 5** (Faces with normals parallel to one-hot encodings). To get some intuition how to represent face base measures, let us consider face normals  $\theta \in \{\lambda \cdot \epsilon_l [L] : l \in [p], \lambda \in \mathbb{R}/\{0\}\}$ . We use basis encodings of the coordinates  $\mathcal{S}_l$  of the statistic  $\mathcal{S}$ , with head variables  $X_{\mathcal{S}_l}$  with dimension  $m_{\mathcal{S}_l}$  enumerating the image  $\text{im}(\mathcal{S}_l) \subset \mathbb{R}$  in an ascending order. If  $\theta [L] = \lambda \cdot \epsilon_l [L]$  with  $\lambda > 0$ , then  $\text{argmax}_{x_{[d]}} \langle \theta, \mathcal{S}(x_{[d]}) \rangle [\emptyset]$  consists of states  $x_{[d]}$  with minimal statistic  $\mathcal{S}_l [X_{[d]}] = x_{[d]}$ , that is

$$\nu^{\mathcal{S},\lambda \cdot \epsilon_l} [X_{[d]}] = \left\langle \beta^{\mathcal{S}_l} [X_{[d]}, X_{\mathcal{S}_l}], \epsilon_{m_{\mathcal{S}_l}-1} [X_{\mathcal{S}_l}] \right\rangle [X_{[d]}] .$$

If  $\theta [L] = \lambda \cdot \epsilon_l [L]$  with  $\lambda < 0$ , then at the states with minimal statistic  $\mathcal{S}_l [X_{[d]}] = x_{[d]}$ , that is

$$\nu^{\mathcal{S},\lambda \cdot \epsilon_l} [X_{[d]}] = \left\langle \beta^{\mathcal{S}_l} [X_{[d]}, X_{\mathcal{S}_l}], \epsilon_0 [X_{\mathcal{S}_l}] \right\rangle [X_{[d]}] .$$



**Theorem 34.** For the maximal graph  $\mathcal{G}^{\max} = ([p], \{[p]\})$ , which has a single hyperedge containing all head variables we have

$$\mathcal{M}_{\mathcal{S}, \nu} = \left\{ \langle \mathbb{P} [X_{[d]}], \sigma^{\mathcal{S}} [X_{[d]}, L] \rangle [X_{[d]}], \mathbb{P} \in \Lambda^{\mathcal{S}, \mathcal{G}^{\max}} \right\}.$$

*Proof.* It is enough show, that for any output tuples  $\tilde{\nu}, \tilde{\mathcal{S}}$  of the Base Measure Refinement Algorithm 3 we have

$$\Gamma^{\tilde{\nu}, \tilde{\mathcal{S}}} \subset \Lambda^{\mathcal{S}, \mathcal{G}^{\max}}.$$

We notice, that the normalization of any face base measure is realizable by  $\Lambda^{\mathcal{S}, \mathcal{G}^{\max}}$ . Providing a more technical argument, we have

$$\mathbb{I}_{\arg\max_{x_{[d]}} \langle \theta, \mathcal{S}(x_{[d]}) \rangle [\emptyset]} [X_{[d]}] = \left\langle \beta^{\mathcal{S}}, \sum_{\mathcal{S}(x_{[d]}) : x_{[d]} \in \arg\max_{x_{[d]}} \langle \theta, \mathcal{S}(x_{[d]}) \rangle [\emptyset]} \epsilon_{I(\mathcal{S}(x_{[d]})})} [Y_{[p]}] \right\rangle [X_{[d]}].$$

Since during the execution of Algorithm 3,  $\tilde{\mathcal{S}}$  is a subset of  $\mathcal{S}$ , we can find a corresponding  $\theta_i$  extending the face normal by vanishing coordinates to  $\mathcal{S}$ . We then have, that

$$\tilde{\nu} = \left\langle \{ \beta^{\mathcal{S}} [Y_{[p]}, X_{[d]}] \} \cup \left\{ \sum_{\mathcal{S}(x_{[d]}) : x_{[d]} \in \arg\max_{x_{[d]}} \langle \theta_i, \mathcal{S}(x_{[d]}) \rangle [\emptyset]} \epsilon_{I(\mathcal{S}(x_{[d]})})} [Y_{[p]}] : i \in [n] \right\} \right\rangle [Y_{[p]}]$$

represents the output base measure, where  $i \in [n]$  label the faces chosen during in the loop of Algorithm 3. Now, any member  $\mathbb{P}^{\tilde{\nu}, \theta, \tilde{\mathcal{S}}} \in \Gamma^{\tilde{\nu}, \tilde{\mathcal{S}}}$  can be represented by a member of  $\Lambda^{\mathcal{S}, \mathcal{G}^{\max}}$ , by contracting these base measure representing cores with the activation cores  $\bigotimes_{l \in [p]} \alpha^{\mathcal{S}_l, \theta[L=l]} [Y_l]$ .  $\square$

### 6.5.5 Mean parameters by soft and hard constraints

We provide further intuition on the position of the mean parameter, when interpreting  $\nu$  and  $\mathcal{S}$  as soft and hard constraint mechanisms on a distribution. To this end, let  $\mathcal{X}$  be an arbitrary finite state set, where we so far had the set  $\times_{k \in [d]} [m_k]$  for factored representations.

We recall the statistics encoding of states as a map

$$\sigma^{\mathcal{S}} : \mathcal{X} \rightarrow \mathbb{R}^p, \quad \sigma^{\mathcal{S}}(x) = \sigma^{\mathcal{S}} [X_{[d]} = x_{[d]}, L].$$

The mean polytope is then the subset of the polytope  $\mathcal{M}_{\mathcal{S}, \mathbb{I}}$ ,

$$\mathcal{M}_{\mathcal{S}, \nu} = \text{conv} (\sigma^{\mathcal{S}} [X_{[d]} = x_{[d]}, L] : \nu [X_{[d]} = x_{[d]}] = 1)$$

The mean parameter of a base measure  $\nu$ , which normalization is understood as a distribution is then

$$\begin{aligned} \mu [\mathcal{S}, 0, \nu] &= \frac{1}{\langle \nu \rangle [\emptyset]} \sum_{x_{[d]} : \nu [X_{[d]} = x_{[d]}] = 1} \sigma^{\mathcal{S}} [X_{[d]} = x_{[d]}, L] \\ &= \langle \langle \nu \rangle [X_{[d]} | \emptyset], \sigma^{\mathcal{S}} [X_{[d]}, L] \rangle [L]. \end{aligned}$$

We understand  $\mu [\mathcal{S}, 0, \nu]$  as the weighted center of  $\mathcal{M}_{\mathcal{S}, \nu}$ , where weights are allocated on the image of the statistics encoding by the cardinality of their pre-images. These concepts are sketched in Figure 16 in blue.

The choice of a base measure  $\nu$  can be regarded as the hard structure of a distribution. Given  $\nu$ , let us now implement soft structure in form of distributions in exponential families  $\Gamma^{\mathcal{S}, \nu}$ . While the normalized base measure  $\nu$  reproduces only the weighted center of  $\mathcal{M}_{\mathcal{S}, \nu}$ , canonical parameters can be chosen to alter to any interior point in  $\mathcal{M}_{\mathcal{S}, \nu}$ . This alternations are sketched in Figure 16 in red.

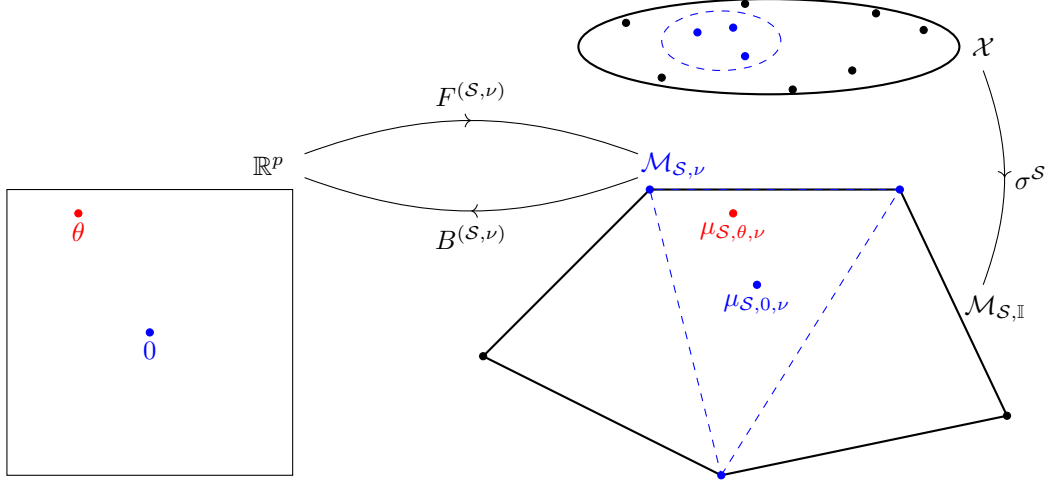


Figure 16: Sketch of hard constraints (blue) and soft constraints (red) determining the position of the mean parameter in  $\mu_{S,\mathbb{I}}$ . In the set of states  $\mathcal{X}$ , the support of a base measure  $\nu$  is marked blue, and the base measure is the corresponding subset encoding. The normalized  $\nu$  has the mean parameter  $\mu_{S,0,\nu}$ , and the mean polytope  $\mathcal{M}_{S,\nu}$  is the convex hull of the image to the statistics encoding  $\sigma^S$  restricted on the support set of  $\nu$ . Given a base measure  $\nu$ , the corresponding forward and backward maps implement soft constraints depending on canonical parameters  $\theta$ . The forward mapping maps the zero canonical parameter to  $\mu_{S,0,\nu}$ , correspond with the normalized base measure. Generic parameters  $\theta$  (red) are mapped onto the interior  $\mathcal{M}_{S,\nu}^\circ$ .

## 6.6 Forward Mapping in Exponential Families

### 6.6.1 Mode queries by annealing

The mode of a distribution is related to the forward mapping of  $\beta \cdot \theta$  in the limit  $\beta \rightarrow \infty$  of low temperatures. To sketch this relation, we recall the variational formulation of mode queries by

$$\text{conv} \left( \sigma^S [X_{[d]} = x_{[d]}, L] : x_{[d]} \in \text{argmax}_{x_{[d]}} \langle \theta, \mathcal{S}(x_{[d]}) \rangle [\emptyset] \right) = \text{argmax}_{\mu \in \mathcal{M}_{S,\nu}} \langle \mu, \theta \rangle [\emptyset] .$$

Further, for any by the inverse temperature  $\beta \neq 0$  annealed canonical parameter  $\theta [L]$  (see Sect. 6.2.3) we have

$$\text{argmax}_{\mu \in \mathcal{M}_{S,\nu}} \langle \mu, \beta \cdot \theta \rangle [\emptyset] + \mathbb{H} [\mathbb{P}^\mu] = \text{argmax}_{\mu \in \mathcal{M}_{S,\nu}} \langle \mu, \theta \rangle [\emptyset] + \frac{1}{\beta} \cdot \mathbb{H} [\mathbb{P}^\mu] .$$

In the annealing limit, that is for large  $\beta$ , the entropy term becomes negligible and the forward mapping tends to the convex hull

$$\text{conv} \left( \sigma^S [X_{[d]} = x_{[d]}, L] : x_{[d]} \in \text{argmax}_{x_{[d]}} \langle \theta, \mathcal{S}(x_{[d]}) \rangle [\emptyset] \right) = \text{argmax}_{\mu \in \mathcal{M}_{S,\nu}} \langle \mu, \theta \rangle [\emptyset] .$$

For a more detailed discussion of this relation, see Theorem 8.1 in Wainwright and Jordan (2008).

## 6.7 Mean Field Methods

Mean field methods are approximation schemes for forward mappings, designed for efficient inference. To introduce them we turn the maximization over the mean parameter polytope in the the variational principle of forwards mappings into a maximization over the reproducing distributions, that is

$$\max_{\mu \in \mathcal{M}_{S,\nu}} \langle \mu, \theta \rangle [\emptyset] + \mathbb{H} [\mathbb{P}^\mu] = \max_{\mathbb{P} \in \Gamma^{\delta,\nu}} \langle \phi, \mathbb{P} \rangle [\emptyset] + \mathbb{H} [\mathbb{P}]$$

where

$$\phi = \langle \sigma^S, \theta \rangle [X_{[d]}] .$$

Mean field methods now provide lower bounds on this maximization by restricting the distribution optimized over to a tractable subset of distributions.

Let  $\Gamma$  be a subset of  $\Gamma^{\delta, \mathbb{I}}$ , we state the mean field problem as

$$\operatorname{argmax}_{\mathbb{P} \in \Gamma^{\mathcal{S}^{\mathcal{G}}, \mathbb{I}}} \langle \phi, \mathbb{P} \rangle [\emptyset] + \mathbb{H} [\mathbb{P}] \quad (\mathbb{P}_{\Gamma, \mathbb{P}^\phi}^I)$$

We show next, that the mean field problem is an instance of an information projection, as indicated by the notation.

**Theorem 35.** *For any hypergraph  $\mathcal{G}$  and energy tensor  $\phi$  we have for  $\mathbb{P}^\phi [X_{[d]}] = \langle \phi \rangle [X_{[d]} | \emptyset]$*

$$\operatorname{argmax}_{\mathbb{P} \in \Gamma} \langle \phi, \mathbb{P} \rangle [\emptyset] + \mathbb{H} [\mathbb{P}] = \operatorname{argmin}_{\mathbb{P} \in \Gamma} D_{\text{KL}} [\mathbb{P} || \mathbb{P}^\phi]$$

*Problem  $\mathbb{P}_{\Gamma, \mathbb{P}^\phi}^I$  is thus the information projection of a distribution  $\mathbb{P}^\phi [X_{[d]}] = \langle \phi \rangle [X_{[d]} | \emptyset]$  onto the exponential family  $\Gamma^{\mathcal{S}^{\mathcal{G}}, \mathbb{I}}$ .*

*Proof.* The cross entropy between a  $\mathbb{P} \in \Gamma$  and  $\mathbb{P}^\phi$  is

$$\begin{aligned} \mathbb{H} [\mathbb{P}, \mathbb{P}^\phi] &= \langle \mathbb{P} [X_{[d]}], -\ln [\mathbb{P}^\phi [X_{[d]}]] \rangle [\emptyset] \\ &= \langle \mathbb{P} [X_{[d]}], -\phi [X_{[d]}] \rangle [\emptyset] + \langle \mathbb{P} [X_{[d]}] \rangle [\emptyset] \cdot \ln [\langle \exp [\phi [X_{[d]}]] \rangle [\emptyset]] \end{aligned}$$

Together we have, that

$$\begin{aligned} D_{\text{KL}} [\mathbb{P} || \mathbb{P}^\phi] &= \mathbb{H} [\mathbb{P}, \mathbb{P}^\phi] - \mathbb{H} [\mathbb{P}] \\ &= -\langle \mathbb{P} [X_{[d]}], \phi [X_{[d]}] \rangle [\emptyset] - \mathbb{H} [\mathbb{P}] + \ln [\langle \exp [\phi [X_{[d]}]] \rangle [\emptyset]] \end{aligned}$$

Since the last term is constant among  $\mathbb{P} \in \Gamma$ , it holds that

$$\operatorname{argmax}_{\mathbb{P} \in \Gamma} \langle \phi, \mathbb{P} \rangle [\emptyset] + \mathbb{H} [\mathbb{P}] = \operatorname{argmax}_{\mathbb{P} \in \Gamma} -D_{\text{KL}} [\mathbb{P} || \mathbb{P}^\phi] = \operatorname{argmax}_{\mathbb{P} \in \Gamma} D_{\text{KL}} [\mathbb{P} || \mathbb{P}^\phi] .$$

□

### 6.7.1 Naive Mean Field Method

In the naive mean field method, we choose the approximating set  $\Gamma$  by the exponential family of Markov Networks  $\Gamma^{\mathcal{S}^{\mathcal{G}^{\text{EL}}, \mathbb{I}}}$  on the elementary graph

$$\mathcal{G}^{\text{EL}} = ([d], \{\{k\} : k \in [d]\}) .$$

Markov Networks on this graph are represented by normed leg cores

$$\mathbb{P} [X_{[d]}] = \bigotimes_{k \in [d]} \rho^k [X_k]$$

Problem  $\mathbb{P}_{\Gamma, \mathbb{P}^\phi}^I$  is for this instance of the form

$$\operatorname{argmax}_{\rho^k [X_k] : \langle \rho^k \rangle [\emptyset] = 1, k \in [d]} \langle \{\phi\} \cup \{\rho^k : k \in [d]\} \rangle [\emptyset] + \sum_{k \in [d]} \mathbb{H} [\rho^k] .$$

We approximately solve this problem in Algorithm 4 by alternation through the leg cores and performing a locally optimal leg core update. The optimal update equations are derived as the next theorem.

**Theorem 36** (Update equations for the naive mean field approximation). *For any  $k \in [d]$  and leg cores  $\{\rho^{\tilde{k}} [X_{\tilde{k}}] : \tilde{k} \in [d], \tilde{k} \neq k\}$  the local problem*

$$\operatorname{argmax}_{\rho^k [X_k] : \langle \rho^k \rangle [\emptyset] = 1} \langle \{\phi\} \cup \{\rho^{\tilde{k}} [X_{\tilde{k}}] : \tilde{k} \in [d]\} \rangle [\emptyset] + \sum_{k \in [d]} \mathbb{H} [\rho^k]$$

*is solved at*

$$\rho^k [X_k] = \left\langle \exp \left[ \left\langle \{\phi[X_{[d]}]\} \cup \{\rho^{\tilde{k}} [X_{\tilde{k}}] : \tilde{k} \neq k\} \right\rangle [X_{[d]}] \right] \right\rangle [X_k | \emptyset] .$$

*Proof.* We have

$$\frac{\partial \mathbb{H} [\rho^k]}{\partial \rho^k} = -\ln [\rho^k [X_k]] + \mathbb{I} [X_k]$$

and by multilinearity of tensor contractions

$$\frac{\partial \left\langle \{\phi\} \cup \{\rho^{\tilde{k}} : \tilde{k} \in [d]\} \right\rangle [\emptyset]}{\partial \rho^k} = \left\langle \{\phi\} \cup \{\rho^{\tilde{k}} : \tilde{k} \in [d], \tilde{k} \neq k\} \right\rangle [X_k] .$$

Combining both, the condition

$$0 = \frac{\partial \left( \left\langle \{\phi\} \cup \{\rho^{\tilde{k}} : \tilde{k} \in [d]\} \right\rangle [\emptyset] + \sum_{k \in [d]} \mathbb{H}[\rho^k] \right)}{\partial \rho^k}$$

is equal to

$$\ln [\rho^k [X_k]] = \mathbb{I}[X_k] + \left\langle \{\phi\} \cup \{\rho^{\tilde{k}} : \tilde{k} \in [d], \tilde{k} \neq k\} \right\rangle [X_k] .$$

Together with the condition  $\langle \rho^k \rangle [\emptyset] = 1$  this is satisfied at

$$\rho^k [X_k] = \left\langle \exp \left[ \left\langle \{\phi\} \cup \{\rho^{\tilde{k}} : \tilde{k} \neq k\} \right\rangle [X_k] \right] \right\rangle [X_k | \emptyset] .$$

□

Algorithm 4 is the alternation of legwise updates until a stopping criterion is met.

---

**Algorithm 4** Naive Mean Field Approximation

---

**Require:** Energy tensor  $\phi[\phi]$

**Ensure:** Tensor Network  $\{\rho^k[] : k \in [d]\}$  approximating  $\langle \phi[X_{[d]}] \rangle [X_{[d]} | \emptyset]$

---

**for**  $k \in [d]$  **do**

$$\rho^k [X_k] \leftarrow \langle \mathbb{I} \rangle [X_k | \emptyset]$$

**end for**

**while** Stopping criterion is not met **do**

**for**  $k \in [d]$  **do**

$$\rho^k [X_k] \leftarrow \left\langle \exp \left[ \left\langle \{\phi[X_{[d]}]\} \cup \{\rho^{\tilde{k}}[X_{\tilde{k}}] : \tilde{k} \neq k\} \right\rangle [X_k] \right] \right\rangle [X_k | \emptyset]$$

**end for**

**end while** **return**  $\{\rho^k[] : k \in [d]\}$

---

### 6.7.2 Structured Variational Approximation

We now generalize the naive mean field method towards generic families of Markov Networks. Let  $\mathcal{G}$  be any hypergraph, the structured variational approximation method is Problem  $P_{\Gamma, \mathbb{P}^\phi}^I$  with

$$\Gamma = \Gamma^{\mathcal{S}^{\mathcal{G}}, \mathbb{I}} .$$

We approximate the solution of this problem again by an alternating algorithm, which iteratively updates the cores of the approximating Markov Network.

**Theorem 37** (Update equations for the structured variational approximation). *The Markov Network  $\tau^{\mathcal{G}}$  with hypercores  $\{\tau^e[X_e] : e \in \mathcal{E}\}$  is a stationary point for structured variational approximation, if for all  $e \in \mathcal{E}$  we find a  $\lambda > 0$  with*

$$\tau^e[X_e] = \lambda \cdot \exp \left[ \frac{\left\langle \{\phi\} \cup \{\tau^{\tilde{e}} : \tilde{e} \neq e\} \right\rangle [X_e]}{\left\langle \{\tau^{\tilde{e}} : \tilde{e} \neq e\} \right\rangle [X_e]} - \sum_{\tilde{e} \neq e} \frac{\left\langle \{\ln[\tau^{\tilde{e}}]\} \cup \{\tau^{\tilde{e}} : \tilde{e} \neq \hat{e}\} \right\rangle [X_e]}{\left\langle \{\tau^{\tilde{e}} : \tilde{e} \neq \hat{e}\} \right\rangle [X_e]} \right] .$$

Here, the quotient denotes the coordinatewise quotient.

*Proof.* We proof the theorem by the first order condition on the objective

$$O(\tau^G) = \langle \phi, \langle \tau^G \rangle [X_V | \emptyset] \rangle [\emptyset] + \mathbb{H} [\langle \tau^G \rangle [X_V | \emptyset]]$$

We further use Lem. 31, which shows a characterization of the derivative of functions dependent on tensors.

We have

$$\langle \phi, \langle \tau^G \rangle [X_{[d]} | \emptyset] \rangle [\emptyset] = \frac{\langle \{\phi\} \cup \tau^G \rangle [\emptyset]}{\langle \tau^G \rangle [\emptyset]}.$$

Further we have

$$\mathbb{H} [\langle \tau^G \rangle [X_{[d]} | \emptyset]] = \left( \sum_{\tilde{e} \in \mathcal{E}} \langle -\ln [\tau^{\tilde{e}}], \langle \tau^G \rangle [X_{[d]} | \emptyset] \rangle [\emptyset] \right) + \ln [\langle \tau^G \rangle [\emptyset]]$$

We define the tensor

$$\tilde{\tau}[X_V] = \phi[X_V] - \sum_{\tilde{e} \neq e} \ln [\tau^{\tilde{e}} [X_{\tilde{e}}]] \otimes \mathbb{I} [X_{V/\tilde{e}}]$$

and notice, that  $\tilde{\tau}$  does not depend on  $\tau^e$ .

The objective has then a representation as

$$O(\tau^G) = \langle \tilde{\tau}[X_V], \langle \tau^G \rangle [X_V | \emptyset] \rangle [\emptyset] - \langle \ln [\tau^e], \langle \tau^G \rangle [X_V | \emptyset] \rangle [\emptyset] + \ln [\langle \tau^G \rangle [\emptyset]]$$

Let us now differentiate all terms. With Lem. 31 we now get

$$\begin{aligned} \frac{\partial}{\partial \tau^e [Y_e]} \langle \tilde{\tau}[X_V], \langle \tau^G \rangle [X_V | \emptyset] \rangle [\emptyset] &= \left\langle \tilde{\tau}[X_V], \delta [Y_e, X_e], \frac{\langle \tau^G \rangle [X_e]}{\tau^e [X_e]}, \langle \tau^G \rangle [X_{V/e} | X_e] \right\rangle [Y_e, X_V] \\ &\quad - \langle \tilde{\tau}[X_V], \langle \tau^G \rangle [X_V | \emptyset] \rangle [\emptyset] \otimes \left\langle \frac{\langle \tau^G \rangle [Y_e]}{\tau^e [Y_e]} \right\rangle [Y_e]. \end{aligned}$$

Further we have

$$\begin{aligned} \frac{\partial}{\partial \tau^e [Y_e]} \langle \ln [\tau^e], \langle \tau^G \rangle [X_V | \emptyset] \rangle [\emptyset] &= \left\langle \ln [\tau^e [X_e]], \delta [Y_e, X_e], \frac{\langle \tau^G \rangle [X_e]}{\tau^e [X_e]}, \langle \tau^G \rangle [X_{V/e} | X_e] \right\rangle [Y_e, X_V] \\ &\quad - \langle \ln [\tau^e [X_e]], \langle \tau^G \rangle [X_V | \emptyset] \rangle [\emptyset] \otimes \left\langle \frac{\langle \tau^G \rangle [Y_e]}{\tau^e [Y_e]} \right\rangle [Y_e] \\ &\quad - \left\langle \frac{1}{\tau^e [X_e]}, \langle \tau^G \rangle [X_V | \emptyset] \right\rangle [\emptyset] \end{aligned}$$

and (see Proof of 30)

$$\frac{\partial}{\partial \tau^e [Y_e]} \ln [\langle \tau^G \rangle [\emptyset]] = \frac{\frac{\partial}{\partial \tau^e [Y_e]} \langle \tau^G \rangle [\emptyset]}{\langle \tau^G \rangle [\emptyset]} = \frac{\langle \tau^G \rangle [Y_e]}{\tau^e [Y_e]}.$$

Together, the first order condition

$$0 = \frac{\partial}{\partial \tau^e [Y_e]} O(\tau^G)$$

is equal to all  $y_e$  satisfying

$$\begin{aligned} 0 &= \frac{\langle \tau^G \rangle [Y_e = y_e]}{\tau^e [Y_e = y_e]} \left( \langle \tilde{\tau} [X_{V/e}, X_e = y_e], \langle \tau^G \rangle [X_{V/e} | X_e = y_e] \rangle [\emptyset] \right. \\ &\quad - \langle \tilde{\tau} [X_V], \langle \tau^G \rangle [X_V | \emptyset] \rangle [\emptyset] \\ &\quad - \langle \ln [\tau^e [X_e = y_e]], \langle \tau^G \rangle [X_{V/e} | X_e = y_e] \rangle [\emptyset] \\ &\quad \left. + \langle \ln [\tau^e [X_e]], \langle \tau^G \rangle [X_V | \emptyset] \rangle [\emptyset] \right). \end{aligned}$$

We notice, that by normalization

$$\langle \ln [\tau^e [X_e = y_e]], \langle \tau^G \rangle [X_{V/e} | X_e = y_e] \rangle [\emptyset] = \ln [\tau^e [X_e = y_e]]$$

and that the scalar

$$\lambda_1 = \langle \tilde{\tau} [X_V], \langle \tau^G \rangle [X_V | \emptyset] \rangle [\emptyset] - \langle \ln [\tau^e [X_e]], \langle \tau^G \rangle [X_V | \emptyset] \rangle [\emptyset]$$

is the constant for all  $y_e$ .

The first order condition is therefore equal to the existence of a  $\lambda_1 \in \mathbb{R}$  such that for all  $y_e$

$$\ln [\tau^e [X_e = y_e]] = \langle \tilde{\tau} [X_{V/e}, X_e = y_e], \langle \tau^G \rangle [X_{V/e} | X_e = y_e] \rangle [\emptyset] + \lambda_1.$$

The claim follows when applying the exponential on both sides and with the observation, that

$$\langle \tilde{\tau} [X_{V/e}, X_e = y_e], \langle \tau^G \rangle [X_{V/e} | X_e = y_e] \rangle [\emptyset] = \frac{\langle \{\tilde{\tau}\} \cup \{\tau^{\tilde{e}} : \tilde{e} \neq e\} \rangle [X_e = y_e]}{\langle \{\tau^{\tilde{e}} : \tilde{e} \neq e\} \rangle [X_e = y_e]}$$

and reparametrization of  $\lambda_1$  to

$$\lambda = \exp [\lambda_1].$$

□

## 6.8 Backward Map in Exponential Families

Let us now continue with the discussion of the backward map, which calculates to a mean parameter  $\mu [L]$  a canonical parameter  $\theta [L]$ , such that the corresponding member of the exponential family reproduced  $\mu [L]$ .

### 6.8.1 Variational Formulation

We now provide a variational characterization of the backward map.

**Theorem 38.** *Let there be a statistic  $\mathcal{S}$ , which is minimal with respect to a boolean base measure  $\nu$ . The map  $B^{(\mathcal{S}, \nu)} : (\mathcal{M}_{\mathcal{S}, \nu})^\circ \rightarrow \mathbb{R}^p$  defined as*

$$B^{(\mathcal{S}, \nu)}(\mu) = \operatorname{argmax}_{\theta \in \mathbb{R}^p} \langle \mu, \theta \rangle [\emptyset] - A^{(\mathcal{S}, \nu)}(\theta).$$

*is a backward mapping.*

*Proof.* We show the claim can be shown by the first order condition on the objective. It holds that

$$\begin{aligned} \frac{\partial}{\partial \theta [L]} A^{(\mathcal{S}, \nu)}(\theta) &= \frac{\partial}{\partial \theta [L]} \ln [\langle \exp [\langle \sigma^{\mathcal{S}}, \theta \rangle [X_{[d]}]] \rangle [\emptyset]] \\ &= \frac{\partial}{\partial \theta [L]} \frac{\langle \sigma^{\mathcal{S}} [L], \exp [\langle \sigma^{\mathcal{S}}, \theta \rangle [X_{[d]}]] \rangle [\emptyset]}{\langle \exp [\langle \sigma^{\mathcal{S}}, \theta \rangle [X_{[d]}]] \rangle [\emptyset]} \\ &= F^{(\mathcal{S}, \nu)}(\theta) [L] \end{aligned}$$

and thus

$$\frac{\partial}{\partial \theta [L]} \left( \langle \mu, \theta \rangle [\emptyset] - A^{(\mathcal{S}, \nu)}(\theta) \right) = \mu [L] - F^{(\mathcal{S}, \nu)}(\theta) [L].$$

The first order condition is therefore

$$\mu [L] = F^{(\mathcal{S}, \nu)}(\theta) [L]$$

and any  $\theta$  satisfies this condition exactly when  $\theta = B^{(\mathcal{S}, \nu)}(\mu)$  for a backward map. We conclude the proof by noticing, that since  $\mu [L]$  is in the interior  $(\mathcal{M}_{\mathcal{S}, \nu})^\circ$  we find by The. 24 a  $\theta$  such that the first order condition is met. □

### 6.8.2 Interpretation as a moment projection

Backward maps coincide with the Maximum Likelihood Estimation Problem  $(P_{\Gamma, \mathbb{P}^D}^M)$ , when we take the hypothesis to be exponential family  $\Gamma^{\mathcal{S}, \nu}$ . We show this as a more general statement for moment projections onto exponential families.

**Theorem 39.** *Let there be any exponential family, a mean parameter vector  $\mu^* \in \text{im}(F^{(\mathcal{S}, \nu)})$  and a backward map  $B^{(\mathcal{S}, \nu)}$ . Then  $\hat{\theta} = B^{(\mathcal{S}, \nu)}(\mu^*)$  is the canonical parameter to the solution of the moment projection Problem  $\mathbb{P}_{\Gamma, \mathbb{P}^*}^M$  of any  $\mathbb{P}^*$  with*

$$\langle \sigma^{\mathcal{S}}, \mathbb{P}^* \rangle [L] = \mu^*[L]$$

*onto the exponential family, if*

$$\mathbb{P}^{(\mathcal{S}, \hat{\theta}, \nu)} \in \text{argmax}_{\mathbb{P} \in \Gamma^{\mathcal{S}, \nu}} \mathbb{H}[\mathbb{P}^*, \mathbb{P}] .$$

*Proof.* We exploit the variational characterization of the backward map by The. 38, and first show that the objective coincides with the cross entropy between the distribution  $\mathbb{P}^*$  and the respective member of the exponential family. For any  $\mathbb{P}^*$  and  $\theta$  we have with Example 4

$$\mathbb{H}[\mathbb{P}^*, \mathbb{P}^{(\mathcal{S}, \theta, \nu)}] = \langle \mathbb{P}^*, \sigma^{\mathcal{S}}, \theta \rangle [\emptyset] - A^{(\mathcal{S}, \nu)}(\theta) .$$

We use that by assumption  $\langle \mathbb{P}^*, \sigma^{\mathcal{S}} \rangle [L] = \mu^*[L]$  and thus

$$\mathbb{H}[\mathbb{P}^*, \mathbb{P}^{(\mathcal{S}, \theta, \nu)}] = \langle \mu^*, \theta \rangle [\emptyset] - A^{(\mathcal{S}, \nu)}(\theta) .$$

This shows, that the backward map coincides with the moment projection onto  $\Gamma = \Gamma^{\mathcal{S}, \nu}$ .  $\square$

In particular, if we choose  $\mu[L] = \langle \sigma^{\mathcal{S}}[X_{[d]}, L], \mathbb{P}^D \rangle [L]$  for an empirical distribution  $\mathbb{P}^D$ , the backward map is a maximum likelihood estimator. This holds, since by Lem. 3 maximum likelihood estimation is a special instance of moment projection, when projecting  $\mathbb{P}^D$ . In that case, we choose  $\mu[L] = \langle \sigma^{\mathcal{S}}[X_{[d]}, L], \mathbb{P}^D \rangle [L]$ .

### 6.8.3 Approximation by alternating algorithms

While the forward map always has a representation in closed form by contraction of the probability tensor, the backward map in general fails to have a closed form representation. Computation of the Backward map can instead be performed by alternating algorithms, as we show here. We alternate through the coordinates of the statistics and adjust  $\theta[L = l]$  to a minimum of the likelihood, i.e. where for any  $l \in [p]$

$$0 = \frac{\partial}{\partial \theta[L = l]} \mathcal{L}_D(\mathbb{P}^{(\mathcal{S}, \theta, \nu)}) .$$

This condition is equal to the collection of moment matching equations

$$\langle \mathbb{P}^{(\mathcal{S}, \theta, \nu)}, \sigma^{\mathcal{S}} \rangle [L = l] = \langle \mathbb{P}^D, \sigma^{\mathcal{S}} \rangle [L = l] .$$

**Lemma 8.** *For any sufficient statistic  $\mathcal{S}$  a parameter vector  $\theta$  and an index  $l \in [p]$  we define*

$$\tau[X_{\mathcal{S}_l}] = \langle \{\beta^{\mathcal{S}}\} \cup \{\alpha^{\tilde{l}} : \tilde{l} \in [p], \tilde{l} \neq l\} \rangle [X_{\mathcal{S}_l}] .$$

*Then the moment matching condition for  $\mathcal{S}_l$  relative to  $\theta$  and  $\mu$  is satisfied for any  $\theta[L = l]$  with*

$$\langle \alpha^l, \text{Id}_{|\text{im}(\mathcal{S}_l)}, \tau[X_{\mathcal{S}_l}] \rangle [\emptyset] = \langle \alpha^l, \tau[X_{\mathcal{S}_l}] \rangle [\emptyset] \cdot \mu[L = l] .$$

*Proof.* We have

$$\mathbb{P}^{(\mathcal{S}, \theta, \nu)} = \frac{\langle \alpha^l, \tau \rangle [X_{[d]}]}{\langle \alpha^l, \tau \rangle [\emptyset]}$$

and

$$\langle \mathbb{P}^{(\mathcal{S}, \theta, \nu)}, \mathcal{S}_l \rangle [\emptyset] = \frac{\langle \alpha^l, \text{Id}_{|\text{im}(\mathcal{S}_l)}, \tau \rangle [X_{[d]}]}{\langle \alpha^l, \tau \rangle [\emptyset]} .$$

Here we used

$$\mathcal{S}_l = \langle \alpha^l, \text{Id}_{|\text{im}(\mathcal{S}_l)} \rangle [X_{[d]}]$$

and redundancies of copies of basis encodings. It follows that

$$\langle \mathbb{P}^{(\mathcal{S}, \theta, \nu)}, \mathcal{S}_l \rangle [\emptyset] = \langle \mathbb{P}^D, \mathcal{S}_l \rangle [\emptyset]$$

is equal to

$$\langle \alpha^l, \text{Id}_{|\text{im}(\mathcal{S}_l)}, \tau[X_{\mathcal{S}_l}] \rangle [\emptyset] = \langle \alpha^l, \tau[X_{\mathcal{S}_l}] \rangle [\emptyset] \cdot \mu[L = l] .$$

$\square$

The steps have to be alternated until sufficient convergence, since matching the moment to  $l$  by modifying  $\theta [L = l]$  will in general change other moments, which will have to be refit.

An alternating optimization is the coordinate descent of the negative likelihood, seen as a function of the coordinates of  $\theta$ , see Algorithm 5. Since the log likelihood is concave, the algorithm converges to a global minimum.

---

**Algorithm 5** Alternating Moment Matching for the Backward Map

---

**Require:** Empirical distribution  $\mathbb{P}^D$ , statistic  $\mathcal{S}$  and base measure  $\nu$

**Ensure:** Canonical parameter  $\theta [L]$ , such that  $\mathbb{P}^{(\mathcal{S}, \theta, \nu)}$  is the (approximative) moment projection of  $\mathbb{P}^D$  onto  $\Gamma^{\mathcal{S}, \nu}$

---

Set  $\theta [L] = 0 [L]$

Compute  $\mu_D [L] = \langle \mathbb{P}^D, \sigma^{\mathcal{S}} \rangle [L]$

**while** Stopping criterion is not met **do**

**for**  $l \in [p]$  **do**

    Compute

$$\tau^l [X_{S_l}] \leftarrow \left\langle \{\beta^{\mathcal{S}}\} \cup \{\alpha^{\tilde{l}} : \tilde{l} \in [p], \tilde{l} \neq l\} \right\rangle [X_{S_l}]$$

  Set  $\theta [L = l]$  to a solution of

$$\langle \alpha^l, \text{Id}|_{\text{im}(S_l)}, \tau^l \rangle [\emptyset] = \langle \alpha^l, \tau^l \rangle [\emptyset] \cdot \mu_D [L = l] .$$

**end for**

**end while** **return**  $\theta [L]$

---

In general, if  $\text{im}(S_l)$  contains more than two elements, there exists no closed form solutions. We will investigate the case of binary images, where there are closed form expressions, later in Sect. 12.3.

The computation of  $\tau^l$  in Algorithm 5 can be intractable and be replaced by an approximative procedure based on message passing schemes.

### 6.8.4 Second order Methods

The Hesse matrix of  $A^{(\mathcal{S}, \nu)}$  at  $\theta$  is the covariance of the features with respect to  $\mathbb{P}^{(\mathcal{S}, \theta, \nu)}$ , as we show next.

**Lemma 9.** *At any  $\theta \in \mathbb{R}^p$  we have*

$$\begin{aligned} \nabla_{\tilde{\theta}[\tilde{L}]} \nabla_{\tilde{\theta}[L]} |_{\theta} A^{(\mathcal{S}, \nu)}(\tilde{\theta}) &= \left\langle \mathbb{P}^{(\mathcal{S}, \theta, \nu)} [X_{[d]}], \sigma^{\mathcal{S}} [X_{[d]}, L], \sigma^{\mathcal{S}} [X_{[d]}, \tilde{L}] \right\rangle [\tilde{L}, L] \\ &\quad - \left\langle \mathbb{P}^{(\mathcal{S}, \theta, \nu)} [X_{[d]}], \sigma^{\mathcal{S}} [X_{[d]}, \tilde{L}] \right\rangle [\tilde{L}] \otimes \left\langle \mathbb{P}^{(\mathcal{S}, \theta, \nu)} [X_{[d]}], \sigma^{\mathcal{S}} [X_{[d]}, L] \right\rangle [L] . \end{aligned}$$

*Proof.* By Lem. 4 we have

$$\nabla_{\tilde{\theta}[L]} |_{\theta} A^{(\mathcal{S}, \nu)}(\tilde{\theta}) = \left\langle \mathbb{P}^{(\mathcal{S}, \theta, \nu)} [X_{[d]}], \sigma^{\mathcal{S}} [X_{[d]}, L] \right\rangle [L] .$$

It further holds

$$\nabla_{\tilde{\theta}[\tilde{L}]} |_{\theta} \mathbb{P}^{(\mathcal{S}, \theta, \nu)} [X_{[d]}] = \left\langle \mathbb{P}^{(\mathcal{S}, \theta, \nu)} [X_{[d]}], \sigma^{\mathcal{S}} [X_{[d]}, \tilde{L}] \right\rangle [X_{[d]}, \tilde{L}] - \left\langle \mathbb{P}^{(\mathcal{S}, \theta, \nu)} [X_{[d]}], \sigma^{\mathcal{S}} [X_{[d]}, \tilde{L}] \right\rangle [\tilde{L}] \otimes \mathbb{P}^{(\mathcal{S}, \theta, \nu)} [X_{[d]}] .$$

Combining both equations, we get the claim.  $\square$

With this characterization, we can perform second order optimization algorithms such as the Newton Method, see Algorithm 6 to solve the backward map.

## 6.9 Discussion

The forward and backward map have a correspondence with each other in terms of convex duality Rockafellar (1997). As such, the cumulant function  $A^{(\mathcal{S}, \nu)}$  (see Def. 24) and its conjugate dual  $(A^{(\mathcal{S}, \nu)})^*$  represent the forward and backward map by their gradient.

Further approximation schemes arise for the forward and backward map in exponential families arise from loopy message passing algorithms. For inference on Markov Network families, they are derived based on outer bounds on



**Algorithm 6** Newton Method for the Backward Map**Require:** Empirical distribution  $\mathbb{P}^D$ , statistic  $\mathcal{S}$  and base measure  $\nu$ **Ensure:** Canonical parameter  $\theta[L]$ , such that  $\mathbb{P}^{(\mathcal{S}, \theta, \nu)}$  is the (approximative) moment projection of  $\mathbb{P}^D$  onto  $\Gamma^{\mathcal{S}, \nu}$ Set  $\theta[L] = 0[L]$ Compute  $\mu_D[L] = \langle \mathbb{P}^D, \sigma^{\mathcal{S}} \rangle [L]$ **while** Convergence criterion not met **do**    Calculate  $\nabla_{\tilde{\theta}[\tilde{L}]} \nabla_{\tilde{\theta}[L]} |_{\theta A^{(\mathcal{S}, \nu)}(\tilde{\theta})}$  and  $\nabla_{\tilde{\theta}[L]} |_{\theta A^{(\mathcal{S}, \nu)}(\tilde{\theta})}$  as in Lem. 9.

Solve the linear equation

$$\left( \nabla_{\tilde{\theta}[\tilde{L}]} \nabla_{\tilde{\theta}[L]} |_{\theta A^{(\mathcal{S}, \nu)}(\tilde{\theta})} \right) \Delta[L] = \mu[L] - \nabla_{\tilde{\theta}[L]} |_{\theta A^{(\mathcal{S}, \nu)}(\tilde{\theta})}$$

Update the canonical parameter

$$\theta[L] \leftarrow \theta[L] - \Delta[L]$$

**end while**

the mean polytope in terms of the local consistency polytope and approximations on the entropy term by local terms (see Chapter 3 in Wainwright and Jordan (2008)). In particular, different approximations of the polytope and the entropies are made in the Bethe and Kickuchi method. We will discuss message passing schemes in Chapter 20 with a focus on exact computation of generic contractions by local contractions.

## 7 Propositional Logic

Propositional logics describes systems with  $d$  boolean variables, which are called atoms and denoted by  $X_k$  for  $k \in [d]$ . Indices  $x_k \in [2]$  to the atoms  $k \in [d]$  enumerate the  $2^d$  states of these systems, which are called worlds. In each world indexed by  $x_{[d]} = x_0, \dots, x_{d-1}$  the indices  $X_k$  encode whether the corresponding variable is True.

The epistemological commitments of propositional logics are whether the state is True or False reflected by the coordinate of the one-hot encoding being 1 or 0. Intuitively this describes, whether a specific world can be the state of a factored system. Propositional logic amounts to reason about boolean variables, which are categorical variables with 2 possible values. Such boolean tensors have already appeared as base measures in the representation of probability distributions in Chapter 5.

Before discussing the semantics and syntax of propositional formulas, we first investigate how Boolean can be represented by vectors in order to mechanize their processing based on contractions.

### 7.1 Encoding of Booleans

Booleans are variables valued by  $\{\text{False}, \text{True}\}$  and consist a basic data structure.

#### 7.1.1 Representation by coordinates

To represent Booleans by categorical variables  $X$  with two states we use the index interpretation function

$$I : [2] \rightarrow \{\text{False}, \text{True}\}$$

defined as

$$I(1) = \text{True} \quad \text{and} \quad I(0) = \text{False}.$$

In Def. 80 in Part III will define encodings of arbitrary sets based on index interpretation maps.

One motivation for this particular choice of the interpretation function  $I$  is the effective execution of the conjunction as we show in the next Lemma.

**Lemma 10.**  *$I$  is a homomorphism between the groups*

$$(\{0, 1\}, \cdot) \quad \text{and} \quad (\{\text{False}, \text{True}\}, \wedge).$$

*Proof.* It suffices to notice, that for arbitrary  $z_0, z_1 \in \{0, 1\}$  we have

$$I(z_0 \cdot z_1) = I(z_0) \wedge I(z_1).$$

□

Based on this homomorphism, contractions of boolean tensors, in which all variables are kept open, can be regarded as parallel calculations of the conjunction  $\wedge$  encoded by  $I$ . This homomorphism is further applied in type conversion in dynamically-typed languages (e.g. in python Foundation (2025)).

Operations like the negation fail to be linear and are only affine linear, since for  $z \in \{\text{False}, \text{True}\}$  we have

$$I^{-1}(\neg z) = 1 - I^{-1}(z). \quad (9)$$

Since any logical connective can be represented as a composition of conjunctions and negations, any logical connective corresponds with an affine linear function on the interpreted truth values. Direct applications of this insight to execute logical calculus will be discussed later in Sect. 17.6. For our purposes here, we would like to execute logical connective based on single contractions and avoid summations over them. This is why we call the negation representation as in (9) the affine representation problem, which we in the following want to resolve.

While in this work, we will always encode boolean states by  $I$ , other index interpretation functions could be chosen. For example, the interpretation

$$I_{\vee} : \{0, 1\} \rightarrow \{\text{False}, \text{True}\}$$

defined as

$$I_{\vee}(0) = \text{True} \quad \text{and} \quad I_{\vee}(1) = \text{False},$$

results is a homomorphism between the groups

$$(\{0, 1\}, \cdot) \quad \text{and} \quad (\{\text{False}, \text{True}\}, \vee).$$

While placing the disjunction  $\vee$  as the logical connective effectively executed by contractions, the negation will for arbitrary interpretations mapping onto  $\{0, 1\}$  remain the function

$$I_{\vee}^{-1}(\neg z) = 1 - I_{\vee}^{-1}(z).$$

Thus, the problem of affine linear operations cannot be resolved by a clever choice of an interpretation function with image in  $\{0, 1\}$ .

### 7.1.2 Representation by basis vectors

While contractions can just perform conjunctions, we need a representation trick to extend the contraction expressivity to arbitrary connectives and resolve the affine representation problem. To this end we now compose  $I$  with the one-hot encoding  $\epsilon$  and get an encoding

$$\epsilon \circ I^{-1} : \{\text{False}, \text{True}\} \rightarrow \{\epsilon_0[X], \epsilon_1[X]\},$$

where  $X$  is a categorical variable with  $m = 2$ . For any  $z \in \{\text{False}, \text{True}\}$  we have

$$\epsilon \circ I^{-1}(z) = \begin{bmatrix} I^{-1}(\neg z) \\ I^{-1}(z) \end{bmatrix}.$$

Performing the negation now amounts to switching the coordinates of the encoded vector, which can be performed by contraction with a transposition matrix

$$\beta^{\neg} [Y_{\neg}, X] = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix},$$

where in this notation we always understand the first variable  $X$  as the row index selector and the second variable  $Y_{\neg}$  as the column index selector. We then have

$$\epsilon \circ I^{-1}(\neg z)[Y_{\neg}] = \langle \beta^{\neg} [Y_{\neg}, X], \epsilon \circ I^{-1}(z)[X] \rangle [Y_{\neg}].$$

We therefore arrived at our aim to resolve the affine representation problem and have found a procedure to represent logical negations by a contraction, which is a linear operation. Besides negations, we will show in this chapter, that arbitrary logical formulas can be represented by contractions.

### 7.1.3 Coordinate and Basis Calculus

Our findings on the encoding of booleans hint towards more general schemes to encode information into boolean tensors, which will be explored in more detail in Chapter 16 and Chapter 17. When each coordinate in a boolean tensor represents one in  $\{0, 1\}$  interpreted boolean we call the scheme coordinate calculus. In basis calculus on the other hand, booleans are represented by elements of  $\{\epsilon_0[X], \epsilon_1[X]\}$ . In that scheme, there are pairs of two coordinates (building slice vectors of the tensors), which are restricted to be different from each other. This amounts to posing a global directionality constraint on the boolean tensor, as will be shown in The. 106.

## 7.2 Semantics of Propositional Formulas

We now choose a semantic centric approach to propositional logic, by defining formulas as boolean tensors. Then we investigate the corresponding syntax of formulas as specification of a tensor network decomposition of the basis encoding of formulas.

### 7.2.1 Formulas

Logics is especially useful in interpreting boolean tensors representing Propositional Knowledge Bases, based on connections with abstract human thinking. To make this more precise, we associate each such tensor is associated with a formula  $f$  being a composition of the atomic variables with logical connectives as we proof next.

**Definition 43.** A propositional formula  $f [X_{[d]}]$  depending on  $d$  atoms  $X_k$  is a boolean-valued tensor

$$f [X_{[d]}] : \bigotimes_{k \in [d]} [2] \rightarrow \{0, 1\} \subset \mathbb{R}.$$

We call a state  $x_{[d]} \in \bigotimes_{k \in [d]} [2]$  a model of a propositional formula  $f$ , if

$$f [X_{[d]} = x_{[d]}] = 1.$$

If there is a model to a propositional formula, we say the formula is satisfiable.

The propositional formulas coincide therefore with the boolean tensors (see Def. 4).

Since propositional formulas are binary valued tensors, the generic decomposition of Lem. 25 simplifies to

$$f [X_{[d]}] = \sum_{x_0, \dots, x_{d-1} \in \bigotimes_{k \in [d]} [2]} f [X_0 = x_0, \dots, X_{d-1} = x_{d-1}] \cdot \epsilon_{x_{[d]}} [X_{[d]}] \quad (10)$$

$$= \sum_{x_0, \dots, x_{d-1} \in \bigotimes_{k \in [d]} [2] : f [X_0 = x_0, \dots, X_{d-1} = x_{d-1}] = 1} \epsilon_{x_0, \dots, x_{d-1}} [X_{[d]}]. \quad (11)$$

Thus, any propositional formula is the sum over the one-hot encodings of its models. This is equal to the encoding of the set of models, which will be introduced in Chapter 17 (see Def. 80).

We depict this decomposition in the diagrammatic notation by

$$\begin{array}{c} \boxed{f} \\ \hline X_0 \mid X_1 \mid \dots \mid X_{d-1} \end{array} = \sum_{\substack{x_0, \dots, x_{d-1} \in \bigotimes_{k \in [d]} [2] \\ f(x_0, \dots, x_{d-1}) = 1}} \begin{array}{c} \boxed{\epsilon_{x_0}} \\ \hline \downarrow X_0 \end{array} \dots \begin{array}{c} \boxed{\epsilon_{x_{d-1}}} \\ \hline \downarrow X_{d-1} \end{array}$$

We here chose a semantic approach to propositional logic in contrary to the standard syntactical approach. Instead of defining formulas by connectives acting on atomic formulas, we define them here as binary valued functions of the states of a factored system. They are interpreted by marking possible states as models, given the knowledge of  $f$ . The syntactical side will then be introduced later by studying decompositions of formulas.

### 7.2.2 Basis encoding of formulas

There are two ways to represent formulas by tensors. One way is to understand  $[2]$  as subset of  $\mathbb{R}$  and interpreting the formula directly as a tensor (as in Def. 43). Another way is to understand  $[2]$  as the possible values of a categorical variable. Following this second perspective, formulas are maps between factored systems, where the image system is the factored systems of atoms and the target system the atomic system defined by a variable  $Y_f$  representing the formula satisfaction. We can then build the basis encoding (Def. 14) of that map to represent the formula (see Figure 17).

Given a factored system with  $d$  atoms  $X_{[d]}$  and a propositional formula  $f$ , the basis encoding of  $f$  (see Def. 14) is the tensor

$$\beta^f [Y_f, X_{[d]}] \in \left( \bigotimes_{k \in [d]} \mathbb{R}^2 \right) \otimes \mathbb{R}^2$$

decomposable as

$$\beta^f [Y_f, X_{[d]}] = \sum_{x_{[d]} \in \bigotimes_{k \in [d]} [2]} \epsilon_{x_{[d]}} [X_{[d]}] \otimes \epsilon_{f[X_{[d]}=x_{[d]}]} [Y_f]. \quad (12)$$

We can build basis encodings more generally of any tensors, where we identify the image of the tensor with the states of a categorical variable. Exactly for propositional formulas, this construction will lead to Boolean image variables.

**Lemma 11.** *For any formula  $f$  we have*

$$\beta^f [Y_f, X_{[d]}] = f [X_{[d]}] \otimes \epsilon_1 [Y_f] + \neg f [X_{[d]}] \otimes \epsilon_0 [Y_f] .$$

*In particular*

$$f [X_{[d]}] = \langle \beta^f [Y_f, X_{[d]}], \epsilon_1 [Y_f] \rangle [X_{[d]}] .$$

*Proof.* We can decompose basis encodings of formulas into the sum (see Figure 17)

$$\beta^f [Y_f, X_{[d]}] = \epsilon_0 [Y_f] \otimes \left( \sum_{x_{[d]} : f[x_{[d]}]=0} \epsilon_{x_{[d]}} [X_{[d]}] \right) \quad (13)$$

$$+ \epsilon_1 [Y_f] \otimes \left( \sum_{x_{[d]} : f[x_{[d]}]=1} \epsilon_{x_{[d]}} [X_{[d]}] \right) \quad (14)$$

where the second term sums up the models of  $f$  and the first one the models of  $\neg f$ .  $\square$

Compared with the direct interpretation of a formula as a tensor and the decomposition into models in Equation 10, we notice that the basis encoding also represents encoding of worlds where the formula is not satisfied. This representation is required to represent arbitrary propositional formulas by contracted tensor networks of its components, as will be investigated in the following sections.

The basis encoding  $\beta^f$  has slices

$$\langle \beta^f, \epsilon_{x_{[d]}} \rangle [Y_f] \beta^f [X_{[d]} = x_{[d]}, Y_f] = \begin{cases} \epsilon_1 [Y_f] & \text{if the world } x_{[d]} \text{ is a model of } f \\ \epsilon_0 [Y_f] & \text{else .} \end{cases}$$

The contractions of the basis encoding therefore calculate whether an assignment of atoms is a model of the formula, using basis calculus (see The. 107).

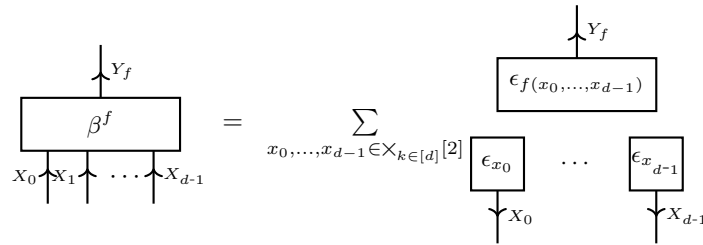


Figure 17: basis encoding of a propositional formula. The encoding is a sum of the one hot encodings of all states of the factored system in a tensor product with basis vectors, which encode whether the state is a model of the formula. The tensor is directed, since any contraction with an encoded state results in the basis vector evaluating the formula, which we called basis calculus.

### 7.3 Syntax of Propositional Formulas

basis encodings of propositional formulas are especially useful when representing function compositions by the representation of their components (see The. 108). In propositional logics, the syntax of defining propositional formulas is oriented on compositions of formulas by connectives. We in this section investigate the decomposition schemes of basis encodings into tensor networks of component encodings for binary tensors following propositional logic syntax.

### 7.3.1 Atomic Formulas

We call atomic formulas the most granular formulas, which are not splitted into compositions of other formulas. Our syntactic decomposition of propositional formulas will then investigate, how any propositional formula can be represented by these.

**Definition 44.** The tensors  $f_k [X_{[d]}]$  defined for  $x_{[d]} \in \times_{k \in [d]} [2]$  as

$$f_k [X_{[d]} = x_{[d]}] = x_k$$

are called atomic formulas.

Atomic formulas and their basis encodings have an especially compelling representation.

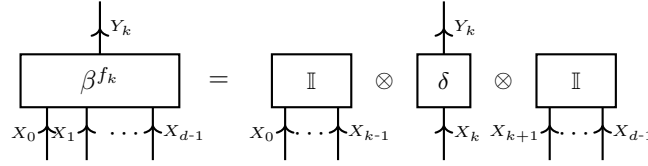
**Theorem 40.** Any atomic formula  $f_k [X_{[d]}]$  is represented as

$$f_k [X_{[d]} = x_{[d]}] = \langle \epsilon_1 [X_k] \rangle [X_{[d]}] = \epsilon_1 [X_k] \otimes \mathbb{I} [X_{[d]}/\{k\}] .$$

The basis encoding of any atomic formula  $X_k [X_{[d]}]$  has a tensor decomposition by

$$\beta^{X_k} [Y_d, X_{[d]}] = \langle \delta [X_k, Y_k] \rangle [X_{[d]}] = \delta [X_k, Y_k] \otimes \mathbb{I} [X_{[d]}/\{k\}] .$$

The decomposition is depicted in a network diagram as



*Proof.* We have by definition

$$\begin{aligned} \beta^{X_k} [Y_k, X_{[d]}] &= \sum_{x_0, \dots, x_{d-1} \in \times_{k \in [d]} [2]} \epsilon_{x_0, \dots, x_{d-1}} [X_{[d]}] \otimes \epsilon_{f_k [X_0=x_0, \dots, X_{d-1}=x_{d-1}]} [Y_k] \\ &= (\epsilon_{0,0} [X_k, Y_k] + \epsilon_{1,1} [X_k, Y_k]) \otimes \mathbb{I} [X_l : l \neq k] \\ &= \langle \delta [X_k, Y_k] \rangle [X_{[d]}, Y_k] . \end{aligned}$$

□

### 7.3.2 Syntactical combination of formulas

Propositional formulas are elements of tensor spaces with  $d$  axis. The number of coordinates thus grows exponentially with the number of atoms, which is

$$\dim \left( \bigotimes_{k \in [d]} \mathbb{R}^2 \right) = 2^d .$$

When the number of atoms is large, the naive representation of formula tensors will be thus intractable. In contrast, typical logical formulas appearing in practical knowledge bases are sparse in the sense that they have short representations in a logical syntax. Motivated by this consideration we now discuss propositional syntax and investigate the sparse decomposition of formula tensors along their formula structure to avoid the curse of dimensionality.

In logical syntax formulas are described by atomic formulas recursively connected via connectives. We show, that representations of logical connectives can be represented by feasible tensor cores  $\beta^\circ$  contracted along a tensor network. Let us first provide in Example 6 unary ( $d = 1$ ) and binary ( $d = 2$ ) connectives.

**Example 6.** We use the following connectives:

- negation  $\neg : [2] \rightarrow [2]$  by the vector

$$\neg[Y_f] = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$$

- *conjunctions*  $\wedge : [2] \times [2] \rightarrow [2]$

$$\wedge[Y_f, Y_h] = \begin{bmatrix} 0 & 0 \\ 0 & 1 \end{bmatrix}$$

- *disjunctions*  $\vee : [2] \times [2] \rightarrow [2]$

$$\vee[Y_f, Y_h] = \begin{bmatrix} 0 & 1 \\ 1 & 1 \end{bmatrix}$$

- *exact disjunction*  $\oplus : [2] \times [2] \rightarrow [2]$

$$\oplus[Y_f, Y_h] = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}$$

- *implications*  $\Rightarrow : [2] \times [2] \rightarrow [2]$

$$\Rightarrow[Y_f, Y_h] = \begin{bmatrix} 1 & 1 \\ 0 & 1 \end{bmatrix}$$

- *biimplication*  $\Leftrightarrow : [2] \times [2] \rightarrow [2]$

$$\Leftrightarrow[Y_f, Y_h] = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

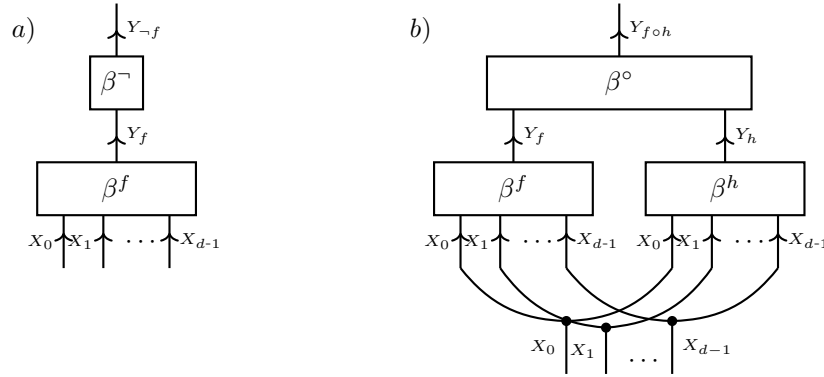


Figure 18: a) Basis encoding of a negated formula  $f$  as a tensor network of the encoded formula and the encoded connective  $\neg$ . b) Basis encoding of a composition of formulas  $f, h$  by a connective  $\circ \in \{\wedge, \vee, \oplus, \Rightarrow, \Leftrightarrow\}$ . The encoding is a contraction of encodings to  $f, h$  and  $\circ$ .

We now show how formulas consisting of connectives acting on other formulas can be represented by basis calculus. Let there be formulas  $f$  and  $h$  depending on categorical variables  $X_{[d]}$  and a binary connective

$$\circ : [2] \times [2] \rightarrow [2] .$$

Then we can show as a special case of the next theorem, that (see Figure 18)

$$\beta^{f \circ h} [X_{[d]}, X_{f \circ h}] = \langle \beta^\circ [Y_f, Y_h, X_{f \circ h}], \beta^f [X_{[d]}, Y_f], \beta^h [X_{[d]}, Y_h] \rangle [X_{[d]}, X_{f \circ h}] .$$

For any unary connective  $\circ : [2] \rightarrow [2]$  we have

$$\beta^{\circ f} [X_{[d]}, X_{\circ f}] = \langle \beta^\circ [Y_f, X_{\circ f}], \beta^f [X_{[d]}, Y_f] \rangle [X_{[d]}, X_{\circ f}] .$$

Let us now generalize this observation to arbitrary arity of connectives and provide a proof of its correctness.

**Theorem 41** (Composition of Formulas). *Let there be a formula  $f [X_{[d]}]$ , which has a syntactical decomposition into connectives  $\{\circ_l [Y_{\mathcal{V}^l}] : l \in [p]\}$  taking their inputs by variables  $Y_{\mathcal{V}^l} \subset Y_{\mathcal{V}}$  and output by a variable  $Y_{\circ_l}$ . We here denote by  $\mathcal{F}$  the set of sub-formulas and use a boolean variable  $Y_h$  for each  $h \in \mathcal{F}$ . In particular, we denote for each atom in  $\mathcal{F}$  the corresponding boolean variable by  $Y_k$ . It then holds*

$$\beta^f [Y_f, X_{[d]}] = \langle \{\beta^{\circ_l} [Y_{\circ_l}, Y_{\mathcal{V}^l}] : l \in [p]\} \cup \{\delta [Y_k, X_k] : k \in [d]\} \rangle [Y_f, X_{[d]}] .$$

*Proof.* When a variable in  $Y_{\mathcal{F}}$  appears multiple times as input to connectives, we replace it by a set of copies (which won't change the contraction, since all tensors are binary and The. 135 can be applied). This follows from an iterative application of The. 108 to be shown in Chapter 17.  $\square$

**Remark 3** ( $d$ -ary connectives such as  $\wedge$  and  $\vee$ ). *Since the decomposition of basis encoding can be applied to generic function compositions (see The. 108), we can also allow for  $d$ -ary connectives*

$$\circ : \bigtimes_{k \in [d]} [2] \rightarrow [2].$$

*The connectives  $\wedge$  and  $\vee$  satisfy associativity and have thus straightforward generalizations to the  $d$ -ary case. This is because associativity can be exploited to represent the basis encoding by any tree-structured composition of binary  $\wedge$  and  $\vee$  connectives.*

Propositional syntax consists in the application of connectives on atomic formulas, and recursively on the results of such constructions. When passed towards connective cores, atomic formula tensors act trivial on the legs and just identify the corresponding atomic formula index  $x_{X_k}$  with  $x_k$ . This is due to the fact, that contractions with the trivial tensor  $\mathbb{I}$  leaves any tensor invariant, and the contraction with the elementary matrix  $\delta$  identifies indices with each other. We can thus safely ignore the atomic formula tensors appearing in the decomposition of formula tensors to non-atomic formulas. An example of such a decomposition is depicted in Figure 19.

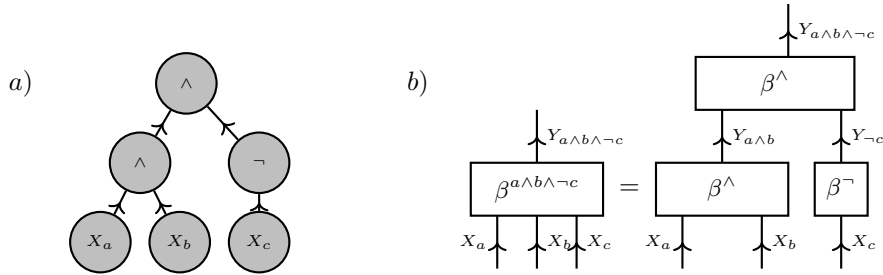


Figure 19: Decomposition of the formula tensor to  $f = a \wedge b \wedge \neg c$  into unary (matrix) and binary (third order tensor) cores. a) Visualization of  $f$  as a graph. b) Tensor Network decomposition of  $f$ . We can make use of the invariance of a Hadamard product with a constant tensor  $\mathbb{I}$  and thus not draw axis to atoms not affected by a formula.

**Remark 4** (Tensor Network Decomposition of Formulas). *The decomposition of the propositional into a tensor network is a hierarchical decomposition of the formula tensor, which we will describe in more detail in Sect. 17.5. Of special interest are tree hypergraphs, where the format is called Hierarchical Tucker. At each decomposition of a formula into sub-formulas, two subspaces spanned by the respective atomic spaces are selected.*

### 7.3.3 Syntactical decomposition of formulas

We have seen how the decomposition of complex formulas into connectives acting on the component formulas can be exploited to find effective representations of the semantics by tensor networks. Here the question arises here, how to perform such decompositions in case of a missing syntactical representation of a formula. By Def. 43 any binary tensor is a formula. We show in the following, how we can find a syntactic specification of a formula given its tensor.

**Definition 45** (Terms and Clauses). *Given two disjoint subsets  $\mathcal{V}^0$  and  $\mathcal{V}^1$  of  $[d]$ , the corresponding term is the formula defined on the indices  $x_{[d]} \in \bigtimes_{k \in [d]} [2]$  by*

$$Z_{\mathcal{V}^0, \mathcal{V}^1}^{\wedge} [X_{[d]}] = \left( \bigwedge_{k \in \mathcal{V}^0} \neg f_k \right) \wedge \left( \bigwedge_{k \in \mathcal{V}^1} f_k \right)$$

*and the corresponding clause is the formula defined on the indices  $x_0, \dots, x_{d-1} \in \bigtimes_{k \in [d]} [2]$  by*

$$Z_{\mathcal{V}^0, \mathcal{V}^1}^{\vee} [X_{[d]}] = \left( \bigvee_{k \in \mathcal{V}^0} f_k \right) \vee \left( \bigvee_{k \in \mathcal{V}^1} \neg f_k \right),$$

*where by  $\bigwedge_{k \in \mathcal{V}}$  and  $\bigvee_{k \in \mathcal{V}}$  we refer to the  $n$ -ary connectives  $\wedge$  and  $\vee$ . We call the term a minterm and the clause a maxterm, if  $\mathcal{V}^0 \cup \mathcal{V}^1 = [d]$ .*

Terms and Clauses have for any index tuple  $x_{[d]}$  a polynomial representation by

$$Z_{\mathcal{V}^0, \mathcal{V}^1}^\wedge [X_{[d]} = x_{[d]}] = \left( \prod_{k \in \mathcal{V}^0} (1 - x_k) \right) \left( \prod_{k \in \mathcal{V}^1} x_k \right)$$

and

$$Z_{\mathcal{V}^0, \mathcal{V}^1}^\vee [X_{[d]} = x_{[d]}] = 1 - \left( \prod_{k \in \mathcal{V}^0} (1 - x_k) \right) \left( \prod_{k \in \mathcal{V}^1} x_k \right).$$

**Lemma 12.** *Terms are contractions of one-hot encodings, that is for any disjoint subsets  $\mathcal{V}^0, \mathcal{V}^1 \subset [d]$  we have*

$$Z_{\mathcal{V}^0, \mathcal{V}^1}^\wedge [X_{[d]}] = \langle \epsilon_{\{x_k=0:k \in \mathcal{V}^0\} \cup \{x_k=1:k \in \mathcal{V}^1\}} \rangle [X_{[d]}].$$

*Clauses are substractions of one-hot encodings from the trivial tensor, that is for any disjoint subsets  $\mathcal{V}^0, \mathcal{V}^1 \subset [d]$  we have*

$$Z_{\mathcal{V}^0, \mathcal{V}^1}^\vee [X_{[d]}] = \mathbb{I} [X_{[d]}] - \langle \epsilon_{\{x_k=0:k \in \mathcal{V}^0\} \cup \{x_k=1:k \in \mathcal{V}^1\}} \rangle [X_{[d]}].$$

The reference of the formulas in the case  $\mathcal{V}^0 \dot{\cup} \mathcal{V}^1 = [d]$  as minterms and maxterms is due to the fact, that minterms are formulas with unique models and maxterms are formulas with a unique world not satisfying the formula. We use this insight and enumerate maxterms and minterms by the index  $x \in \times_{k \in [d]} [2]$  of the unique world where the minterm is satisfied, respectively the maxterm is not satisfied. For any  $\mathcal{V}^0 \dot{\cup} \mathcal{V}^1 = [d]$  we take the index tuple  $x_0, \dots, x_{d-1}$  where  $x_k = 0$  if  $k \in \mathcal{V}^0$  and  $x_k = 1$  if  $k \in \mathcal{V}^1$  and define

$$Z_{x_0, \dots, x_{d-1}}^\vee = Z_{\mathcal{V}^0, \mathcal{V}^1}^\vee \quad \text{and} \quad Z_{x_0, \dots, x_{d-1}}^\wedge = Z_{\mathcal{V}^0, \mathcal{V}^1}^\wedge.$$

**Corollary 3.** *Minterms are basis elements of the tensor space, that is for any  $x_{[d]} \in \times_{k \in [d]} [2]$  we have*

$$Z_{x_{[d]}}^\wedge = \epsilon_{x_{[d]}} [X_{[d]}]$$

*Maxterms are substraction of basis elements from the trivial tensor, that is for any  $x_{[d]} \in \times_{k \in [d]} [2]$  we have*

$$Z_{x_{[d]}}^\vee = \mathbb{I} [X_{[d]}] - \epsilon_{x_{[d]}} [X_{[d]}].$$

*Proof.* Follows from Lem. 12, since when  $\mathcal{V}^0 \cup \mathcal{V}^1 = [d]$  the contraction of the one-hot encodings coincides with the one-hot encoding of a fully specified state.  $\square$

Based on this insight, we can decompose any propositional formula into a conjunction of maxterms or a disjunction of minterms as we show next.

**Theorem 42.** *For any boolean tensor  $\tau [X_{[d]}] \in \otimes_{k \in [d]} \mathbb{R}^2$  with leg-dimensions two we have*

$$\tau [X_{[d]}] = \left( \bigvee_{x_{[d]} : \tau [X_{[d]} = x_{[d]}] = 1} Z_{\{k:x_k=0\}, \{k:x_k=0\}}^\wedge \right) [X_{[d]}]$$

and

$$\tau [X_{[d]}] = \left( \bigwedge_{x_{[d]} : \tau [X_{[d]} = x_{[d]}] = 0} Z_{\{k:x_k=0\}, \{k:x_k=0\}}^\vee \right) [X_{[d]}].$$

*Proof.* To show the representation by minterms we use the decomposition

$$\tau [X_{[d]}] = \sum_{x_{[d]} : \tau [X_{[d]} = x_{[d]}] = 1} \epsilon_{x_{[d]}} [X_{[d]}]$$

and notice that each term in the disjunction modifies the formula by adding respective world  $x_{[d]}$  to the models of the formula. To show the representation by maxterms we use the decomposition

$$\tau [X_{[d]}] = \mathbb{I} [X_{[d]}] - \sum_{x_{[d]} : \tau [X_{[d]} = x_{[d]}] = 0} \epsilon_{x_{[d]}} [X_{[d]}]$$

and notice that each term in the conjunction modifies the formula by removing the respective world  $x_{[d]}$  from the models of the formula. Thus, both decompositions are propositional formulas with the same set of models as the formula  $\tau$  and are thus identical to  $\tau$ .  $\square$



The decompositions found in The. 42 are also called canonical normal forms to propositional formulas  $\tau [X_{[d]}]$ .

**Remark 5** (Efficient Representation in Propositional Syntax). *The decomposition in The. 42 is a basis CP decomposition of the binary tensor and will further be investigated in Chapter 18. The formulas constructed in the proof of The. 42 are however just one possibility to represent a formula tensor in propositional syntax. Typically there are much sparser representations for many formula tensors, in the sense that less connectives and atomic symbols are required. Having such a sparser syntactical description of a propositional formula can be exploited to find a shorter conjunctive normal form of the formula and construct a sparse polynomial based on similar ideas as in The. 42. We will provide such constructions in Chapter 18, where we show that dropping the demand of directionality and investigating binary CP Decompositions will improve the sparsity of the polynomial formula representation.*

## 7.4 Outlook

While we in this chapter investigated representation schemes for single propositional formulas, we will further study the representation of knowledge bases consisting in multiple formulas in Sect. 11.2. Further, we will build hybrid models bridging the concepts of probability distributions and propositional logics in Sect. 11.3. Propositional formulas will therein serve as features and base measures for exponential families.

## 8 Logical Inference

We approach logical inference by defining probability distributions based on propositional formulas and then apply the methodology introduced in the more generic situation of probabilistic inference. Logical approaches pay here special attention to situations of certainty, where a state of a variable has probability 1. In this situation, we say that the corresponding formula is entailed.

We start the discussion with the derivation of contraction criteria for logical entailment. We interpret formulas by distributions and extend logical entailment towards probabilistic reasoning.

### 8.1 Entailment in Propositional Logics

Entailment is the central consequence relation among logical formulas. Let us define this relation first based on the models of a knowledge base and a test formula.

**Definition 46** (Entailment of propositional formulas). *Given two propositional formulas  $\mathcal{KB}$  and  $f$  we say that  $\mathcal{KB}$  entails  $f$ , denoted by  $\mathcal{KB} \models f$ , if any model of  $\mathcal{KB}$  is also a model of  $f$ , that is*

$$\forall x_{[d]} \in \bigtimes_{k \in [d]} [2] : (\mathcal{KB} [X_{[d]} = x_{[d]}] = 1) \Rightarrow (f [X_{[d]} = x_{[d]}] = 1).$$

*If  $\mathcal{KB} \models \neg f$  holds, we say that  $\mathcal{KB}$  contradicts  $f$ .*

To use the tensor network formalism for the decision of entailment, we will in the following develop three equivalent criteria for entailment.

#### 8.1.1 Deciding Entailment by contractions

First of all, we can decide entailment based on vanishing contractions with the negated test formula.

**Theorem 43** (Contraction Criterion of Entailment). *We have  $\mathcal{KB} \models f$  if and only if*

$$\langle \mathcal{KB}, \neg f \rangle [\emptyset] = 0.$$

*Proof.* " $\Leftarrow$ ": If for a  $x_{[d]} \in \bigtimes_{k \in [d]} [2]$  we have  $\mathcal{KB} [X_{[d]} = x_{[d]}] = 1$  but not  $(f [X_{[d]} = x_{[d]}] = 1)$ , we would have  $(\neg f [X_{[d]} = x_{[d]}] = 1)$  and

$$\langle \mathcal{KB}, \neg f \rangle [\emptyset] = \sum_{x_{[d]} \in \bigtimes_{k \in [d]} [2]} \mathcal{KB} [X_{[d]} = x_{[d]}] \cdot f [X_{[d]} = x_{[d]}] > 1.$$

Thus, whenever the contraction vanishes, we have

$$\forall x_{[d]} \in \bigtimes_{k \in [d]} [2] : (\mathcal{KB} [X_{[d]} = x_{[d]}] = 1) \Rightarrow (f [X_{[d]} = x_{[d]}] = 1).$$

" $\Rightarrow$ ": Conversely, if the contraction  $\langle \mathcal{KB}, \neg f \rangle [\emptyset]$  does not vanish, we would find  $x_{[d]} \in \times_{k \in [d]} [2]$  with  $\mathcal{KB} [X_{[d]} = x_{[d]}] = 1$  and  $\neg f [X_{[d]} = x_{[d]}] = 1$ , therefore  $f [X_{[d]} = x_{[d]}] = 0$ . It follows that  $\mathcal{KB} \models f$  does not hold.  $\square$

The contraction criterion can be extended to the decision of contradiction as well, since  $\mathcal{KB} \models \neg f$  is equivalent to  $\langle \mathcal{KB}, f \rangle [\emptyset] = 0$ . Therefore, entailment and contradiction can be decided simultaneously by a single contraction, as we state next.

**Theorem 44.** *Given propositional formulas  $\mathcal{KB}$  and  $f$  we build*

$$\tau [Y_f] = \langle \mathcal{KB} [X_{[d]}], f [X_{[d]}, Y_f] \rangle [Y_f] .$$

*Then  $\mathcal{KB} \models f$  is equivalent to  $\tau [Y_f = 0] = 0$ , and  $\mathcal{KB} \models \neg f$  is equivalent to  $\tau [Y_f = 1] = 0$  .*

*Proof.* This follows from The. 43 using that

$$\langle \mathcal{KB}, \neg f \rangle [\emptyset] = \tau [Y_f = 0]$$

and

$$\langle \mathcal{KB}, f \rangle [\emptyset] = \tau [Y_f = 1] .$$

$\square$

### 8.1.2 Deciding Entailment by partial ordering

Logical entailment can be understood by subset relations of the models of the respective formulas. This perspective can be applied with subset encodings in Chapter 17. The subset relation corresponds with partial ordering of its encoded tensors, as will be shown in The. 104. For two propositional formulas, we denote to this end  $f \prec h$  (see Def. 81), if and only if for all  $x_{[d]} \in \times_{k \in [d]} [2]$

$$f [X_{[d]} = x_{[d]}] \leq h [x_{[d]}] .$$

**Theorem 45** (Partial Ordering Criterion of Entailment). *We have  $\mathcal{KB} \models f$  if and only if  $\mathcal{KB} [X_{[d]}] \prec f [X_{[d]}]$ .*

*Proof.* Since both  $\mathcal{KB}$  and  $f$  are boolean tensors, we have for any  $x_{[d]} \in \times_{k \in [d]} [2]$  that

$$\mathcal{KB} [X_{[d]} = x_{[d]}], f [X_{[d]} = x_{[d]}] \in \{0, 1\} .$$

Thus,

$$\forall x_{[d]} \in \times_{k \in [d]} [2] : \mathcal{KB} [X_{[d]} = x_{[d]}] \leq f [X_{[d]} = x_{[d]}]$$

is equivalent to

$$\forall x_{[d]} \in \times_{k \in [d]} [2] : (\mathcal{KB} [X_{[d]} = x_{[d]}] = 1) \Rightarrow (f [X_{[d]} = x_{[d]}] = 1) .$$

This states that  $\mathcal{KB} [X_{[d]}] \prec f [X_{[d]}]$  is equivalent to  $\mathcal{KB} \models f$ .  $\square$

### 8.1.3 Redundancy of entailed formulas

Another interpretation of entailment is by redundancy of a formula in a Knowledge Base. This is especially interesting for the sparse representation of Knowledge Bases.

**Theorem 46** (Redundancy Criterion of Entailment). *If and only if  $\mathcal{KB} \models f$  we have*

$$\mathcal{KB} [X_{[d]}] = \langle \mathcal{KB}, f \rangle [X_{[d]}] .$$

*Proof.* For any formula  $f$  we have

$$\mathbb{I} [X_{[d]}] = f [X_{[d]}] + \neg f [X_{[d]}]$$

and thus

$$\begin{aligned} \mathcal{KB} [X_{[d]}] &= \langle \mathcal{KB} [X_{[d]}], \mathbb{I} [X_{[d]}] \rangle [X_{[d]}] \\ &= \langle \mathcal{KB} [X_{[d]}], f [X_{[d]}] \rangle [X_{[d]}] + \langle \mathcal{KB} [X_{[d]}], \neg f [X_{[d]}] \rangle [X_{[d]}] . \end{aligned}$$

Now, by The. 43 we have  $\mathcal{KB} \models f$ , if and only if  $\langle \mathcal{KB} [X_{[d]}], \neg f [X_{[d]}] \rangle [X_{[d]}] = 0$ , which is thus equal to

$$\mathcal{KB} [X_{[d]}] = \langle \mathcal{KB} [X_{[d]}], f [X_{[d]}] \rangle [X_{[d]}] .$$

$\square$

		$\mathcal{KB} \models f$	
		False	True
$\mathcal{KB} \models \neg f$	False	"Contingent"	"Entailed"
	True	"Contracticted"	"Inconsistent"

Figure 20: Table of possible logical relations between a knowledge base  $\mathcal{KB}$  and  $f$ , based on whether the knowledge base entails the formula ( $\mathcal{KB} \models f$ ) and its negation ( $\mathcal{KB} \models \neg f$ ).

#### 8.1.4 Contraction Knowledge Base

We exploit the contraction and redundancy criteria of entailment to sketch an implementation of a propositional Knowledge Base in Algorithm 7. Here the function  $\text{ASK}(f)$  returns, whether a formula  $f$  is entailed or contradicted by a Knowledge Base. If the formula is neither entailed or contradicted, we say it is contingent. If it is both, we have  $\mathcal{KB}[X_{[d]}] = 0$  and thus an inconsistent Knowledge Base. Exploiting The. 44 we decide these situations based on a single contraction.

The function  $\text{TELL}(f)$  incorporates an additional formula  $f$  into a Knowledge Base  $\mathcal{KB}$ . Here we exploit The. 46 and do not add a formula, which is entailed in order to maintain a sparse representation. The function further refuses to add a formula, which would make the Knowledge Base inconsistent (returns Refused) and only changes the Knowledge Base in case of a contingent formula (returns Added).

---

#### Algorithm 7 Contraction Knowledge Base with operations ASK and TELL

---

##### ASK( $\mathcal{KB}, f$ )

**Require:** Knowledge base  $\mathcal{KB}$ , query formula  $f$

**Ensure:** Decision which relation between  $\mathcal{KB}$  and  $f$  holds (see Figure 20)

---

```

 $\tau[Y_f] \leftarrow \langle \{h[X_{[d]}] : h \in \mathcal{KB}\}, \beta^f[Y_f, X_{[d]}] \rangle[Y_f]$ 
if  $\tau[Y_f = 0] = 0$  and  $\tau[Y_f = 1] = 0$  then
  return "Inconsistent"
else if  $\tau[Y_f = 0] = 0$  then
  return "Entailed"
else if  $\tau[Y_f = 1] = 0$  then
  return "Contradicted"
else
  return "Contingent"
end if

```

---

##### TELL( $\mathcal{KB}, f$ )

**Require:** Knowledge base  $\mathcal{KB}$ , query formula  $f$

**Ensure:** Decision whether a formula is added to the knowledge base ("Added"), or an exception in ("Inconsistent", "Redundant" or "Refused") is raised.

---

```

answer  $\leftarrow \text{ASK}(f)$ 
if answer is "Inconsistent": then
  return "Inconsistent"
else if answer is "Entailed": then
  return "Redundant"
else if answer is "Contradicted": then
  return "Refused"
else if answer is "Contingent": then
   $\mathcal{KB} \leftarrow \mathcal{KB} \cup \{f\}$ 
  return "Added"
end if

```

---

## 8.2 Formulas as Random Variables

In order to present logical entailment as extreme cases of more generic probabilistic reasoning, we now provide probabilistic interpretations of propositional formulas. In the next sections, we will investigate two ways of interpreting basis encodings of formulas as conditional probabilities. The atom centric one, which understands the atomic legs as conditions and calculates the truth of the formula, leads to a direct interpretation of  $\beta^f$  as a conditional probability distribution. When instead taking the formula itself centric, we get uniform distributions of its models and the complement, when conditioning on the satisfaction of the formula.

### 8.2.1 Probabilistic queries by formulas

Let  $\mathbb{P}[X_{[d]}]$  be a joint distribution of atomic variables  $X_k$ , where  $k \in [d]$ , taking variables in  $m_k = 2$ . Let us then ask a query in the formalism of Def. 35, where the query function is assumed to be a propositional formula. The joint distribution can be extended to a variable  $Y_f$  representing the satisfaction of a formula  $f$  given an assignment to the atoms, by adding its basis encoding as

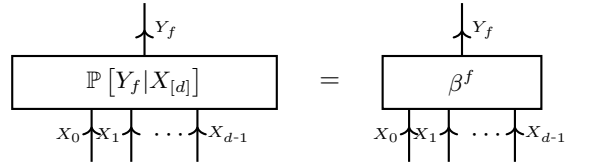
$$\mathbb{P}[Y_f, X_{[d]}] = \langle \beta^f[Y_f, X_{[d]}], \mathbb{P}[X_{[d]}] \rangle [X_{[d]}] .$$

Let us note, that this is a normed probability distribution, since  $\langle \beta^f[Y_f, X_{[d]}] \rangle [X_{[d]}] = \mathbb{I}[X_{[d]}]$  and  $\mathbb{P}[X_{[d]}]$  is normed.

Conditioning this probability distribution on the atoms, we get

$$\mathbb{P}[Y_f|X_{[d]}] = \beta^f[X_{[d]}] .$$

We thus interpret the basis encoding of a formula as a conditional probability of  $f$  given the assignments to the atoms  $X_{[d]}$  and depict this by



To be more precise, we have for any  $x_{[d]}$

$$\mathbb{P}[Y_f|X_{[d]} = x_{[d]}] = \begin{cases} \epsilon_0[Y_f] & \text{if } f[X_{[d]} = x_{[d]}] = 0, \text{ i.e. } x_{[d]} \text{ is not a model of } f \\ \epsilon_1[Y_f] & \text{if } f[X_{[d]} = x_{[d]}] = 1, \text{ i.e. } x_{[d]} \text{ is a model of } f \end{cases} .$$

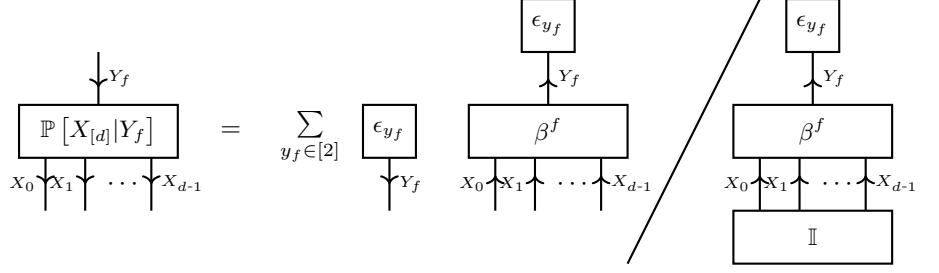
Since the conditional query  $\mathbb{P}[Y_f|X_{[d]}]$  provides an interpretation of  $\beta^f$  as a conditional probability, we interpret  $\mathbb{P}[Y_f]$  as a marginal distribution inherited by  $\mathbb{P}[X_{[d]}]$ . This is also reflected in the fact that both  $\mathbb{P}[Y_f|X_{[d]}]$  and  $\beta^f[Y_f, X_{[d]}]$  are directed, since the first is a normalization by Def. 35 and the second an basis encoding of a formula. Probabilistic queries (see Def. 35), which functions are propositional formulas are thus answered by the satisfaction rate of a propositional formula given a joint distribution of the corresponding atoms.

### 8.2.2 Uniform distributions of the models

Let us now converse the order of conditioning from  $\mathbb{P}[Y_f|X_{[d]}]$  to  $\mathbb{P}[X_{[d]}|Y_f]$ . In this way, we understand a propositional formula as a definition of a joint probability distributions of the atoms, instead of a formulation of a probabilistic query against a joint distribution. To this end, we define by the single tensor core  $\{\beta^f[Y_f, X_{[d]}]\}$  a Markov Network  $\mathbb{P}^{\{f\} \cup [d]}[Y_f, X_{[d]}]$ . By definition we have

$$\mathbb{P}^{\{f\} \cup [d]}[X_{[d]}|Y_f] = \langle \beta^f \rangle [X_{[d]}|Y_f] .$$

We depict this construction by:



Let us further investigate the slices of  $\mathbb{P}[X_{[d]}|f]$  with respect to  $f$ , which define distributions of the states of the factored system. To this end, let us condition on the event of  $f = 1$ , for which we have the distribution

$$\mathbb{P}[X_{[d]}|Y_f = 1] = \frac{1}{\langle f \rangle [\emptyset]} \sum_{x_{[d]} \in \times_{k \in [d]} [2] : f[X_{[d]} = x_{[d]}] = 1} \epsilon_{x_{[d]}} [X_{[d]}}. \quad (15)$$

With  $\langle f \rangle [\emptyset]$  being the number of models of  $f$ , this is the uniform distribution among the models of  $f$ . Conversely, when conditioning on the event  $Y_f = 0$  we get a uniform distribution of the models of  $\neg f$ .

The probability distribution in Equation (15) is well defined except for the case that  $\langle f \rangle [\emptyset] = 0$ . In that case we would have  $f[X_{[d]}] = 0[X_{[d]}]$  and call  $f$  unsatisfiable, since it has no models.

From an epistemological point of view, probability theory is a generalization of logics, since we allow for probability values in the interval  $[0, 1]$ . The set of distributions being constructed by conditioning on propositional formulas as in Equation (15) correspond within the set of probability distributions with those being constant on their support. While the distributions build a  $2^d - 1$ -dimensional manifold, the formulas parametrize by this construction  $2^{(2^d)}$ .

### 8.2.3 Probability of a formula given a Knowledge Base

We now combine the ideas of the previous two subsections and define probabilities of formulas  $f$  given the satisfaction of another formula  $\mathcal{KB}$ , which we call a knowledge base. We have

$$\begin{aligned} \mathbb{P}[Y_f|Y_{\mathcal{KB}}] &= \langle \mathbb{P}[Y_f|X_{[d]}], \mathbb{P}[X_{[d]}|Y_{\mathcal{KB}}] \rangle [Y_f, Y_{\mathcal{KB}}] \\ &= \langle \beta^f, \beta^{\mathcal{KB}} \rangle [Y_f|Y_{\mathcal{KB}}]. \end{aligned}$$

We notice, that we have to assume a satisfiable knowledge base  $\mathcal{KB}$  for this construction to be well-defined.

Of special interest is the conditional probability of  $Y_f$  given that  $Y_{\mathcal{KB}}$  is satisfied, that is

$$\begin{aligned} \mathbb{P}[Y_f|Y_{\mathcal{KB}} = 1] &= \langle \{\beta^f, \mathcal{KB}\} \rangle [Y_f|\emptyset] \\ &= \frac{\langle \{\beta^f, \mathcal{KB}\} \rangle [Y_f]}{\langle \{\mathcal{KB}\} \rangle [\emptyset]}. \end{aligned}$$

This conditional probability establishes a connection with the entailment relation of propositional formulas, as we show next.

**Theorem 47.** *Given a satisfiable formula  $\mathcal{KB}$ , we have  $\mathcal{KB} \models f$ , if and only if*

$$\mathbb{P}[Y_f = 0|Y_{\mathcal{KB}} = 1] = 0.$$

*Proof.* Since  $\mathcal{KB}$  is satisfiable, we have  $\langle \mathcal{KB} \rangle [\emptyset] > 0$  and

$$\mathbb{P}[Y_f = 0|Y_{\mathcal{KB}} = 1] = \frac{\langle \neg f, \mathcal{KB} \rangle [\emptyset]}{\langle \mathcal{KB} \rangle [\emptyset]}.$$

This term vanishes if and only if  $\langle \neg f, \mathcal{KB} \rangle [\emptyset]$  vanish. Now, by The. 43 we have  $\mathcal{KB} \models f$  if and only if  $\langle \mathcal{KB}, \neg f \rangle [\emptyset] = 0$ , which is therefore equal to  $\mathbb{P}[Y_f = 0|Y_{\mathcal{KB}} = 1] = 0$ .  $\square$

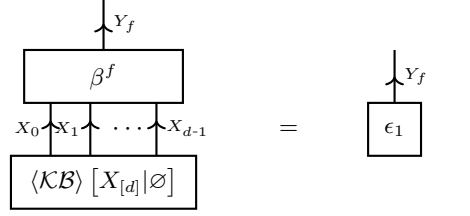
Since any conditional distribution is directed, we have

$$\mathbb{P}[Y_f|Y_{\mathcal{KB}} = 1] = \epsilon_0[Y_f] \text{ if } \mathcal{KB} \models \neg f \quad (16)$$

$$\epsilon_1[Y_f] \text{ if } \mathcal{KB} \models f \quad (17)$$

$$\notin \{\epsilon_0[Y_f], \epsilon_1[Y_f]\} \text{ else.} \quad (18)$$

We depict the case of entailment  $\mathcal{KB} \models f$  by the contraction diagram



We can further omit the normalization by  $\langle \mathcal{KB} \rangle [\emptyset]$  when deciding entailment, and thus drop the assumption of satisfiability of  $\mathcal{KB}$ , as we state next.

**Theorem 48.** *Given a formula  $\mathcal{KB}$ , we have  $\mathcal{KB} \models f$  (respectively  $\mathcal{KB} \models \neg f$ ), if and only if*

$$\langle \mathcal{KB}, \beta^f \rangle [Y_f = 0] = 0 \quad (\text{respectively } \langle \mathcal{KB}, \beta^f \rangle [Y_f = 1] = 0).$$

*Proof.* This follows from The. 43 using that

$$\beta^f [Y_f = 0, X_{[d]}] = \neg f [X_{[d]}] \quad \text{and} \quad \beta^f [Y_f = 1, X_{[d]}] = f [X_{[d]}].$$

□

Relating entailment to probability distributions motivates an extension of the entailment provided by Def. 46 to arbitrary probability distributions.

**Definition 47.** *For any propositional formula  $f [X_{[d]}]$  we say that a probability distribution  $\mathbb{P}^{X_{[d]}}$  probabilistically entails  $f$ , denoted as  $\mathbb{P} \models f$ , if*

$$\langle \mathbb{P} [X_{[d]}], \beta^f [Y_f, X_{[d]}] \rangle [Y_f = 0] = 0.$$

*If  $\mathbb{P} \models \neg f$ , that is  $\langle \mathbb{P} [X_{[d]}], \beta^f [Y_f, X_{[d]}] \rangle [Y_f = 1] = 0$ , we say that  $\mathbb{P}$  probabilistically contradicts  $f$ .*

We note, that when choosing for a formula  $\mathcal{KB}$  the uniform distribution

$$\mathbb{P} [X_{[d]}] = \langle X_{[d]} \rangle [Y_{\mathcal{KB}} = 1 | \emptyset]$$

among its models, then probabilistic entailment  $\mathbb{P} \models f$  of a propositional formula  $f$  is by The. 47 equivalent to  $\mathcal{KB} \models f$ .

### 8.2.4 Knowledge Bases as Base Measures for Probability Distributions

Let us now further relate the probabilistic entailment provided by Def. 47 with logical entailment, by constructing a corresponding propositional formula to an arbitrary distribution. Given a generic probability distribution  $\mathbb{P}$  we can build a Knowledge Base by

$$\mathcal{KB}^{\mathbb{P}} = \mathbb{I}_{\neq 0} \circ \mathbb{P},$$

where  $\mathbb{I}_{\neq 0} : \mathbb{R} \rightarrow \mathbb{R}$  denotes the indicator function of the support defined as

$$\mathbb{I}_{\neq 0}(x) = \begin{cases} 0 & \text{if } x = 0 \\ 1 & \text{else} \end{cases}. \quad (19)$$

Probabilistic entailment with respect to  $\mathbb{P}$  is then equivalent to entailment with respect to  $\mathcal{KB}^{\mathbb{P}}$ , as we show next.

**Theorem 49.** *Any probability distribution  $\mathbb{P} [X_{[d]}]$  probabilistically entails a formula  $f [X_{[d]}]$ , if and only if  $\mathcal{KB}^{\mathbb{P}} \models f$ .*

*Proof.* Whenever  $\mathbb{P}$  does not entail  $f$  probabilistically we find a state  $x_{[d]} \in \times_{k \in [d]} [2]$  such that

$$\mathbb{P} [X_{[d]} = x_{[d]}] > 0 \quad \text{and} \quad f [X_{[d]} = x_{[d]}] = 0.$$

We further have  $\mathbb{P} [X_{[d]} = x_{[d]}] > 0$  if and only if  $\mathcal{KB}^{\mathbb{P}} [X_{[d]} = x_{[d]}] = 1$ . Therefore the statement

$$(\mathcal{KB}^{\mathbb{P}} [X_{[d]} = x_{[d]}] = 1) \Rightarrow (f [X_{[d]} = x_{[d]}] = 1)$$

is not satisfied. Together,  $\mathbb{P} \models f$  does not holds if and only if

$$\forall x_{[d]} \in \prod_{k \in [d]} [2] : \left( \mathcal{KB}^{\mathbb{P}}[X_{[d]} = x_{[d]}] = 1 \right) \Rightarrow (f[X_{[d]} = x_{[d]}] = 1)$$

is not satisfied. Therefore, probabilistic entailment of  $f$  by  $\mathbb{P}$  is equivalent to logical entailment of  $f$  by  $\mathcal{KB}^{\mathbb{P}}$ .  $\square$

Let us use this to connect the entailment formalism with the representability (see Def. 17) and positivity (see Def. 18) of distributions with respect to boolean base measures.

**Theorem 50.** *Let  $\mathbb{P}$  be a distribution of boolean variables and let  $\nu$  be a boolean base measure. Then,  $\mathbb{P}$  is representable with respect to  $\nu$ , if and only if  $\mathbb{I}_{\neq 0} \circ \mathbb{P} \models \nu$ . Further,  $\mathbb{P}$  is positive with respect to  $\nu$ , if and only if  $\nu = \mathbb{I}_{\neq 0} \circ \mathbb{P}$ .*

*Proof.* To show the first claim, let  $\mathbb{P}$  be a distribution and  $\nu$  be a base measure. With Def. 17,  $\mathbb{P}$  is representable with respect to  $\nu$ , if and only if

$$\forall x_{[d]} \in \prod_{k \in [d]} [2] : (\nu[X_{[d]} = x_{[d]}] = 0) \Rightarrow (\mathbb{P}[X_{[d]} = x_{[d]}] = 0)$$

This is equal to

$$\forall x_{[d]} \in \prod_{k \in [d]} [2] : (\mathbb{I}_{\neq 0} \circ \mathbb{P}[X_{[d]} = x_{[d]}] = 1) \Rightarrow (\nu[X_{[d]} = x_{[d]}] = 1)$$

and by definition Def. 46 equal to  $\nu \models \mathbb{I}_{\neq 0} \circ \mathbb{P}$ .

To proof the second claim, we show that when  $\mathbb{P}$  is in addition positive with respect to  $\nu$ , then also  $\nu \models \mathbb{I}_{\neq 0} \circ \mathbb{P}$  and thus  $\nu = \mathbb{I}_{\neq 0} \circ \mathbb{P}$ . Let  $\mathbb{P}$  be a distribution, which is representable with respect to  $\nu$ . Then  $\mathbb{P}$  is positive with respect to  $\nu$ , if and only if

$$\forall x_{[d]} \in \prod_{k \in [d]} [2] : (\nu[X_{[d]} = x_{[d]}] = 1) \Rightarrow (\mathbb{P}[X_{[d]} = x_{[d]}] > 0)$$

This is equal to

$$\forall x_{[d]} \in \prod_{k \in [d]} [2] : (\nu[X_{[d]} = x_{[d]}] = 1) \Rightarrow (\mathbb{I}_{\neq 0} \circ \mathbb{P}[X_{[d]} = x_{[d]}] = 1)$$

and thus  $\nu \models \mathbb{I}_{\neq 0} \circ \mathbb{P}$ .  $\square$

### 8.3 Constraint Satisfaction Problems

Let us now explore a more general class of logical inference problems and discuss probabilistic entailment within that class. We then provide further examples based on categorical constraints. Following Chapter 5 in Russell and Norvig (2021), we now define Constraint Satisfaction Problems.

**Definition 48.** *Let there be a hypergraph  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$  and  $\tau^{\mathcal{G}}$  be a tensor network of boolean constraint tensors  $\tau^e[X_e]$  to each  $e \in \mathcal{E}$ , that is*

$$\tau^{\mathcal{G}} = \{\tau^e[X_e] : e \in \mathcal{E}\}.$$

*The Constraint Satisfaction Problem (CSP) to  $\tau^{\mathcal{G}}$  is the decision whether there is a state  $x_{\mathcal{V}}$  such that*

$$\langle \tau^{\mathcal{G}} \rangle [X_{\mathcal{V}} = x_{\mathcal{V}}] = 1.$$

*We say the CSP is satisfiable, when there is such a state, and unsatisfiable if not.*

#### 8.3.1 Deciding entailment on Markov Networks

Deciding entailment on Markov Networks is a general class of constraint satisfaction problems. Here, any factor tensor in the Markov Networks produces a constraint tensor in the respective CSP.

**Theorem 51.** Let  $\mathbb{P}^{\mathcal{G}}$  be a Markov Network to the Tensor Network  $\tau^{\mathcal{G}} = \{\tau^e[X_e] : e \in \mathcal{E}\}$  on a hypergraph  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ . For each  $e \in \mathcal{E}$  we build the factor constraint cores

$$\tilde{\tau}^e[e] = \mathbb{I}_{\neq 0} \circ \tau^e[X_e] .$$

Let further  $f[X_{\tilde{\mathcal{V}}}]$  be a formula depending on the variables  $\tilde{\mathcal{V}}$ , and build  $\tilde{\mathcal{G}} = (\mathcal{V}, \mathcal{E} \cup \{\tilde{\mathcal{V}}\})$ . Then we have that  $\mathbb{P}^{\mathcal{G}} \models f$  if and only if the constraint satisfaction problem of  $\tilde{\mathcal{G}}$  to the constraint tensors

$$\{\tilde{\tau}^e : e \in \mathcal{E}\} \cup \{\neg f\}$$

is unsatisfiable.

*Proof.* We first show, that

$$\mathbb{I}_{\neq 0} \circ \mathbb{P}^{\mathcal{G}}[X_{\mathcal{V}}] = \langle \{\tilde{\tau}^e : e \in \mathcal{E}\} \rangle [X_{\mathcal{V}}] .$$

To this end, let  $x_{\mathcal{V}} \in \times_{v \in \mathcal{V}} [m_v]$  be arbitrary. We have  $\mathbb{P}^{\mathcal{G}}[X_{\mathcal{V}} = x_{\mathcal{V}}] = 0$  if and only if at there is an edge  $e \in \mathcal{E}$  with  $\tau^e[X_e = x_e]$ . But this is equivalent to

$$\langle \{\tilde{\tau}^e : e \in \mathcal{E}\} \rangle [X_{\mathcal{V}}] .$$

We thus have for any  $x_{\mathcal{V}} \in \times_{v \in \mathcal{V}} [m_v]$

$$\mathbb{I}_{\neq 0} \circ \mathbb{P}^{\mathcal{G}}[X_{\mathcal{V}} = x_{\mathcal{V}}] = \langle \{\tilde{\tau}^e : e \in \mathcal{E}\} \rangle [X_{\mathcal{V}} = x_{\mathcal{V}}] .$$

To continue, we have  $\mathbb{P}^{\mathcal{G}} \models f$  if and only if

$$\langle \tau^{\mathcal{G}}[X_{[d]}], \neg f[X_{[d]}] \rangle [\emptyset] = 0$$

which is equal to

$$\langle \mathbb{I}_{\neq 0} \circ \tau^{\mathcal{G}}[X_{[d]}], \neg f[X_{[d]}] \rangle [\emptyset] = 0 .$$

We notice, that this is the unsatisfiability of the claimed Constraint Satisfaction Problem. □

For any positive tensor  $\tau$  we have

$$\mathbb{I}_{\neq 0} \circ \tau[X_e] = \mathbb{I}[X_e] ,$$

which does not influence the distribution and can be omitted from the Markov Network. By The. 51, when deciding entailment, we can reduce all tensors of a Markov Network to their support and omit those with full support. Since the support indicating tensors  $\mathbb{I}_{\neq 0} \circ \tau[X_e]$  are boolean, each is a propositional formula and the Markov Network is turned into a Knowledge Base of their conjunctions. Deciding probabilistic entailment is thus traced back to logical entailment.

Exponential families have a tensor network representation by a Markov Network (see The. 11). However, all factors corresponding with a coordinate of the statistic  $\mathcal{S}$  have a trivial support, and therefore do not influence the support of the distribution. The only tensors with non-trivial support are those to the boolean base measure  $\nu$ .

### 8.3.2 Categorical Constraints

We so far in this chapter made the assumption that all categorical variables in factored systems to be represented by propositional logics take binary values (i.e.  $m = 2$ ). In cases where a categorical variable  $X$  takes multiple values we define for each  $x$  an atomic formula  $X_x$  representing whether  $X$  is assigned by  $x$  in a specific state. Following this construction we have the constraint that exactly one of the atoms  $X_x$  is 1 at each state.

**Definition 49** (Categorical Constraint and Atomization Variables). Given a list  $X_0, \dots, X_{m-1}$  of boolean variables and a categorical variable  $X$  with dimension  $m$  a categorical constraint is a tensor  $Z[X, X_{[m]}]$  defined as

$$Z[X_{[m]} = x_{[m]}, X = x] = \begin{cases} 1 & \text{if } x_{[m]} = \epsilon_x \quad \left( \text{i.e. } \forall k \in [m] (x = k) \Leftrightarrow (x_k = 1) \right) \\ 0 & \text{else.} \end{cases}$$

We then call the variables  $X_0, \dots, X_{m-1}$  the atomization variables to the categorical variable  $X$ .



With The. 126 the basis encoding  $\beta^Z$  decomposes in a basis CP format (see Figure 21b) of if its coordinate maps  $Z^k$ , where  $k \in [m]$ , defined as

$$Z^k [X_k = x_k, X = x] = \begin{cases} 1 & \text{if } x = k \\ 0 & \text{else.} \end{cases}$$

Their basis encoding are decomposed as

$$\beta^{Z^k} [X_k, X] = \epsilon_1 [X_k] \otimes \epsilon_k [X] + \epsilon_0 [X_k] \otimes (\mathbb{I}[X] - \epsilon_k [X]). \quad (20)$$

We further have by The. 126

$$\beta^Z [X_{[m]}, X] = \left\langle \{\beta^{Z^k} [X_k, X] : k \in [m]\} \right\rangle [X, X_0, \dots, X_{m-1}].$$

In the next theorem we show how a categorical constraint can be enforced in a tensor network by adding the tensor  $Z$  to a contraction.

**Theorem 52.** For any tensor  $\tau [X_{[d]}]$  and a categorical constraint defined by an ordered subset  $X_A \subset X_{[d]}$ , a variable  $X \in X_{[d]}$  we have

$$\langle \tau [X_{[d]}], Z [X_A, X] \rangle [X_0 = x_0, \dots, X_{d-1} = x_{d-1}] = \begin{cases} \tau [X_0 = x_0, \dots, X_{d-1} = x_{d-1}] & \text{if } x_A = \epsilon_x \\ 0 & \text{else.} \end{cases}$$

Here by  $x_A$  we denote the restriction of  $x_{[d]}$  on the set  $A$ .

*Proof.* For any  $x_{[d]}$  we have

$$\langle \tau [X_{[d]}], Z \rangle [X_0 = x_0, \dots, X_{d-1} = x_{d-1}] = \tau [x_{[d]}] \cdot Z [X_A = x_A, X = x].$$

If  $x_A = \epsilon_x$  we have  $Z [X_A = x_A, X = x] = 1$  and thus

$$\langle \tau [X_{[d]}], Z \rangle [X_0 = x_0, \dots, X_{d-1} = x_{d-1}] = \tau [x_{[d]}].$$

If  $x_A \neq \epsilon_x$  then  $Z [X_A = x_A, X = x] = 0$  and

$$\langle \tau [X_{[d]}], Z \rangle [X_0 = x_0, \dots, X_{d-1} = x_{d-1}] = 0.$$

□

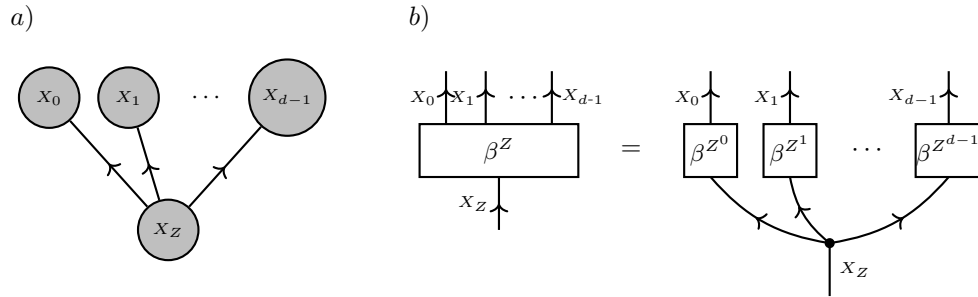


Figure 21: Representation of a categorical constraint in a CP Format tensor network. a) Representation of the dependency of the graphical model. b) Tensor Representation with further network decomposition.

**Remark 6** (Constraint Satisfaction Problems of Categorical Constraints). We can define CSPs by collection of categorical constraints. An example, where the corresponding Constraint Satisfaction Problem is unsatisfiable are the categorical constraints to the three sets

$$\{X_0, X_1, X_2, X_3\}, \{X_0, X_1\}, \{X_2, X_3\}.$$

$X_{0,0}$	$X_{0,1}$	$X_{0,2}$	$X_{0,3}$	$X_{0,4}$	$X_{0,5}$	$X_{0,6}$	$X_{0,7}$	$X_{0,8}$
$X_{1,0}$	$X_{1,1}$	$X_{1,2}$	$X_{1,3}$	$X_{1,4}$	$X_{1,5}$	$X_{1,6}$	$X_{1,7}$	$X_{1,8}$
$X_{2,0}$	$X_{2,1}$	$X_{2,2}$	$X_{2,3}$	$X_{2,4}$	$X_{2,5}$	$X_{2,6}$	$X_{2,7}$	$X_{2,8}$
$X_{3,0}$	$X_{3,1}$	$X_{3,2}$	$X_{3,3}$	$X_{3,4}$	$X_{3,5}$	$X_{3,6}$	$X_{3,7}$	$X_{3,8}$
$X_{4,0}$	$X_{4,1}$	$X_{4,2}$	$X_{4,3}$	$X_{4,4}$	$X_{4,5}$	$X_{4,6}$	$X_{4,7}$	$X_{4,8}$
$X_{5,0}$	$X_{5,1}$	$X_{5,2}$	$X_{5,3}$	$X_{5,4}$	$X_{5,5}$	$X_{5,6}$	$X_{5,7}$	$X_{5,8}$
$X_{6,0}$	$X_{6,1}$	$X_{6,2}$	$X_{6,3}$	$X_{6,4}$	$X_{6,5}$	$X_{6,6}$	$X_{6,7}$	$X_{6,8}$
$X_{7,0}$	$X_{7,1}$	$X_{7,2}$	$X_{7,3}$	$X_{7,4}$	$X_{7,5}$	$X_{7,6}$	$X_{7,7}$	$X_{7,8}$
$X_{8,0}$	$X_{8,1}$	$X_{8,2}$	$X_{8,3}$	$X_{8,4}$	$X_{8,5}$	$X_{8,6}$	$X_{8,7}$	$X_{8,8}$

Figure 22: Sudoku grid of basic categorical variables  $X_{i,j}$ , here drawn in the standard case of  $n = 3$ , each with dimension  $m = n^2 = 9$ . Each basic categorical variables has  $n^2$  corresponding atomization variables, which are further atomization variables to the row, column and squares constraints. Instead of depicting those constraints by hyperedges in a variable dependency graph, we here just indicate their existence through row, column and squares blocks.

**Example 7 (Sudoku).** An interesting example, where categorical constraints are combined is Sudoku, the game of assigning numbers to a grid (see for example Section 5.2.6 in Russell and Norvig (2021)). The basic variables therein are  $X_{i,j}$ , with  $m_{i,j} = n^2$  and  $i, j \in [n^2]$ . By understanding  $i$  as a line index and  $j$  as a column index, they are ordered in a grid as sketched in Figure 22 in the case  $n = 3$ .

For a  $n \in \mathbb{N}$  we further define the atomization variables  $X_{i,j,k}$  where  $i, j, k \in [n^2]$  and  $m_{i,j,k} = 2$ . These  $n^6$  variables are the booleans indicating whether a specific position has a specific number assigned. The consistency of the atomization variables to the basic variables is then for each  $i, j \in [n^2]$  ensured by the constraints

$$\{X_{i,j,k} : k \in [n^2]\}.$$

We further have  $3 \cdot n^2$  constraints by the

- Row constraints: Each number  $k$  appears exactly once in each row  $i \in [n^2]$ , captured by the constraints

$$\{X_{i,j,k} : j \in [n^2]\}.$$

- Column constraints: Each number  $k$  appears exactly once in each column  $j \in [n^2]$ , captured by the constraints

$$\{X_{i,j,k} : i \in [n^2]\}.$$

- Square constraints: Each number appears exactly once in each square  $s, r \in [n]$ , captured by the constraints

$$\{X_{i+n \cdot s, j+n \cdot r, k} : i, j \in [n]\}.$$

In total we have  $3 \cdot n^2 + n^4$  constraints for  $n^6$  variables.

Deciding whether a Sudoku has a solution is a Constraint Satisfaction Problem Simonis (2005), which is NP-hard Agerbeck and Hansen (2008). Let us notice, that due to this large number of variables and constraints, direct solution of the problem by a global contraction is not feasible. For efficient algorithmic solutions, we instead refer to Sect. 8.4.

## 8.4 Deciding Entailment by local contractions

When having a Constraint Satisfaction Problem on a large number of variables, which are densely connected by constraint tensors, direct exploitation of the global entailment criterion in The. 43 will be infeasible. An alternative to deciding entailment by global operations is the use of local operations. Here we interpret a part of the network (for example a single core) as an own knowledge base (with atomic formulas being the roots of the directed subgraph, that is potentially differing with the atoms in the global perspective) and perform entailment with respect to that.

### 8.4.1 Monotonicity of entailment

Vanishing local contractions provide sufficient but not necessary criterion to decide entailment, as we show in the next theorem.

**Theorem 53** (Monotonicity of Entailment). *For any Markov Network on the decorated hypergraph  $\mathcal{G}$  and any subgraph  $\tilde{\mathcal{G}}$ , we have for any formula that  $\mathbb{P}^{\mathcal{G}} \models f$  if  $\mathbb{P}^{\tilde{\mathcal{G}}} \models f$ .*

To prove the theorem, we first establish the following lemma that states if a contraction of non-negative tensors vanishes, the vanishing of a contraction over a subset of these tensors is a sufficient criterion.

**Lemma 13.** *For any non-negative tensor network  $\tau^{\mathcal{G}}$  on  $\mathcal{G}$  and  $\tilde{\mathcal{E}} \subset \mathcal{E}$  we have the following. For  $\tilde{\mathcal{G}} = (\tilde{\mathcal{V}}, \tilde{\mathcal{E}})$  with  $\tilde{\mathcal{V}} = \cup_{e \in \tilde{\mathcal{E}}} e$  and the tensor network  $\tau^{\tilde{\mathcal{G}}}$  with tensors coinciding on  $\tilde{\mathcal{E}}$  with those in  $\tau^{\mathcal{G}}$  we have*

$$\langle \tau^{\mathcal{G}} \rangle [\emptyset] = 0$$

$$\text{if } \langle \tau^{\tilde{\mathcal{G}}} \rangle [\emptyset] = 0.$$

*Proof.* Since the tensor network  $\tau^{\tilde{\mathcal{G}}}$  is non-negative, we have whenever  $\langle \tau^{\tilde{\mathcal{G}}} \rangle [\emptyset] = 0$  that

$$\langle \tau^{\tilde{\mathcal{G}}} \rangle [X_{\tilde{\mathcal{V}}}] = 0 [X_{\tilde{\mathcal{V}}}] .$$

It follows with the commutation of contractions (see The. 131 in Chapter 20), that

$$\begin{aligned} \langle \tau^{\mathcal{G}} \rangle [X_{\mathcal{V}}] &= \left\langle \{ \tau^e : e \in \mathcal{E} / \tilde{\mathcal{E}} \} \cup \{ \langle \tau^{\tilde{\mathcal{G}}} \rangle [X_{\tilde{\mathcal{V}}}] \} \right\rangle [X_{\mathcal{V}}] \\ &= \left\langle \{ \tau^e : e \in \mathcal{E} / \tilde{\mathcal{E}} \} \cup \{ 0 [X_{\tilde{\mathcal{V}}}] \} \right\rangle [X_{\mathcal{V}}] \\ &= 0 \end{aligned}$$

Thus, also the contraction of  $\tau^{\mathcal{G}}$  vanishes in this case. □

*Proof of The. 53.* We use Lem. 13 on the subset  $\tau^{\tilde{\mathcal{G}}}$  of the cores  $\tau^{\mathcal{G}}$  to the Markov Network  $\mathbb{P}^{\mathcal{G}}$ , which itself defines the Markov Network  $\mathbb{P}^{\tilde{\mathcal{G}}}$ . Whenever  $\mathbb{P}^{\tilde{\mathcal{G}}} \models f$  for a formula  $f$ , then we have by The. 43

$$\langle \tau^{\tilde{\mathcal{G}}} \cup \{ \neg f \} \rangle [\emptyset] = 0 .$$

It follows with Lem. 13 that also

$$\langle \tau^{\mathcal{G}} \cup \{ \neg f \} \rangle [\emptyset] = 0 .$$

and therefore  $\mathbb{P}^{\mathcal{G}} \models f$ . □

**Remark 7.** *To make use of The. 53 we can exploit any entailment criterion. However, there is no general statement about entailment possible, when the local entailment does not hold. The. 53 therefore just provides a sufficient but not necessary criterion of entailment with respect to  $\mathbb{P}^{\mathcal{G}}$ .*

### 8.4.2 Knowledge Cores

To store preliminary conclusions, we define auxiliary knowledge cores storing constraints on variables  $e \in \mathcal{V}$ . They are understood as logical formulas to the atomization variables ( $X_e = x_e$ ) of the respective formulas

$$\kappa^e [X_e] = \bigvee_{x_e : \kappa^e [X_e = x_e]} \bigwedge_{v \in e} (X_e = x_e) .$$

**Definition 50.** Let  $\tau^{\mathcal{G}}$  be a constraint satisfaction problem. We say that a knowledge core  $\kappa^e[X_e]$  is sound for  $\tau^{\mathcal{G}}$ , if

$$\mathbb{I}_{\neq 0} \circ \langle \tau^{\mathcal{G}} \rangle [X_e] \prec \kappa^e[X_e]$$

and complete for  $\tau^{\mathcal{G}}$  if in addition

$$\mathbb{I}_{\neq 0} \circ \langle \tau^{\mathcal{G}} \rangle [X_e] = \kappa^e[X_e] .$$

### 8.4.3 Knowledge Propagation

We now provide a solution algorithm for constraint satisfaction problems by propagating local contractions. The dynamic programming paradigm is implemented by the storage of partial entailment results in Knowledge Cores. We then iterate over local entailment checks, where we recursively add further entailment checks to be redone due to additional knowledge. This local entailment scheme is called Knowledge Propagation and described in a generic way in Algorithm 8.

---

#### Algorithm 8 Knowledge Propagation

---

**Require:** Boolean Tensor Network  $\tau^{\mathcal{G}}$  on  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ , domain edges  $\mathcal{E}^k$  and a set  $\mathcal{U}$  of subsets of  $\mathcal{E}$  for local propagation

**Ensure:** Knowledge cores  $\kappa^e[X_e]$  for  $e \in \mathcal{E}^k$  with  $\langle \tau^{\mathcal{G}} \rangle [X_e] \prec \kappa^e[X_e]$

---

Initialize for all  $e \in \mathcal{E}^k$ :

$$\kappa^e[X_e] = \mathbb{I}[X_e]$$

Initialize a queue

$$\mathcal{Q} = \mathcal{U}$$

**while**  $\mathcal{Q}$  is not empty **do**

    Choose a set of edges from the queue

$$\tilde{\mathcal{E}} \leftarrow \mathcal{Q}.\text{pop}()$$

**for**  $e \in \mathcal{E}^k$  with  $e \cap \bigcup_{\tilde{e} \in \tilde{\mathcal{E}}} \tilde{e} \neq \emptyset$  **do**

        Contract

$$\tau[X_e] = \mathbb{I}_{\neq 0} \circ \left\langle \{ \tau^{\tilde{e}}[X_{\tilde{e}}] : \tilde{e} \in \tilde{\mathcal{E}} \} \cup \{ \kappa^e[X_e] : e \in \mathcal{E}^k, e \cap \bigcup_{\tilde{e} \in \tilde{\mathcal{E}}} \tilde{e} \neq \emptyset \} \right\rangle [X_e]$$

**if**  $\tau[X_e] \neq \kappa^e[X_e]$  **then**

$$\kappa^e[X_e] \leftarrow \tau[X_e]$$

**for**  $\tilde{\mathcal{E}} \in \mathcal{U}$  with  $e \cap \bigcup_{\tilde{e} \in \tilde{\mathcal{E}}} \tilde{e} \neq \emptyset$  **do**

$$\mathcal{Q}.\text{push}(\tilde{\mathcal{E}})$$

**end for**

**end if**

**end for**

**end while**

**return**  $\{ \kappa^e[X_e] : e \in \mathcal{E}^k \}$

---

Each chosen subset  $\tilde{\mathcal{E}} \in \mathcal{U}$  is understood as a local knowledge base, which is then applied for local entailment. The knowledge cores are understood as messages, which propagate information from different regions of a tensor network (see Chapter 20).

There are different ways of implementing Algorithm 8, by choosing the set  $\mathcal{U}$  of constraint sets  $\tilde{\mathcal{E}}$  and domain  $\mathcal{E}^k$ . The AC-3 algorithm (see Mackworth (1977)) is a specific instance, where knowledge cores are assigned to single variables and propagation is performed on single constraint cores.

**Theorem 54.** At any state of the Knowledge Propagation Algorithm 8, we have that each knowledge core  $\kappa^{\tilde{e}}$  is sound for  $\tau^{\mathcal{G}}$ . After each update in Algorithm 8,  $\kappa^{\tilde{e}}$  is further monotonically decreasing with respect to the partial ordering.

*Proof.* We show the first claim by induction over the update steps in Algorithm 8. At the start, where  $\kappa^e[X_e] = \mathbb{I}[X_e]$ , we trivially have

$$\langle \tau^{\mathcal{G}} \cup \{\kappa^e[X_e] : e \in \mathcal{E}^k\} \rangle [X_{\mathcal{V}}] = \langle \tau^{\mathcal{G}} \cup \{\mathbb{I}[X_e] : e \in \mathcal{E}^k\} \rangle [X_{\mathcal{V}}] = \langle \tau^{\mathcal{G}} \rangle [X_{\mathcal{V}}] .$$

Let us now assume, that for a state of cores  $\{\kappa^e : e \in \mathcal{E}^k\}$  the first claim holds and let  $\tilde{\mathcal{E}} \subset \mathcal{E}$  be chosen for the update of  $\kappa^{\tilde{e}}$ . By the invariance under adding the support of subcontractions, which we will proof in more detail as The. 134 in Chapter 20, we have for the update

$$\tilde{\kappa}^{\tilde{e}}[X_{\tilde{e}}] = \mathbb{I}_{\neq 0} \circ \left\langle \{\tau^e : e \in \tilde{\mathcal{E}}\} \cup \{\kappa^e : e \in \mathcal{E}^k, e \cap \bigcup_{\tilde{e} \in \tilde{\mathcal{E}}} \tilde{e} \neq \emptyset\} \right\rangle [X_{\tilde{e}}]$$

that

$$\langle \tau^{\mathcal{G}} \cup \{\kappa^e[X_e] : e \in \mathcal{E}^k\} \rangle [X_{\mathcal{V}}] = \langle \tau^{\mathcal{G}} \cup \{\kappa^e[X_e] : e \in \mathcal{E}^k\} \cup \{\tilde{\kappa}^{\tilde{e}}[X_{\tilde{e}}]\} \rangle [X_{\mathcal{V}}] .$$

Thus, the first claim holds also after the update of the core to  $\tilde{e}$ .

We further have with the monotonicity of boolean contraction (see The. 134) that for any update of  $\kappa^{\tilde{e}}$  by  $\tilde{\kappa}^{\tilde{e}}$

$$\tilde{\kappa}^{\tilde{e}}[X_{\tilde{e}}] = \mathbb{I}_{\neq 0} \circ \left\langle \{\tau^e : e \in \tilde{\mathcal{E}}\} \cup \{\kappa^e : e \in \mathcal{E}^k, e \cap \bigcup_{e \in \tilde{\mathcal{E}}} e \neq \emptyset\} \right\rangle [X_{\tilde{e}}] \prec \kappa^{\tilde{e}}[X_{\tilde{e}}] .$$

Thus, each Knowledge Core is monotonously decreasing at each update, with respect to the partial tensor ordering.

From the first claim we further have for any  $\tilde{e} \in \tilde{\mathcal{E}}$

$$\langle \{\kappa^{\tilde{e}}[X_{\tilde{e}}]\} \cup \tau^{\mathcal{G}} \cup \{\kappa^e : e \in \mathcal{E}^k / \{\tilde{e}\}\} \rangle [X_{\tilde{e}}] = \langle \tau^{\mathcal{G}} \rangle [X_{\tilde{e}}]$$

And thus in combination with the monotonicity of boolean contraction (see The. 134) that

$$\mathbb{I}_{\neq 0} (\langle \tau^{\mathcal{G}} \rangle [X_{\tilde{e}}]) \prec \kappa^{\tilde{e}}[X_{\tilde{e}}] .$$

□

Let us now show that Knowledge Propagation always terminates. We can further characterize the knowledge cores at termination.

**Definition 51.** We say that a set of knowledge cores  $\{\kappa^e : e \in \mathcal{E}^k\}$  is consistent with a set  $\{\tau^e : e \in \tilde{\mathcal{E}}\}$ , if for any  $e \in \mathcal{E}^k$

$$\kappa^e[X_e] = \mathbb{I}_{\neq 0} \circ \left\langle \{\tau^e : e \in \tilde{\mathcal{E}}\} \cup \{\kappa^e : e \in \mathcal{E}^k\} \right\rangle [X_e] .$$

This property is similar to the completeness of a knowledge core, when interpreting the other knowledge cores and the constraints  $\{\tau^e : e \in \tilde{\mathcal{E}}\}$  as posing a Constraint Satisfaction Problem.

**Theorem 55.** Knowledge Propagation Algorithm 8 always terminates. At termination we further for each  $\tilde{\mathcal{E}} \in \mathcal{U}$  and  $e \in \mathcal{E}^k$  with  $e \cap \bigcup_{e \in \tilde{\mathcal{E}}} e \neq \emptyset$ , that the knowledge cores  $\{\kappa^e : e \in \mathcal{E}^k, e \cap \bigcup_{e \in \tilde{\mathcal{E}}} e \neq \emptyset\}$  are consistent with  $\{\tau^e : e \in \tilde{\mathcal{E}}\}$ .

*Proof.* For each knowledge core, there are finitely many boolean tensor precessing it with respect to the partial order. Therefore, since they are monotonously decreasing, each knowledge core can only be varied finitely many times during the algorithm. In total the algorithm can run only finitely many times in the second for loop, where new sets of edges are pushed into the queue. Therefore the while loop will always terminate.

When after a single pass through the while loop with chosen  $\tilde{\mathcal{E}} \in \mathcal{U}$ , the set  $\tilde{\mathcal{E}}$  is not pushed back into  $\mathcal{Q}$ , we have for any  $e \in \mathcal{E}^k$  with  $e \cap \bigcup_{e \in \tilde{\mathcal{E}}} e \neq \emptyset$  that

$$\kappa^e[X_e] = \mathbb{I}_{\neq 0} \circ \left\langle \{\tau^e : e \in \tilde{\mathcal{E}}\} \cup \{\kappa^e : e \in \mathcal{E}^k, e \cap \bigcup_{e \in \tilde{\mathcal{E}}} e \neq \emptyset\} \right\rangle [X_e] .$$

Whenever the contraction on the right hand side changes during the algorithm, the set  $\tilde{\mathcal{E}}$  is pushed into  $\mathcal{Q}$ . At termination of the algorithm,  $\mathcal{Q}$  is empty, and the claimed consistency therefore has to hold. □

We can exploit the Knowledge Propagation Algorithm 8 for the solution of Constraint Satisfaction Problems, by taking  $\tau^G$  as the tensor network of constraint tensors. Whenever a knowledge core vanishes, we can conclude that the Constraint Satisfaction Problems is not satisfiable, as we show next.

**Corollary 4.** *Let us for a Constraint Satisfaction Problem encoded by  $\tau^G$  run Knowledge Propagation Algorithm 8. Whenever for any  $\tilde{e} \in \mathcal{E}$  we have  $\kappa^{\tilde{e}}[X_{\tilde{e}}] = 0[X_{\tilde{e}}]$ , then the Constraint Satisfaction Problem is not satisfiable.*

*Proof.* Whenever  $\kappa^{\tilde{e}}[X_{\tilde{e}}] = 0[X_{\tilde{e}}]$ , then we have by The. 54

$$\mathbb{I}_{\neq 0}(\langle \tau^G \rangle [X_{\tilde{e}}]) \prec 0[X_{\tilde{e}}]$$

and therefore

$$\langle \tau^G \rangle [\emptyset] = 0.$$

□

When the Knowledge Propagation Algorithm 8 converges in a given implementation and no knowledge core vanishes, we can however not conclude that the Constraint Satisfaction Problem is not satisfiable. However, for any index tuple  $x_V$  to be a solution of the CSP to  $\tau^G$ , we have the necessary condition

$$\forall \tilde{e} \in \tilde{\mathcal{E}} : \kappa^{\tilde{e}}[X_{\tilde{e}} = x_V|_{\tilde{e}}] = 1,$$

where by  $x_V|_{\tilde{e}}$  we denote the restriction of the index tuple  $x_V$  to the variables included in  $\tilde{e}$ . One can use this insight as a starting point for backtracking search, where the assignments to variables  $X_{\tilde{v}}$  are iteratively guessed, based on the restriction that each constraint is locally satisfiable, i.e. .

$$\forall \tilde{e} \in \tilde{\mathcal{E}} : \langle \kappa^{\tilde{e}}[X_{\tilde{e} \cap \tilde{v}} = x_V|_{\tilde{e} \cap \tilde{v}}, X_{\tilde{e}/\tilde{v}}] \rangle [\emptyset] \neq 0.$$

One can understand the guess of an assignment  $x_v$  to a variable  $X_v$ , as it is done during backtracking search, as an inclusion of a constraint

$$\kappa^{\{v\}}[X_v] = \epsilon_{x_v}[X_v].$$

Therefore, Knowledge Propagation Algorithm 8 can be integrated with backtracking search, with iterations between propagations of knowledge and guessing of additional variables.

#### 8.4.4 Applications

Let us exemplify the usage of Knowledge Propagation on Constraint Satisfaction Problems posed by entailment queries on Markov Networks.

**Corollary 5.** *Let Algorithm 8 be run on the cores  $\tau^G \cup \{\beta^f\}$  with an arbitrary design of  $\tilde{\mathcal{E}}$ . Whenever for a formula  $f[X_{\tilde{v}}]$  and a  $\kappa^e$  we have*

$$\langle \kappa^e, \beta^f \rangle [Y_f = 0] = 0$$

*then the Markov Network  $\tau^G$  probabilistically entails  $f$ . If on the contrary*

$$\langle \kappa^e, \beta^f \rangle [Y_f = 1] = 0$$

*then the Markov Network  $\tau^G$  probabilistically entails  $\neg f$ , that is probabilistically contradicts  $f$ .*

*Proof.* This follows from The. 54 ensuring the soundness of Knowledge Propagation and the sufficiency of local entailment. □

**Example 8** (Batch decision of entailment). *Let  $\mathcal{F}$  be a set of formulas and  $\mathbb{P}^G[X_{[d]}]$  a Markov Network, for which it shall be decided, which formulas in  $\mathcal{F}$  are entailed, contradicted or contingent. We can in addition to the cores of the Markov Network create the cores  $\{\beta^f[Y_f, X_{[d]}] : f \in \mathcal{F}\}$  and prepare the knowledge cores*

$$\kappa^{\{f\}}[Y_f].$$

*To decide entailment batchwise, Knowledge Propagation Algorithm 8 can be run. Whenever during the algorithm we have that for a  $f$ , then Cor. 5 implies that if*

$$\kappa^{\{f\}}[Y_f] = \begin{cases} \epsilon_1[Y_f] & \text{then, the formula is entailed by } \mathbb{P}^G. \\ \epsilon_0[Y_f] & \text{then, the formula is contradicted by } \mathbb{P}^G. \\ \mathbb{I}[Y_f] & \text{then no conclusion can be drawn.} \end{cases}$$

*Note, that  $\kappa^{\{f\}}[Y_f] = 0[Y_f]$  can not happen, since this would mean that  $\mathbb{I}_{\neq 0}(\mathbb{P}^G)$  is inconsistent. Thus, at any stage of Algorithm 8, one of the three holds.*

#### 8.4.5 Mimiking Inference Rules by Propagation

While so far we have discussed semantic based entailment, there are inference rules exploiting only logical syntax to infer entailed statements. We here show, that they can be captured by the knowledge propagation scheme, if the sets  $\mathcal{U}$  and  $\mathcal{E}^k$  are chosen properly.

Whenever

$$\bigvee_{f \in \mathcal{F}} f \models h$$

then

$$\epsilon_1[Y_h] = \langle \{f[X_{[d]}] : f \in \mathcal{F}\} \cup \{\beta^h[Y_h, X_{[d]}]\} \rangle[Y_h],$$

that is the inference rule can be performed in when  $\mathcal{F}$  are in  $\mathcal{U}$ .

**Example 9.** *Modul Ponens* For example, when for two formulas  $f, h \in \mathcal{F}$  we have  $f \models h$ , then when  $\kappa^{\{f\}}[Y_f] = \epsilon_1[Y_f]$  we have

$$\epsilon_1[Y_h] = \langle (f \Rightarrow h)[Y_f, Y_h], \kappa^{\{f\}}[Y_f] \rangle[Y_h],$$

that is entailment of  $h$  can be concluded using a single update.

When we have a Knowledge Base of horn clauses, we run Knowledge Propagation with each horn clause being a constraint core and a knowledge core for any variable. Algorithm 8 therefore resembles the forward chaining algorithm of propositional logics (see Figure 7.15 in Russell and Norvig (2021)). It is known, that forward chaining is complete for Horn Logic. Thus, the knowledge cores returned in that case by Algorithm 8 are complete for the Knowledge Base as a CSP.

#### 8.5 Discussion

**Remark 8** (Interpretation of Contractions in Logical Reasoning). *The coordinates of contracted boolean tensor networks describe whether the by the coordinate indexed world is a model of the Knowledge Base at hand. Contractions, which only leave a part variables open, store the counts of the world respecting conditions given by the choice of slices. When contracting without open variables, we thus get the total world count.*

*This is consistent with the probabilistic interpretation of contractions, when applying the frequentist interpretation of probability and defining normed worldcounts as probabilities.*

**Remark 9.** *Tradeoff between generality and efficiency* While generic entailment decision algorithms (those by the full network) can decide any entailment, local algorithms as presented here can only perform some, but therefore more effectively as operating batchwise (dynamically deciding entailment for many leg variables). This is a typical phenomenon in logical reasoning and related to decidability.

Local contraction approaches in inference, especially when orchestrated by a Knowledge Propagation algorithm, mimik inference rules in syntax-based prove approaches.

## Part II

# Neuro-Symbolic Approaches

## 9 Introduction into Part II

After having explored tensor approach to the classical logical and probabilistic reasoning in Part I, we now turn to neuro-symbolic applications of tensor networks in artificial intelligence.

We start in Chapter 10 with the design of explainable and efficient learning architectures, which we frame formula selecting networks. These architectures will be utilized in the representation and reasoning on graphical models in Chapter 11 and Chapter 12.

While the focus is on propositional logic as an explainable framework in machine learning, we show extensions towards more expressive first-order logics in Chapter 14.

Besides that, probabilistic guarantees on the success of the learning problems are derived in Chapter 13.

## 10 Formula Selecting Networks

In this chapter we will investigate efficient schemes to represent collections of formulas with similar structure in one tensor network.

**Definition 52.** Given a set of  $p$  formulas  $\{f_l : l \in [p]\}$ , the formula selecting map is the map

$$\mathcal{H} : \bigtimes_{k \in [d]} [2] \rightarrow \bigtimes_{l \in p} [2]$$

defined for  $x_{[d]} \in \bigtimes_{k \in [d]} [2]$  as

$$\mathcal{H}(x_{[d]}) = \bigtimes_{l \in p} f_l[x_{[d]}] .$$

A tensor representation of a formula selecting map is provided by the selection encoding (see Def. 15)

$$\sigma^{\mathcal{H}}[X_{[d]}, L]$$

where the selection variable  $L$  takes values in  $[p]$  and selects specific formulas in the set  $\{f_l : l \in [p]\}$ . By definition, we have for any  $x_{[d]} \in \bigtimes_{k \in [d]} [2]$  and  $l \in [p]$

$$\sigma^{\mathcal{H}}[X_{[d]} = x_{[d]}, L = l] = f_l[X_{[d]} = x_0, \dots, x_{d-1}] .$$

This selection encoding is thus the sum

$$\sigma^{\mathcal{H}}[X_{[d]}, L] = \sum_{l \in [p]} f_l[X_{[d]}] \otimes \epsilon_l[L] .$$

Such a representation scheme requires linear resources in the number of formulas. We will show in the following, that we can exploit common structure in formulas to drastically reduce this resource consumption. Central to these sparse representation scheme are basis encodings  $\beta^{\mathcal{H}}$  of the selection encodings  $\sigma^{\mathcal{H}}$ , which we depict in Figure 23.

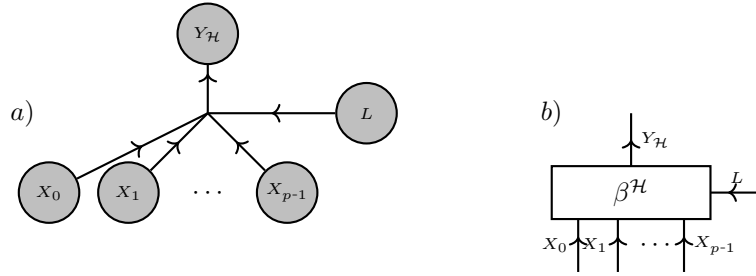


Figure 23: Representation of the basis encoding  $\beta^{\mathcal{H}}$  to a selection encoded formula selecting map as an a) Factor of a Graphical Model with a selection variable  $L$  and a computed variable  $Y_{\mathcal{H}}$ . b) Decorating Tensor Core with selection variable corresponding with an additional axis.

### 10.1 Construction schemes

Let us now investigate efficient schemes to define sets of formulas to be used in the definition of  $\mathcal{H}$ . We will motivate the folding of the selection variable into multiple selection variables by compositions of selection maps.

#### 10.1.1 Connective Selecting Maps

We represent choices over connectives with a fixed number of arguments by adding a selection variable to the cores and defining each slice by a candidate connective.

**Definition 53.** Let  $\{\circ_0, \dots, \circ_{p_C-1}\}$  be a set of connectives with  $d$  arguments. The associated connective selection map is

$$\mathcal{H}_C : \bigtimes_{k \in [d]} [2] \rightarrow \bigtimes_{l \in [p_C]} [2]$$

defined for  $x_{[d]} \in \bigtimes_{k \in [d]} [2]$  as

$$\mathcal{H}_C(x_{[d]}) = \bigtimes_{l \in [p_C]} \circ_l[X_{[d]} = x_{[d]}] .$$



### 10.1.2 Variable Selecting Maps

Choices of connectives can be combined with selections of variables assigned building the arguments of a connective. To this end, we introduce variable selecting maps.

**Definition 54.** *The selection of one out of  $p$  variables in a list  $X_{[p]}$  is done by variable selecting maps*

$$\mathcal{H}_V : \bigtimes_{l \in [p]} [2] \rightarrow \bigtimes_{l \in [p]} [2]. \quad (21)$$

defined as the identity map

$$\mathcal{H}_V (X_{[p]}) = X_{[p]}.$$

The selection encoding of the variable selecting map is the tensor  $\sigma^{\mathcal{H}_V} [X_{[p]}, L_V]$

$$\sigma^{\mathcal{H}_V} [X_0 = x_0, \dots, X_{p-1} = x_{p-1}, L_V = l_V] = x_{l_V}.$$

Selection encodings of variable selecting maps appear in the literature as multiplex gates (see e.g. Definition 5.3 in Koller and Friedman (2009)).

The basis encoding of the variable selection map has a decomposition

$$\beta^{\mathcal{H}_V} [Y_V, X_{[p_V]}] = \sum_{l_V \in [p_V]} \beta^{X_{l_V}} [Y_V, X_{l_V}] \otimes \epsilon_{l_V} [L_V].$$

This structure is exploited in the next theorem to derive a tensor network decomposition of  $\beta^{\mathcal{H}_V}$ .

**Theorem 56** (Decomposition of Variable Selecting Maps). *Given a list  $X_{[p_V]}$  of variables, we define for each  $l_V \in [p_V]$  the tensors*

$$\tau^{l_V} [X_{l_V}, L_V] = \delta [Y_V, X_{l_V}] \otimes \epsilon_{l_V} [L_V] + \mathbb{I} [Y_V, X_{l_V}] \otimes (\mathbb{I} [L_V] - \epsilon_{l_V} [L_V]).$$

Then we have (see Figure 24)

$$\beta^{\mathcal{H}_V} [Y_V, X_{[p]}, L_V] = \langle \{ \tau^{l_V} [Y_V, X_{l_V}, L_V] : l_V \in [p_V] \} \rangle [Y_V, X_{[p]}, L_V].$$

*Proof.* We show the equivalence of the tensors on an arbitrary coordinates. For  $\tilde{l}_V \in [p_V]$ ,  $Y_V \in [2]$  and  $x_{[p_V]} \in \bigtimes_{k \in [p_V]} [2]$  we have

$$\begin{aligned} & \langle \{ \tau^{l_V} [Y_V, X_{l_V}, L_V] : l_V \in [p_V] \} \rangle [Y_V = y_V, X_{[p]} = x_{[p]}, L_V = \tilde{l}_V] \\ &= \prod_{l_V \in [p_V]} \tau^{l_V} [Y_V = y_V, X_{l_V} = x_{l_V}, L_V = \tilde{l}_V] \\ &= \tau^{\tilde{l}_V} [Y_V = y_V, X_{l_V} = x_{l_V}, L_V = \tilde{l}_V] \\ &= \begin{cases} 1 & \text{if } y_V = x_{l_V} \\ 0 & \text{else} \end{cases} \\ &= \beta^{\mathcal{H}_V} [Y_V = y_V, X_{[p]} = x_{[p]}, L_V = \tilde{l}_V] \end{aligned}$$

In the second equality, we used that the tensor  $\tau^{l_V}$  have coordinates 1 whenever  $\tilde{l}_V \neq l_V$ . □

The decomposition provided by The. 56 is in a CP format (see Chapter 18). The introduced tensors  $\tau^{l_V}$  are Boolean, but not directed and therefore basis encodings of relations but not functions (see Chapter 17).

## 10.2 State Selecting Maps

When the variables to be selected are the atomization variables to the same categorical variable (see Sect. 8.3.2), one can avoid the instantiation of all atomization cores and instead represent the variable selecting map using the categorical variable only. To show this we introduce the state selecting map to a categorical variable  $X$ .

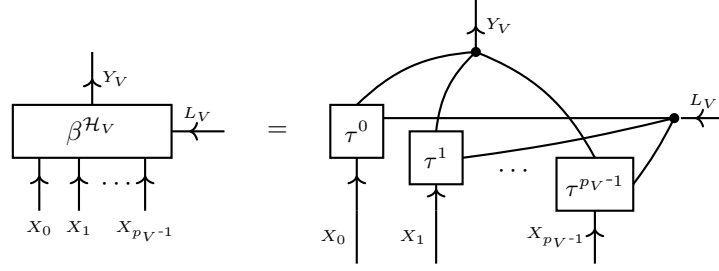


Figure 24: Decomposition of the basis encoding of a variable selecting tensor into a network of tensors defined in The. 56. The decomposition is in a CP format (see Chapter 18).

**Definition 55.** Given a categorical variable  $X_S$  with dimension  $m_S$  and a selection variable  $L_S$  with dimension  $p_S = m_S$  the state selecting map is the map

$$\mathcal{H}_S : [m_S] \rightarrow \bigtimes_{k \in [m_S]} [2]$$

defined for  $x_S \in [m_S]$  as

$$\mathcal{H}_S(x_S) = \epsilon_{x_S}[L_S] .$$

The selection encoding of the state selecting map coincides with the dirac delta tensor, that is for  $x_S \in [m_S]$  and  $l_S \in [p_S]$  we have

$$\sigma^{\mathcal{H}_S}[X_S, L_S] = \begin{cases} 1 & \text{if } x = l_S \\ 0 & \text{else} \end{cases} .$$

The relation of the variable selecting map and the state selecting map is shown in the next lemma.

**Lemma 14.** Let  $X_{[p]}$  be a collection of atomization variables to a categorical variable  $X$  taking values in  $[p]$ . Then we have for

$$\langle \beta^Z[X_{[p]}, X], \sigma^{\mathcal{H}_V}[X_{[p]}, L] \rangle [X, L] = \sigma^{\mathcal{H}_S}[X, L] .$$

*Proof.* For each  $x, l \in [p]$  we have that

$$\begin{aligned} \langle \beta^Z[X_{[p]}, X], \sigma^{\mathcal{H}_V}[X_{[p]}, L] \rangle [X = x, L = l] &= \langle \beta^Z[X_{[p]}, X = x], \sigma^{\mathcal{H}_V}[X_{[p]}, L = l] \rangle [\emptyset] \\ &= \begin{cases} 1 & \text{if } x = l \\ 0 & \text{else} \end{cases} \end{aligned}$$

which is thus equal to  $\sigma^{\mathcal{H}_S}[X = x, L = l]$ . □

Lem. 14 shows that when the variables to be selected are the atomization variables to a categorical variable, the state selecting map can thus be used instead of the variable selecting map. The state selecting map has the advantage, that the instantiation of the tensor  $\beta^Z[X_{[p]}, X]$  enforcing the categorical constraint can be avoided.

### 10.3 Composition of formula selecting maps

We will now parametrize the sets  $\mathcal{F}$  with additional indices and define formula selector maps subsuming all formulas. To handle large sets of formulas, we further fold the selection variable into tuples of selection variables.

**Definition 56.** Let there be a formula  $f_{l_0, \dots, l_{n-1}}$  for each index tuple in  $l_0, \dots, l_{n-1} \in \bigtimes_{s \in [n]} [p_s]$ , where  $n, p_0, \dots, p_{n-1} \in \mathbb{N}$ . The folded formula selection map (see Figure 25) is the map

$$\mathcal{H} : \bigtimes_{k \in [d]} [2] \rightarrow \bigtimes_{l_0 \in [p_0]} \cdots \bigtimes_{l_{n-1} \in [p_{n-1}]} [2]$$

defined as

$$\mathcal{H}(x_{[d]}) = (f_{l_{[n]}}[X_{[d]} = x_{[d]}])_{p_0, \dots, p_{n-1} \in \bigtimes_{s \in [n]} [p_s]} .$$

Folding the selection variable into multiple selection variables is especially useful to find efficient decomposition schemes of the formula selecting maps. In the reminder of this section we will provide an example, where each selection variable is constructed to parameterize a local change to the formula with the result that the basis encoding of the global formula selecting map decomposes into local formula selecting maps.

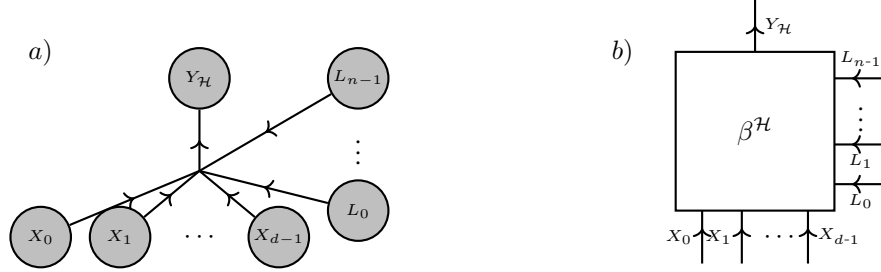


Figure 25: Basis encoding of the folded map  $\mathcal{H}$ .

### 10.3.1 Formula Selecting Neuron

The folding of the selection variable is motivated by the composition of selection maps. We call the composition of a connective selection (see Def. 53) with variable selection maps (see Def. 54) for each argument a formula selecting neuron.

**Definition 57.** Given an order  $n \in \mathbb{N}$  let there be a connective selector  $L_\circ$  selecting connectives of order  $n$  and let  $\mathcal{H}_{V,0}, \dots, \mathcal{H}_{V,n-1}$  be a collection of variable selectors. The corresponding logical neuron is the map

$$\mathcal{N} : \bigtimes_{k \in [d]} [2] \rightarrow \bigtimes_{l_C \in [p_C]} \bigtimes_{l_0 \in [p_0]} \dots \bigtimes_{l_{n-1} \in [p_{n-1}]} [2]$$

defined for  $x_{[d]} \in \bigtimes_{k \in [d]} [2]$  by

$$\mathcal{N}(x_{[d]}) = \left( \circ_{l_C}(X_{l_0}, \dots, X_{l_{n-1}}) \right)_{l_C \in [p_C], l_0 \in [p_0], \dots, l_{n-1} \in [p_{n-1}]}$$

Each neuron has a tensor network decomposition by a connective selector tensor and a variable selector tensor network for each argument, as we state in the next theorem.

**Theorem 57.** Decomposition of formula selecting neurons Let  $\mathcal{N}$  a logical neuron, defined for a connective selector  $L_\circ$  and variable selectors  $\mathcal{H}_{V,0}, \dots, \mathcal{H}_{V,n-1}$ . Then we have (see Figure 26 for the example of  $n = 2$ ):

$$\begin{aligned} & \beta^{\mathcal{N}} [Y_{\mathcal{N}}, X_{[d]}, L_C, L_{V,0}, \dots, L_{V,n-1}] \\ &= \langle \{ \beta^{\mathcal{H}_C} [Y_{\mathcal{N}}, Y_{V,0}, \dots, Y_{V,n-1}], \\ & \quad \beta^{\mathcal{H}_{V,0}} [Y_{V,0}, X_{[d]}, L_{V,0}], \dots, \beta^{\mathcal{H}_{V,n-1}} [Y_{V,n-1}, X_{[d]}, L_{V,n-1}] \} \rangle [Y_{\mathcal{N}}, X_{[d]}, L_C, L_{V,0}, \dots, L_{V,n-1}]. \end{aligned}$$

*Proof.* By composition The. 108. □

### 10.3.2 Formula Selecting Neural Network

Single neurons have a limited expressivity, since for each choice of the selection variables they can just express single connectives acting on atomic variables. The expressivity is extended to all propositional formulas, when allowing for networks of neurons, which can select each others as input arguments.

**Definition 58.** An architecture graph  $\mathcal{G}^A = (\mathcal{V}^A, \mathcal{E}^A)$  is an acyclic directed hypergraph with nodes appearing at most once as outgoing nodes. Nodes appearing only as outgoing nodes are input neurons and are labeled by  $\mathcal{A}^{\text{in}}$  and nodes not appearing as outgoing nodes are the output neurons in the set  $\mathcal{A}^{\text{out}}$  (see Figure 27 for an example).

Given an architecture graph  $\mathcal{G}^A = (\mathcal{V}^A, \mathcal{E}^A)$ , a formula selecting neural network  $\sigma^A$  is a tensor network of logical neurons at each  $\mathcal{N} \in \mathcal{V}^A / \mathcal{A}^{\text{in}}$ , such that each neuron depends on variables  $X_{\text{Pa}(\mathcal{N})}$  and on selection variables  $L_{\mathcal{N}}$ . The collection of all selection variable is notated by  $L_A$ .

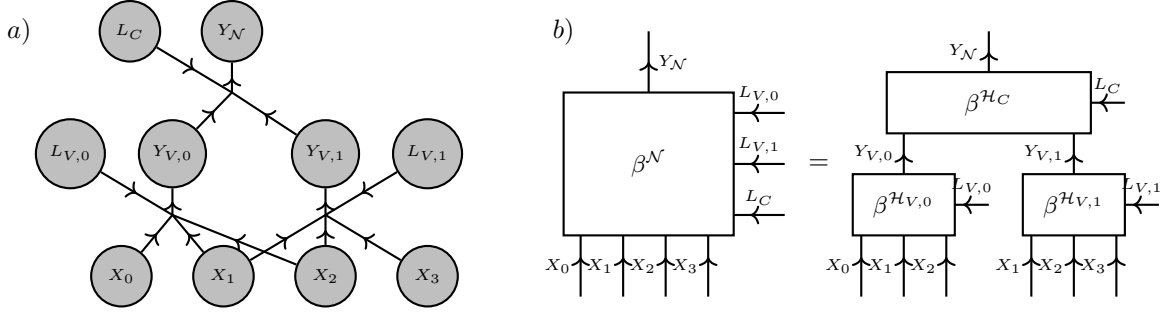


Figure 26: Example of a logical neuron  $\mathcal{N}$  of order  $n = 2$ . a) Selection and categorical variables and their interdependencies visualized in a hypergraph. b) Basis encoding of the logical neuron and tensor network decomposition into variable selecting and connective selecting tensors.

The activation tensor of each neuron  $\mathcal{N} \in \mathcal{V}^{\mathcal{A}}/\mathcal{A}^{\text{in}}$  is

$$\sigma^{\mathcal{N}}[X_{\mathcal{A}^{\text{in}}}, L_{\mathcal{A}}] = \left\langle \{\beta^{\tilde{\mathcal{N}}} : \tilde{\mathcal{N}} \in \mathcal{V}^{\mathcal{A}}/\mathcal{A}^{\text{in}}\} \cup \{\epsilon_1[Y_{\mathcal{N}}]\} \right\rangle [X_{\mathcal{A}^{\text{in}}}, L_{\mathcal{A}}] .$$

The activation tensor of the formula selecting neural network is the contraction

$$\sigma^{\mathcal{A}}[X_{\mathcal{A}^{\text{in}}}, L_{\mathcal{A}}] = \left\langle \{\beta^{\sigma^{\mathcal{N}}} [Y_{\mathcal{N}}, X_{\text{Pa}(\mathcal{N})}, L_{\mathcal{A}}] : \mathcal{N} \in \mathcal{V}^{\mathcal{A}}/\mathcal{A}^{\text{in}}\} \cup \{\epsilon_1[Y_{\mathcal{N}}] : \mathcal{N} \in \mathcal{A}^{\text{out}}\} \right\rangle [X_{\mathcal{A}^{\text{in}}}, L_{\mathcal{A}}] .$$

The expressivity of a formula selecting neural network  $\sigma^{\mathcal{A}}$  is the formula set

$$\mathcal{F}_{\mathcal{A}} = \left\{ \sigma^{\mathcal{A}}[X_{\mathcal{A}^{\text{in}}}, L_{\mathcal{A}} = l_{\mathcal{A}}] : l_{\mathcal{A}} \in \prod_{s \in [n]} [p_s] \right\} .$$

The activation tensor of each neuron depends in general on the activation tensor of its ancestor neurons with respect to the directed graph  $\mathcal{G}^{\mathcal{A}}$ , and thus inherits the selection variables.

We notice that the architecture graph is a scheme to construct the variable dependency graph of the tensor network  $\mathcal{F}_{\mathcal{A}}$ . To this end, we replace each neuron  $\mathcal{N} \in \mathcal{V}^{\mathcal{A}}/\mathcal{A}^{\text{in}}$  by an output variable  $Y_{\mathcal{N}}$  and further add selection variables  $L_{\mathcal{N}}$  to the directed edges, that is to each directed hyperedge  $(\{\mathcal{N}\}, \text{Pa}(\mathcal{N})) \in \mathcal{E}^{\mathcal{A}}$  we construct a directed hyperedge  $(\{Y_{\mathcal{N}}\}, X_{\text{Pa}(\mathcal{N})} \cup L_{\mathcal{N}})$ .

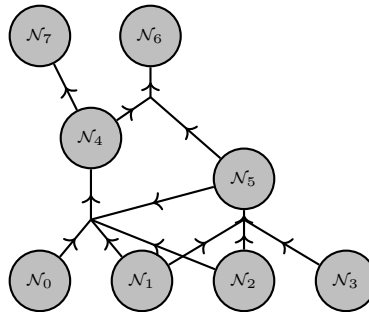


Figure 27: Example of an architecture graph  $\mathcal{G}^{\mathcal{A}}$  with input neurons  $\mathcal{A}^{\text{in}} = \{\mathcal{N}_0, \mathcal{N}_1, \mathcal{N}_2, \mathcal{N}_3\}$  and output neurons  $\mathcal{A}^{\text{out}} = \{\mathcal{N}_6, \mathcal{N}_7\}$

**Theorem 58.** Given fixed selection variables  $L_{\mathcal{A}}$ , the formula selecting neural network is the conjunction of output neurons, that is

$$\sigma^{\mathcal{A}}[X_{\mathcal{A}^{\text{in}}}, L_{\mathcal{A}}] = \left\langle \{\sigma^{\mathcal{N}}[X_{\mathcal{A}^{\text{in}}}, L_{\mathcal{A}}] : \mathcal{N} \in \mathcal{A}^{\text{out}}\} \right\rangle [X_{\mathcal{A}^{\text{in}}}, L_{\mathcal{A}}] .$$

*Proof.* By effective calculus (see The. 116), we have

$$\left\langle \beta^{\wedge} [X_{\wedge}, X_{[d]}], \epsilon_1[X_{\wedge}] \right\rangle [X_{[d]}] = \bigotimes_{k \in [d]} \epsilon_1[X_k]$$

and thus

$$\sigma^{\mathcal{A}} [X_{\mathcal{A}^{\text{in}}}, L_{\mathcal{A}}] = \langle \{\beta^{\mathcal{N}} : \mathcal{N} \in \mathcal{V}^{\mathcal{A}} / \mathcal{A}^{\text{in}}\} \cup \{\beta^{\wedge} [X_{\wedge}, Y_{\mathcal{N}} : \mathcal{N} \in \mathcal{A}^{\text{out}}], \epsilon_1 [X_{\wedge}]\} \rangle [X_{\mathcal{A}^{\text{in}}}, L_{\mathcal{A}}] .$$

□

By the commutation of contractions, we can further use The. 57 to decompose each tensor  $\beta^{\mathcal{N}}$  into connective and variable selecting components to get a sparse representation of a formula selecting neural network  $\sigma^{\mathcal{A}}$ .

## 10.4 Application of Formula Selecting Networks

There are two main applications of formula selecting networks. First, when contracting the selection variables with a weight tensor we get a weighted sum of the parametrized formulas. Second, when contracting the categorical variables with a distribution or a knowledge base, we get a tensor storing the satisfaction rates respectively the world counts of the parametrized formulas.

### 10.4.1 Representation of selection encodings

The main application of formula selecting networks in this work is the efficient representation of selection encodings. This will be exploited in the sparse representation of exponential families by energies and in structure learning. In the next lemma we will show the correspondence of formula selecting networks and selection encodings.

**Lemma 15.** *Given a set  $\{f_{l_0, \dots, l_{n-1}} : l_0, \dots, l_{n-1} \in \times_{s \in [n]} [p_s]\}$  of propositional formulas we define the statistic*

$$\mathcal{F} : x_0, \dots, x_{d-1} \rightarrow (f_{l_0, \dots, l_{n-1}}(x_0, \dots, x_{d-1}))_{l_0, \dots, l_{n-1}} .$$

*and the formula selecting map*

$$\mathcal{H} : x_0, \dots, x_{d-1}, l_0, \dots, l_{n-1} \rightarrow f_{l_0, \dots, l_{n-1}}(x_0, \dots, x_{d-1}) .$$

*Then*

$$\sigma^{\mathcal{F}} [X_{[d]}, L_{[n]}] = \mathcal{H} [X_{[d]}, L_{[n]}] .$$

*Proof.* For any indices  $l_{[n]} \in \times_{s \in [n]} [p_s]$  and  $x_{[d]} \in \times_{k \in [d]} [2]$  we have

$$\sigma^{\mathcal{F}} [X_{[d]} = x_{[d]}, L_{[n]} = l_{[n]}] = f_{l_0, \dots, l_{n-1}}(x_0, \dots, x_{d-1}) = \mathcal{H} [X_{[d]} = x_{[d]}, L_{[n]} = l_{[n]}] .$$

□

Technically, basis encodings have been exploited to derive decompositions based on basis calculus. Selection encodings on the other hand enable the application of formula selecting networks as superpositions of formulas.

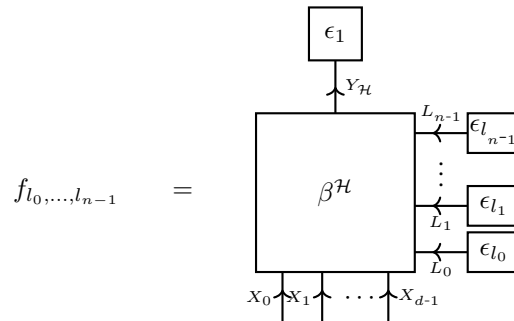
### 10.4.2 Efficient Representation of Formulas

Formula Selecting Neural Networks are means to represent exponentially many formulas with linear (in sum of candidates list lengths) storage. Their contraction with probability tensor networks, is thus a batchwise evaluation of exponentially many formulas. This is possible due to redundancies in logical calculus due to modular combinations of subformulas.

We can retrieve specific formulas by slicing the selection variables, i.e. for  $l_0, \dots, l_{n-1}$  we have

$$f_{l_0, \dots, l_{n-1}} [X_{[d]}] = \mathcal{H} (X_{[d]}, L = l_0, \dots, l_{n-1}) .$$

In a tensor network diagram we depict this by



Another perspective on the efficient formula evaluation by selection tensor networks is dynamic computing. Evaluating a formula requires evaluations of its subformulas, which are done by subcontractions and saved for different subformulas due to the additional selection legs.

However, we need to avoid contracting the tensor with leaving all selection legs open, since this would require exponential storage demand. We can avoid this storage bottleneck by contraction of parameter cores  $\theta$  with efficient network decompositions along the selection variables.

In Gibbs Sampling (Algorithm 2), one can use the energy-based approach to queries The. 22, and contract basis vectors on all but one selection variables.

#### 10.4.3 Batch contraction of parametrized formulas

Given a set  $\mathcal{F}$  of formulas, we build a formula selecting network parametrizing the formulas. The contraction

$$\langle \tau^{\mathcal{G}}, \mathcal{H} \rangle [L_{[n]}]$$

is a tensor containing the contractions of the formulas  $f_{l_{[n]}}$  with an arbitrary tensor network  $\tau^{\mathcal{G}}$  as

$$\langle \tau^{\mathcal{G}}, f_{l_{[n]}} \rangle [\emptyset] = \langle \tau^{\mathcal{G}}, \mathcal{H} \rangle [L_{[n]} = l_{[n]}] .$$

#### 10.4.4 Average contraction of parametrized formulas

We show in the next two examples, how a full contraction of the formula selecting map with a probability distribution or a knowledge base can be interpreted.

**Example 10** (Average satisfaction of formulas). *The average of the formula satisfactions in  $\mathcal{F}$  given a probability tensor  $\mathbb{P}$  is*

$$\frac{1}{\prod_{s \in [n]} p_s} \cdot \langle \mathbb{P}, \sigma^{\mathcal{F}} \rangle [\emptyset] .$$

**Example 11** (Deciding whether any formula is not contradicted). *For example: We want to decide, whether there is a formula in  $\mathcal{F}$  not contradicted by a Knowledge base  $\mathcal{KB}$ . This is the case if and only if*

$$\langle \mathcal{KB}, \sigma^{\mathcal{F}} \rangle [\emptyset] = 0 .$$

*We use Lem. 15 to get that  $\sigma^{\mathcal{F}} = \mathcal{H}$ . When the formulas are representable in a folded scheme, we find tensor network decompositions of  $\mathcal{H}$  and exploit them along efficient representations of  $\mathcal{KB}$  in an efficient calculation of  $\langle \mathcal{KB}, \sigma^{\mathcal{F}} \rangle [\emptyset]$ . This is further equal to*

$$\mathcal{KB} \models \neg \left( \bigvee_{f \in \mathcal{F}} f \right) .$$

### 10.5 Examples of formula selecting neural networks

#### 10.5.1 Correlation

For example (see Figure 28) consider the logical neuron with single activation candidate  $\{\wedge\}$  and two variable selectors selecting  $d$  atomic variables  $X_{[d]}$ . The expressivity of this network is the set of all conjunctions of the atoms

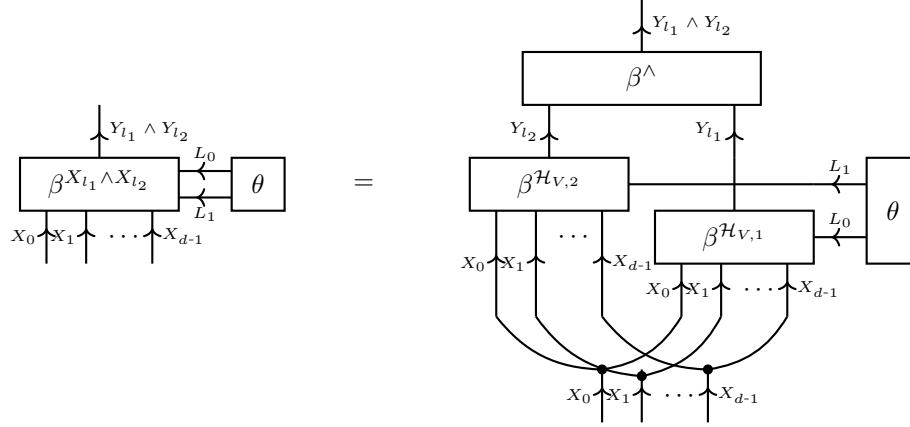
$$\{X_k \wedge X_l : k, l \in [d]\}$$

Contracting with a probability distribution, we use the tensor

$$\tau [L_{V,0}, L_{V,1}] = \langle \sigma^{\mathcal{A}} \rangle [L_{V,0}, L_{V,1}]$$

to read of covariances as

$$\text{Cov}(X_k, X_l) = \tau [L_{V,0} = k, L_{V,1} = l] - \tau [L_{V,0} = k, L_{V,1} = k] \cdot \tau [L_{V,0} = l, L_{V,1} = l] .$$

Figure 28: Superposition of the encoded formulas  $\beta^{X_{l_1} \wedge X_{l_2}}$  with weight  $\theta_{l_1 l_2}$ 

### 10.5.2 Conjunctive and Disjunctive Normal Forms

We can represent any propositional knowledge base by the following scheme: Literal selecting neurons are logical neurons with connective identity/negation (selecting positive/negative literal) and selecting neurons select for each an atom. The single output neuron represents the disjunction, respectively the conjunction, combining the literal selecting neurons. The number of neurons defined by the maximal clause size plus one. Smaller clauses can be covered when adding False as a possible choice (The respective neuron has to choose the identity, otherwise the full clause will be trivial). This architecture will be discussed in more detail in Chapter 19 as CP selecting networks. The parameter core is in the basis CP format and each slice selects a clause of the knowledge base. In combination with polynomial decompositions, which will be provided in Chapter 11, one can exploit this architecture to find sparse formula decompositions.

**Remark 10** (Minterms and Maxterms). *All minterms and maxterms can be represented by a two layer selection tensor networks without variable selection in two layers. The bottom layer has an  $\neg/\text{Id}$  connective selection neuron to each atom and the upper layer consists of a single dary conjunction.*

### 10.6 Extension to variables of larger dimension

While we here restricted on boolean variables, formula selecting networks can be extended to variables of larger cardinality.

- Connective selecting tensors: Can encode arbitrary functions  $h_l$  of discrete variables, but need  $X_{\mathcal{H}_C}$  to be an enumeration of the states, in particular to be of dimension

$$m_{\mathcal{H}_C} = |\cup_{l \in [p]} \text{im}(h_l)|.$$

- Variable selecting tensors can be understood as specific cases of connective selecting tensors and can thus also be generalized in a straight forward manner by

$$m_{\mathcal{H}_C} = |\cup_{l \in [p]} \text{im}(h_l)|.$$

- State selecting tensors are directly defined for larger dimensions

An example of such a more generic usage is a discretization scheme for continuous neurons.

**Example 12** (Discretization of a continuous neuron). *Let there be a neuron by a map of weight vectors and input vectors to  $\mathbb{R}$ , that is*

$$\sigma(w, x) : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}.$$

*We restrict the weights to a subset  $\mathcal{U}^{weight} \subset \mathbb{R}^d$  and the input vectors to  $\mathcal{U}^x \subset \mathbb{R}^d$ , It follows that*

$$|\text{im}(\sigma|_{\mathcal{U}^{weight} \times \mathcal{U}^x})| \leq |\mathcal{U}^{weight}| \cdot |\mathcal{U}^x|.$$

*To discretize the neuron, we use the subset encoding scheme of Def. 80 and define enumeration variables  $O_{weight}$ ,  $O_x$  and  $O_\sigma$  enumerating  $\mathcal{U}^{weight}$ ,  $\mathcal{U}^x$  and  $\text{im}(\sigma|_{\mathcal{U}^{weight} \times \mathcal{U}^x})$ , which are accompanied by respective index interpretation*

functions. Then the basis encoding of the discretized neuron is

$$\beta^\sigma [O_\sigma, O_{weight}, O_x] \cdot = \sum_{O_{weight} \in [r_{weight}], O_x \in [r_x]} \epsilon_{I_\sigma(I_{weight}(O_{weight}), I_x(O_x))}(\sigma) O_\sigma \otimes \epsilon_{O_{weight}} [O_{weight}] \otimes \epsilon_{O_x} [O_x] \cdot$$

If the neuron is of the form

$$\sigma(w, x) = \psi\left(\sum_i w_i \cdot x_i\right)$$

a decomposition into multiplication at each coordinate and summation of the results, with basis encodings for each, can be done. *Here the index interpretation variables are split into a selection enumerated by  $i$  and each variable gets assigned to single cores in the decomposition.*

## 11 Logical Network Representation

Logic networks are graphical models with an interpretation by propositional logics. We first distinguish between Markov Logic Networks, which are an approach to soft logics in the framework of exponential families, and Hard Logic Networks, which correspond with propositional knowledge bases. Then we exploit non-trivial boolean base measures to unify both approaches by Hybrid Logic Networks, which are itself in exponential families.

### 11.1 Markov Logic Networks

Markov Logic Networks exploit the efficiency and interpretability of logical calculus as well as the expressivity of graphical models.

#### 11.1.1 Markov Logic Networks as Exponential Families

We introduce Markov Logic Networks in the formalism of exponential families (see Sect. 5.2).

**Definition 59** (Markov Logic Networks). *Markov Logic Networks are exponential families  $\Gamma^{\mathcal{F}, \mathbb{I}}$  with sufficient statistics by functions*

$$\mathcal{F} : \prod_{k \in [d]} [2] \rightarrow \prod_{f \in \mathcal{F}} [2] \subset \mathbb{R}^{|\mathcal{F}|}$$

*defined coordinatewise by propositional formulas  $f \in \mathcal{F}$ .*

Since the image of each coordinate  $\mathcal{S}_l$  is contained in  $\{0, 1\}$ , each is a propositional formulas (see Def. 43). Thus, any exponential family of distributions of  $\prod_{k \in [d]} [2]$ , such that  $\text{im}(\mathcal{S}_l) \subset \{0, 1\}$  for all  $l \in [p]$  is a set of Markov Logic Networks with fixed formulas.

The sufficient statistics consistent in a map  $\mathcal{F}$  of formulas brings the following advantages:

- Numerical Advantage: The sufficient statistics is decomposable into logical connectives. If the formulas are sparse (in the sense of limited number of connectives necessary in their representation), this gives rise to efficient tensor network decompositions of the basis encoding.
- Statistical Advantage: Since each formula is Boolean valued, the coordinates of the sufficient statistic are Bernoulli variables. Due to their boundedness, they and their averages (by Hoeffdings inequality) are sub-Gaussian variables with favorable concentration properties (absence of heavy tails).

**Remark 11** (Alternative Definitions). *We here defined MLNs on propositional logic, while originally they are defined in FOL. The relation of both frameworks will be discussed further in Chapter 14.*

#### 11.1.2 Tensor Network Representation

Based on the previous discussion on the representation of exponential families by tensor networks in Sect. 5.2 we now derive a representation for Markov Logic Networks.

##### Basis encodings for distributions

**Theorem 59** (Basis encodings for Markov Logic Networks). *A Markov Logic Network to a set of formulas  $\mathcal{F} = \{f_l : l \in [p]\}$  is represented as*

$$\mathbb{P}^{\mathcal{F}, \theta}[X_{[d]}] = \left\langle \{\beta^{f_l}[Y_l, X_{[d]}] : l \in [p]\} \cup \{\alpha^{f_l, \theta[L=l]}[Y_l] : l \in [p]\} \right\rangle [X_{[d]} | \emptyset]$$



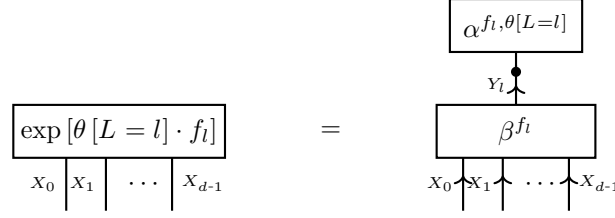


Figure 29: Factor of a Markov Logic Network to a formula  $f_l$ , represented as the contraction of a computation core  $\beta^{f_l}$  and an activation core  $\alpha^{f_l, \theta[L=l]}$ . While the computation core  $\beta^{f_l}$  prepares based on basis calculus a categorical variable representing the value of the statistic formula  $f_l$  dependent on assignments to the distributed variables, the activation core multiplies an exponential weight to coordinates satisfying the formula.

where we denote for each  $l \in [p]$  an activation core

$$\alpha^{f_l, \theta[L=l]} [Y_l = y_l] = \begin{cases} 1 & \text{for } y_l = 0 \\ \exp [\theta [L = l]] & \text{for } y_l = 1 \end{cases}.$$

*Proof.* Markov Logic Networks are exponential families, which base measure is trivial and which statistic consist of boolean features. We apply the tensor network decomposition of more generic exponential families The. 11 to this case and get

$$\mathbb{P}^{\mathcal{F}, \theta} [X_{[d]}] = \left\langle \{ \mathbb{I} [X_{[d]}] \} \cup \{ \beta^{S_l} [Y_l, X_{[d]}] : l \in [p] \} \cup \{ \alpha^{f_l, \theta[L=l]} [Y_l] : l \in [p] \} \right\rangle [X_{[d]} | \emptyset].$$

While the base measure tensor is trivial, it can be ignored in the contraction. Since the image of each feature  $f_l$  is in  $[2]$ , we choose the index interpretation function by the identity  $I : [2] \rightarrow \{0, 1\}$  and get

$$\begin{aligned} \alpha^{f_l, \theta[L=l]} [Y_l = y_l] &= \exp [\theta [L = l] \cdot I_l(y_l)] = \exp [\theta [L = l] \cdot y_l] \\ &= \begin{cases} 1 & \text{for } y_l = 0 \\ \exp [\theta [L = l]] & \text{for } y_l = 1 \end{cases} \end{aligned}$$

□

The. 59 provides a decomposition of markov logic networks by a tensor network of computation cores  $\beta^{f_l}$  and accompanying activation cores  $\alpha^{f_l, \theta[L=l]}$ . Since the head variable  $Y_l$  appears exclusively in these pairs, we can contract each computation core with the corresponding activation core to get a factor, see Figure 29. With this we get the decomposition

$$\mathbb{P}^{\mathcal{F}, \theta} [X_{[d]}] = \left\langle \{ \exp [\theta [L = l] \cdot f_l [X_{[d]}]] : l \in [p] \} \right\rangle [X_{[d]} | \emptyset].$$

More precisely, this transformation of the decomposition holds by The. 131 to be shown in Chapter 20, stating that the contraction of computation and activation cores can be performed before the global contraction of the result.

While in the decomposition of The. 59 the basis encodings of the features carry all distributed variables  $X_{[d]}$ , we now seek towards sparser decompositions. To each  $l \in [p]$  we denote by  $\mathcal{V}^l$  the maximal subset of  $[d]$  such that there is a reduced function  $\tilde{f}_l : \mathcal{X} \rightarrow [2]$  with

$$f_l [X_{[d]}] = \left\langle \tilde{f}_l [X_{\mathcal{V}^l}] \right\rangle [X_{[d]}].$$

We often account for such situations of sparse formulas, when  $f_l$  has a syntactical decomposition involving only the atomic variables  $\mathcal{V}^l$ . As a consequence we have

$$\beta^{f_l} [X_{[d]}] = \beta^{\tilde{f}_l} [X_{\mathcal{V}^l}]$$

and

$$\mathbb{P}^{\mathcal{F}, \theta} [X_{[d]}] = \left\langle \{ \exp [\theta [L = l] \cdot \tilde{f}_l [X_{\mathcal{V}^l}]] : l \in [p] \} \right\rangle [X_{[d]} | \emptyset].$$

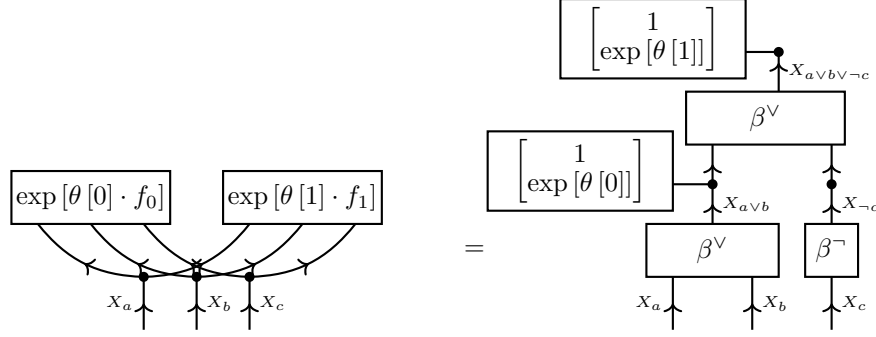


Figure 30: Example of a decomposed Markov Network representation of a Markov Logic Network with formulas  $\{f_0 = a \vee b, f_1 = a \vee b \vee \neg c\}$ . Since both formulas share the subformula  $a \vee b$ , their contracted factors have a representation by a connected tensor network.

Thus, any markov logic network has a sparse representation by a markov network on the graph

$$\mathcal{G}^{\mathcal{F}} = ([d], \{\mathcal{V}^l : l \in [p]\}).$$

This sparsity inducing mechanism is analogous to the decomposition of probability distributions based on conditional independence assumptions, when understanding each formula in the markov logic network as an introduced dependency among the affected variables  $\mathcal{V}^l$ .

A further sparsity introducing mechanism is through exploiting redundancy in the computation of  $f_l$ , when a decomposition of the feature is known. For the propositional formulas  $f_l$  this amounts to a syntactic representation of the formula as a composition of logical connectives is available (see Figure 30). In this case, we exploit the representation by tensor networks of the basis encodings (shown as The. 108 in Chapter 17) Note, that this decomposition scheme introduces further auxiliary variables  $Y$  with deterministic dependence on the distributed variables  $X_{[d]}$ . Such variables are often referred to as hidden.

We can further exploit common syntactical structure in the formulas  $f_l \in \mathcal{F}$  to reduce the number of basis encodings of connectives. This is the case, when the syntax graph of two or more formulas share a subgraph. In that case, the respective syntax graph needs to be represented only once and can be incorporated into the decomposition of all formulas, which share this subgraph. For an example see Figure 30, where the syntactical representation of the formula  $f_0$  is a subgraph of the syntactical representation of  $f_1$ .

To summarize, there are two sparsity mechanisms, originating from graphical models and propositional syntax, providing sparse representations of markov logic network:

- **Dependence Sparsity:** Formulas depend only on subsets of atoms. This exploits the main sparsity mechanism in graphical models, where factors in sparse representations depend only on a subset of variables.
- **Computation Sparsity:** When the features of an exponential family are compositions of smaller formulas, the computation core is decomposed into a tensor network of their basis encodings. This can be regarded as the main sparsity mechanism of propositional logics, where syntactical decompositions of formulas are exploited. Further, when the structure of the smaller formulas is shared among different features, the respective basis encodings need to be instantiated only once.

**Selection encodings for energy tensors** As for generic exponential families, we can represent markov logic networks in terms of their energy tensors

$$\phi^{\mathcal{F}, \theta} [X_{[d]}] = \sum_{l \in [p]} \theta [L = l] \cdot f_l [X_{[d]}] = \langle \sigma^{\mathcal{F}} [X_{[d]}, L], \theta [L] \rangle [X_{[d]}]. \quad (22)$$

The energy tensor provides an exponential representation of the distribution by

$$\mathbb{P}^{\mathcal{F}, \theta} [X_{[d]}] = \langle \exp [\phi^{\mathcal{F}, \theta}] \rangle [X_{[d]} | \emptyset]. \quad (23)$$

In case of a common structure of the formulas in a Markov Logic Network, formula selecting networks (see Chapter 10) can be applied to represent their energies. We represent the superposition of formulas as a contraction with a parameter tensor. Given a factored parametrization of formulas  $f_{l_0, \dots, l_{n-1}}$  with indices  $l_s$  we have the superposition by the network representation:

$$\sum_{l_{[n]} \in \times_{s \in [n]} [p_s]} \theta [L_{[n]} = l_{[n]}] f_{l_{[n]}} =$$

If the number of atoms and parameters gets large, it is important to represent the tensor  $f_{l_0, \dots, l_{n-1}}$  efficiently in tensor network format and avoid contractions. To avoid inefficiency issues, we also have to represent the parameter tensor  $\theta$  in a tensor network format to improve the variance of estimations (see Chapter 13) and provide efficient numerical algorithms.

However, when required to instantiate the probability distribution of a Markov Logic Network as a tensor network, we need to exponentiate and normate the energy tensor, a task for which basis encodings are required. For such tasks, contractions of formula selecting networks are not sufficient and each formula with a nonvanishing weight needs to be instantiated as a factor tensor of a Markov Network.

### 11.1.3 Expressivity

Based on Markov Logic Networks containing only maxterms and minterms (see Def. 45), we now show that any positive probability distribution has a representation by a markov logic network.

**Theorem 60.** *Let there be a positive probability distribution*

$$\mathbb{P} [X_{[d]}] \in \bigotimes_{k \in [d]} \mathbb{R}^2.$$

*Then the Markov Logic Network of minterms (see Def. 45)*

$$\mathcal{F}_{\wedge} = \{Z_{x_0, \dots, x_{d-1}}^{\wedge} : x_0, \dots, x_{d-1} \in \bigtimes_{k \in [d]} [2]\}$$

*with parameters*

$$\theta [L_0 = x_0, \dots, L_{d-1} = x_{d-1}] = \ln \mathbb{P} [X_0 = x_0, \dots, X_{d-1} = x_{d-1}]$$

*coincides with  $\mathbb{P} [X_{[d]}]$ .*

*Further, the Markov Logic Network of maxterms*

$$\mathcal{F}_{\vee} = \{Z_{x_0, \dots, x_{d-1}}^{\vee} : x_0, \dots, x_{d-1} \in \bigtimes_{k \in [d]} [2]\}$$

*with wparameters*

$$\theta [L_0 = x_0, \dots, L_{d-1} = x_{d-1}] = -\ln \mathbb{P} [X_0 = x_0, \dots, X_{d-1} = x_{d-1}]$$

*coincides with  $\mathbb{P} [X_{[d]}]$ .*

*Proof.* It suffices to show, that in both cases of choosing  $\mathcal{F}$  by minterms or maxterms with the respective parameters

$$\phi^{\mathcal{F}, \theta} = \ln \mathbb{P} [X_{[d]}]$$

and therefore

$$\mathbb{P}^{\mathcal{F}, \theta} [X_{[d]}] = \langle \exp [\phi^{\mathcal{F}, \theta}] \rangle [X_{[d]} | \emptyset] = \langle \exp [\phi^{\mathcal{F}, \theta}] \rangle [X_{[d]}] = \mathbb{P} [X_{[d]}].$$

In the case of minterms, we notice that for any  $x_0, \dots, x_{d-1} \in \bigtimes_{k \in [d]} [2]$

$$Z_{x_0, \dots, x_{d-1}}^{\wedge} [X_{[d]}] = \epsilon_{x_0, \dots, x_{d-1}} [X_{[d]}]$$

and thus with the weights in the claim

$$\sum_{x_0, \dots, x_{d-1} \in \bigtimes_{k \in [d]} [2]} (\ln \mathbb{P} [X_0 = x_0, \dots, X_{d-1} = x_{d-1}]) \cdot Z_{x_0, \dots, x_{d-1}}^{\wedge} [X_{[d]}] = \ln \mathbb{P} [X_{[d]}].$$

For the maxterms we have analogously

$$Z_{x_0, \dots, x_{d-1}}^\vee [X_{[d]}] = \mathbb{I} [X_{[d]}] - \epsilon_{x_0, \dots, x_{d-1}} [X_{[d]}]$$

and thus that the maximal clauses coincide with the one-hot encodings of respective states. We thus have

$$\begin{aligned} & \sum_{x_0, \dots, x_{d-1} \in \times_{k \in [d]} [2]} (-\ln \mathbb{P} [X_0 = x_0, \dots, X_{d-1} = x_{d-1}]) \cdot Z_{x_0, \dots, x_{d-1}}^\vee [X_{[d]}] \\ &= \left( \sum_{\nu_0 \subset [d]} (-\ln \mathbb{P} [X_0 = x_0, \dots, X_{d-1} = x_{d-1}]) \cdot \mathbb{I} [X_{[d]}] \right) \\ & \quad + \left( \sum_{\nu_0 \subset [d]} (\ln \mathbb{P} [X_0 = x_0, \dots, X_{d-1} = x_{d-1}]) \cdot \epsilon_{x_0, \dots, x_{d-1}} [X_{[d]}] \right) \\ &= \ln \mathbb{P} [X_{[d]}] + \lambda \cdot \mathbb{I} [X_{[d]}] , \end{aligned}$$

where  $\lambda$  is a constant. □

We note, that there are  $2^d$  maxterms and  $2^d$  minterms, which would have to be instantiated by basis encodings to get a tensor network decomposition. This large number of features originates from the generality of the representation scheme. As a fundamental tradeoff, efficient representations come at the expense of a smaller expressivity of the representation scheme.

Theorem 60 is the analogue in Markov Logic to Theorem 42, which shows that any binary tensor has a representation by a logical formula, to probability tensors. Here we require positive distributions for well-defined energy tensors.

Sparser representation formats based on the same principle as used in The. 60 can be constructed to represent markov networks by markov logic networks. Here, we can separately instantiate the factors by combinations of terms and clauses only involving the variables contained in the factor.

#### 11.1.4 Examples

Let us now provide examples of markov logic networks.

**Distribution of independent variables** We show next, the independent positive distributions are representable by tuning the  $d$  weights of the atomic formulas and keeping all other weights zero.

**Theorem 61.** *Let  $\mathbb{P} [X_{[d]}]$  be a positive probability distribution, such that atomic formulas are independent from each other. Then  $\mathbb{P} [X_{[d]}]$  is the Markov Logic Network of atomic formulas*

$$\mathcal{F}_{[d]} = \{X_k : k \in [d]\}$$

and parameters

$$\theta [L = k] = \ln \left[ \frac{\langle \mathbb{P} \rangle [X_k = 1]}{\langle \mathbb{P} \rangle [X_k = 0]} \right]$$

*Proof.* By Theorem 6 we get a decomposition

$$\mathbb{P} [X_{[d]}] = \bigotimes_{k \in [d]} \mathbb{P}^k [X_k]$$

where

$$\mathbb{P}^k [X_k] = \langle \mathbb{P} \rangle [X_k] .$$

By assumption of positivity, the vector  $\mathbb{P}^k [X_k]$  is positive for each  $k \in [d]$  and the parameter

$$\theta [L = k] = \ln \left[ \frac{\mathbb{P}^k [X_k = 1]}{\mathbb{P}^k [X_k = 0]} \right]$$

well-defined.

We then notice, that

$$\mathbb{P}^{(\{X_k\}, \theta[L=k])} [X_k] = \mathbb{P}^k [X_k]$$

and therefore with the parameter vector of dimension  $p = d$  defined as

$$\theta [L] = \sum_{k \in [d]} \theta [L = k] \cdot \epsilon_k [L]$$

we have

$$\begin{aligned} \mathbb{P}(\{X_k : k \in [d]\}, \theta) [X_{[d]}] &= \bigotimes_{k \in [d]} \mathbb{P}(\{X_k\}, \theta[L=k]) [X_k] \\ &= \bigotimes_{k \in [d]} \mathbb{P}^k [X_k] \\ &= \mathbb{P} [X_{[d]}] . \end{aligned}$$

□

In Theorem 61 we made the assumption of positive distributions. If the distribution fails to be positive, we still get a decomposition into distributions of each variable, but there is at least one factor failing to be positive. Such factors need to be treated by hybrid logic networks, that is they are base measure for an exponential family coinciding with a logical literal (see Sect. 11.2).

All atomic formulas can be selected by a single variable selecting tensor, that is

$$\phi(\{X_k : k \in [d]\}, \theta) [X_{[d]}] = \langle \mathcal{H}_V (X_{[d]}, L), \theta [L] \rangle [X_{[d]}] .$$

In case of negative coordiantes  $\theta [L = k]$  it is convenient to replace  $X_k$  by  $\neg X_k$ , in order to facilitate the interpretation. The probability distribution is left invariant, when also replacing  $\theta [L = k]$  by  $-\theta [L = k]$ .

**Boltzmann machines** A Boltzmann machine is a distribution of boolean variables  $X_{[d]}$  depending on weight tensors

$$W[L_{V,0}, L_{V,1}] \in \mathbb{R}^d \otimes \mathbb{R}^d \quad \text{and} \quad b[L_{V,0}] \in \mathbb{R}^d .$$

Its distribution is

$$\mathbb{P}^{W,b} [X_{[d]}] = \langle \exp [\phi^{W,b} [X_{[d]}]] \rangle [X_{[d]} | \emptyset]$$

where its energy tensor is

$$\phi^{W,b} [X_{[d]} = x_{[d]}] = \sum_{k,l \in [d]} W[L_{V,0} = k, L_{V,1} = l] \cdot x_k \cdot x_l + \sum_{k \in [d]} b[L_{V,0} = k] \cdot x_k .$$

We notice, that this tensor coincides with the energy tensor of a markov logic network with formula set

$$\mathcal{F} = \{X_k \Leftrightarrow X_l : k, l \in [d]\} \cup \{X_k : k \in [d]\}$$

with cardinality  $d^2 + d$ . Each formula is in the expressivity of an architecture consisting of a single binary logical neuron selecting any variable of  $X_{[d]}$  in each argument and selecting connectives  $\{\Leftrightarrow, \triangleleft\}$ , where by  $\triangleleft$  we refer to a connective passing the first argument, defined for  $x_0 \in [m_0], x_1 \in [m_1]$  as

$$\triangleleft [X_0 = x_0, X_1 = x_1] = \mathcal{H}_V (X_0 = x_0, X_1 = x_1, L_V = 0) .$$

When we choose the canonical parameter as

$$\theta [L_C, L_{V,0}, L_{V,1}] = \epsilon_0 [L_C] \otimes W[L_{V,0}, L_{V,1}] + \epsilon_1 [L_C] \otimes b[L_{V,0}] \otimes \epsilon_0 [L_{V,1}] .$$

we have (see Figure 31)

$$\phi^{W,b} [X_{[d]}] = \langle \sigma^A [X_{[d]}, L_C, L_{V,0}, L_{V,1}], \theta [L_C, L_{V,0}, L_{V,1}] \rangle [X_{[d]}] .$$

Therefore, Boltzmann machines are specific markov logic networks with the statistic being biimplications between atoms and atoms itself. Generic markov logic networks are more expressive than Boltzmann machines, by the flexibility to create further features by propositional formulas.

## 11.2 Hard Logic Networks

While exponential families are positive distributions, in logics probability distributions can assign states zero probability. As a consequence, Markov Logic Networks have a soft logic interpretation in the sense that violation of activated formulas have nonzero probability. We here discuss their hard logic counterparts, where worlds not satisfying activated formulas have zero probability.

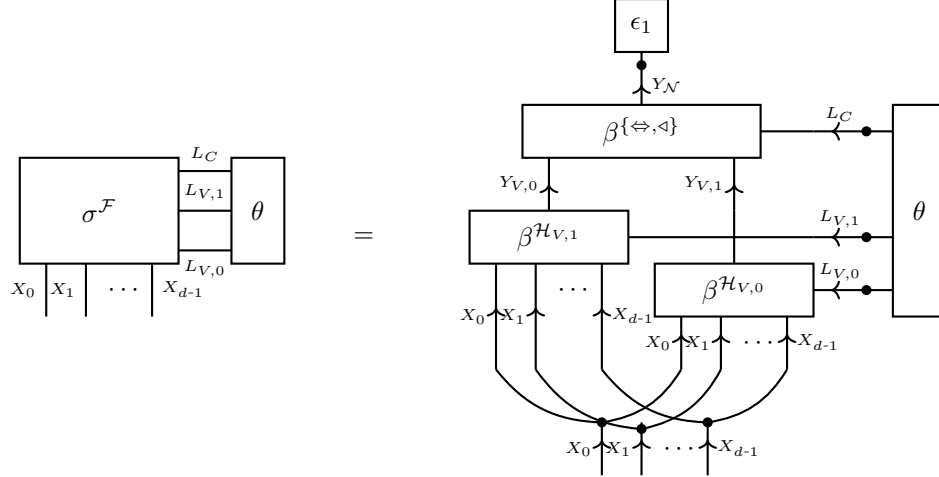


Figure 31: Tensor network representation of the energy of a Boltzmann machine

### 11.2.1 The limit of hard logic

The probability function of Markov Logic Networks with positive weights mimiks the tensor network representation of the knowledge base, which is the conjunction of the formulas. The maxima of the probability function coincide with the models of the corresponding knowledge base, if the latter is satisfiable. However, since the Markov Logic Network is defined as a normed exponentiation of the weighted formula sum, it is a positive distribution whereas uniform distributions among the models of a knowledge base assign zero probability to world failing to be a model. Since both distributions are tensors in the same space to a factored system, we can take the limits of large weights and observe, that Markov Logic Networks indeed converge to normed knowledge bases.

**Lemma 16.** *For any satisfiable formula  $f[X_{[d]}$ ] and a variable weight  $\theta \in \mathbb{R}$ , we have for  $\theta \rightarrow \infty$*

$$\langle \exp[\theta \cdot f[X_{[d]}]] \rangle [X_{[d]}|\emptyset] \rightarrow \langle f \rangle [X_{[d]}|\emptyset]$$

and for  $\theta \rightarrow -\infty$

$$\langle \exp[\theta \cdot f[X_{[d]}]] \rangle [X_{[d]}|\emptyset] \rightarrow \langle \neg f \rangle [X_{[d]}|\emptyset] .$$

Here we denote the understand the convergence of tensors as a convergence of each coordinate.

*Proof.* We have

$$\mathcal{Z}(\mathcal{F}, \theta) = \left( \prod_{k \in [d]} m_k \right) - \langle f \rangle [\emptyset] + \langle f \rangle [\emptyset] \cdot \exp[\theta]$$

and therefore for any  $x_{[d]} \in \times_{k \in [d]} [2]$  with  $f[X_{[d]} = x_{[d]}] = 1$

$$\begin{aligned} \langle \exp[\theta \cdot f] \rangle [X_{[d]} = x_{[d]}|\emptyset] &= \frac{\exp[\theta]}{\left( \prod_{k \in [d]} m_k \right) - \langle f \rangle [\emptyset] + \langle f \rangle [\emptyset] \cdot \exp[\theta]} \\ &\rightarrow \frac{1}{\langle f \rangle [\emptyset]} = \langle f \rangle [X_0 = x_0, \dots, X_{d-1} = x_{d-1}|\emptyset] . \end{aligned}$$

For any  $x_0, \dots, x_{d-1} \in \times_{k \in [d]} [2]$  with  $f[X_{[d]} = x_{[d]}] = 0$  we have on the other side

$$\begin{aligned} \langle \exp[\theta \cdot f] \rangle [X_0 = x_0, \dots, X_{d-1} = x_{d-1}|\emptyset] &= \frac{1}{\left( \prod_{k \in [d]} m_k \right) - \langle f \rangle [\emptyset] + \langle f \rangle [\emptyset] \cdot \exp[\theta]} \\ &\rightarrow 0 = \langle f \rangle [X_0 = x_0, \dots, X_{d-1} = x_{d-1}|\emptyset] . \end{aligned}$$

□

We can by the above Lemma represent both the situation of non-asymptotic weights and the limit for diverging weights by the same computation core  $\beta^f [Y_f, X_{[d]}]$ , with different activation cores, since

$$\langle \exp [\theta \cdot f [X_{[d]}]] \rangle [X_{[d]} | \emptyset] = \langle \beta^f [Y_f, X_{[d]}], \alpha^{f, \theta} \rangle [X_{[d]}]$$

and

$$\langle f \rangle [X_{[d]} | \emptyset] = \langle \beta^f [Y_f, X_{[d]}], \epsilon_1 [Y_f] \rangle [X_{[d]}]$$

respectively

$$\langle \neg f \rangle [X_{[d]} | \emptyset] = \langle \beta^f [Y_f, X_{[d]}], \epsilon_0 [Y_f] \rangle [X_{[d]}] .$$

**Theorem 62.** *Let  $\mathcal{F}$  be a formula set and  $\theta$  a positive parameter vector. If the formula*

$$\mathcal{KB} = \bigwedge_{f \in \mathcal{F}} f$$

*is satisfiable we have in the limit  $\beta \rightarrow \infty$  the coordinatewise convergence*

$$\mathbb{P}^{(\mathcal{F}, \beta \cdot \theta)} [X_{[d]}] \rightarrow \langle \mathcal{KB} \rangle [X_{[d]}] .$$

*Proof.* Since  $\mathcal{KB}$  is satisfiable we find  $x_0, \dots, x_{d-1} \in \times_{k \in [d]} [2]$  with

$$\left\langle \exp \left[ \sum_{f \in \mathcal{F}} \beta \cdot \theta_f \cdot f \right] \right\rangle [X_0 = x_0, \dots, X_{d-1} = x_{d-1}] = \exp \left[ \beta \cdot \sum_{f \in \mathcal{F}} \theta_f \right]$$

and the partition function obeys

$$\left\langle \exp \left[ \sum_{f \in \mathcal{F}} \beta \cdot \theta_f \cdot f \right] \right\rangle [\emptyset] \geq \exp \left[ \beta \cdot \sum_{f \in \mathcal{F}} \theta_f \right] .$$

For any state  $x_0, \dots, x_{d-1} \in \times_{k \in [d]} [2]$  with  $\mathcal{KB}(x_0, \dots, x_{d-1}) = 0$  we find  $h \in \mathcal{F}$  with

$$h(x_0, \dots, x_{d-1}) = 0$$

and have

$$\frac{\left\langle \exp \left[ \sum_{f \in \mathcal{F}} \beta \cdot \theta_f \cdot f \right] \right\rangle [X_0 = x_0, \dots, X_{d-1} = x_{d-1}]}{\left\langle \exp \left[ \sum_{f \in \mathcal{F}} \beta \cdot \theta_f \cdot f \right] \right\rangle [\emptyset]} \leq \frac{\exp \left[ \beta \cdot \sum_{f \in \mathcal{F}: f \neq h} \theta_f \right]}{\exp \left[ \beta \cdot \sum_{f \in \mathcal{F}} \theta_f \right]} = \exp [\beta \cdot \theta_h] \rightarrow 0 .$$

The limit of the distribution has thus support only on the models of  $\mathcal{KB}$ . Since any model of  $\mathcal{KB}$  has same energy at any  $\beta$  the limit is a uniform distribution and coincides therefor with

$$\langle \mathcal{KB} \rangle [X_{[d]} | \emptyset] .$$

□

We here assumed, that the conjunction  $\mathcal{KB}$  of the formulas in  $\mathcal{F}$  is satisfiable, and showed, that the markov logic network converges in the limit of large weights to the uniform distribution of the models of  $\mathcal{KB}$ . In Chapter 12 we will drop this assumption and show that the limit is the face base measure associated with the corresponding face of the mean parameter polytope. The face base measure coincides thereby with  $\mathcal{KB}$ , if  $\mathcal{KB}$  is satisfiable.

### 11.2.2 Tensor Network Representation

Hard Logic Network coincide with Knowledge Bases and are thus representable by contractions of formulas (which can be interpreted as a hybrid calculus scheme, see Sect. 17.6).

**Theorem 63** (Conjunction Decomposition of Knowledge Bases). *For a Knowledge Base*

$$\mathcal{KB} = \bigwedge_{f \in \mathcal{F}} f$$

*we have*

$$\mathcal{KB} [X_{[d]}] = \langle f [X_{[d]}] \rangle [X_{[d]}]$$

*and*

$$\mathcal{KB} [X_{[d]}] = \langle \{ \beta^f [Y_f, X_{[d]}] : f \in \mathcal{F} \} \cup \{ \epsilon_1 [Y_f] : f \in \mathcal{F} \} \rangle [X_{[d]}] .$$

*Proof.* This follows from the representation of conjunctions by contraction (see Sect. 17.6) and

$$f [X_{[d]}] = \langle \beta^f [Y_f, X_{[d]}], \epsilon_1 [Y_f] \rangle [X_{[d]}] .$$

□

We call this representation scheme the  $\wedge$ -symmetry, since we can either represent  $\mathcal{KB}$  by instantiation of  $\beta^{\mathcal{KB}}$ , which involves a basis encoding of the conjunction  $\wedge$ , or by instantiations of a collection of  $\beta^f$ . We use the  $\wedge$  symmetry to represent them as a contraction of the formulas building the Knowledge Base as conjunction.

**Remark 12.**  $\wedge$  symmetry does not generalize to Markov Logic Networks. In Markov Logic Networks, similar decompositions are not possible. For example, consider a MLN with a single formula  $X_0 \wedge X_1$  and nonvanishing weight  $\theta$ . This does not coincide with the distribution of a MLN of two formulas  $X_0$  and  $X_1$ . To see this, we notice that with respect to the distribution of the first MLN, both variables are not independent, while for any MLN constructed by the two atomic formulas they are.

### 11.3 Hybrid Logic Network

Markov Logic Networks are by definition positive distributions. In contrary, Hard Logic Networks model uniform distributions over model sets of the respective Knowledge Base and therefore have vanishing coordinates. We now show how to combine both approaches by defining Hybrid Logic Networks, when understanding Hard Logic Networks as base measures. This trick is known to the field of variational inference, see for Example 3.6 in Wainwright and Jordan (2008).

**Definition 60.** Given a set of formulas  $\mathcal{F}$  with weights  $\theta$  and set  $\mathcal{KB}$  of formulas, which conjunction is satisfiable, the hybrid logic network is the probability distribution

$$\mathbb{P}^{(\mathcal{F}, \theta, \nu^{\mathcal{KB}})} [X_{[d]}] = \langle \{f : f \in \mathcal{KB}\} \cup \{\exp [\theta_f \cdot f] : f \in \mathcal{F}\} [X_{[d]}] \rangle ,$$

which is the member of the exponential family with statistic by  $\mathcal{F}$  and the base measure

$$\nu^{\mathcal{KB}} [X_{[d]}] = \langle \{f : f \in \mathcal{KB}\} [X_{[d]}] \rangle .$$

Given a set of formulas  $\mathcal{F}$ , we define the set of hybrid logic networks realizable with  $\mathcal{F}$  and elementary activation cores as

$$\Lambda^{\mathcal{F}, \text{EL}} = \bigcup_{\tilde{\mathcal{F}} \subset \mathcal{F}, \mu \in \{0,1\}^{|\mathcal{F}|}} \Gamma^{\mathcal{F}/\tilde{\mathcal{F}}, \nu^{\tilde{\mathcal{F}}}, \mu}$$

where we denote base measures

$$\nu^{\tilde{\mathcal{F}}, \mu} [X_{[d]}] = \bigwedge_{f_l \in \tilde{\mathcal{F}}} \neg^{(1-\mu[L=l])} f_l [X_{[d]}] .$$

The assumption of a satisfiable set  $\mathcal{KB}$  is necessary, as we show next.

**Theorem 64.** If any only if  $\bigwedge_{f \in \mathcal{KB}} f$  is satisfiable, the tensor

$$\langle \{f : f \in \mathcal{KB}\} \cup \{\exp [\theta_f \cdot f] : f \in \mathcal{F}\} [X_{[d]}] \rangle$$

is normable.

*Proof.* We need to show that

$$\langle \{f : f \in \mathcal{KB}\} \cup \{\exp [\theta_f \cdot f] : f \in \mathcal{F}\} [\emptyset] \rangle > 0 . \quad (24)$$

Since the conjunction of  $\mathcal{KB}$  is satisfiable we find a  $x_{[d]}$  with  $f [X_{[d]} = x_{[d]}] = 1$  for all  $f \in \mathcal{KB}$ . Then

$$\begin{aligned} \langle \{f : f \in \mathcal{KB}\} \cup \{\exp [\theta_f \cdot f] : f \in \mathcal{F}\} [X_{[d]} = x_{[d]}] \rangle &= \left( \prod_{f \in \mathcal{KB}} f [X_{[d]} = x_{[d]}] \right) \\ &\cdot \left( \prod_{f \in \mathcal{F}} \exp [\theta_f \cdot f] [X_{[d]} = x_{[d]}] \right) \\ &= \left( \prod_{f \in \mathcal{F}} \exp [\theta_f \cdot f] [X_{[d]} = x_{[d]}] \right) \\ &> 0 . \end{aligned}$$

Condition (24) follows from this and the Hybrid Logic Network is well-defined. □



### 11.3.1 Tensor Network Representation

We can employ the formula decompositions to represent both probabilistic facts of the MLN and hard facts (seen as the limit of large weights).

**Theorem 65.** *For any hybrid logic network we have*

$$\mathbb{P}^{(\mathcal{F}, \theta, \mathcal{KB})}[X_{[d]}] = \langle \{\beta^f[Y_f, X_{[d]}] : f \in \mathcal{F} \cup \mathcal{KB}\} \cup \{\epsilon_1[Y_f] : f \in \mathcal{KB}\} \cup \{\alpha^f[Y_f] : f \in \mathcal{F}\} \rangle [X_{[d]} | \emptyset] .$$

*Proof.* By Lem. 11. □

While the statistics computing cores in the basis encoding are shared to compute the soft and the hard logic formulas, their activation cores differ. While probabilistic soft formulas get activation cores (see Theorem 59)

$$\alpha^f[Y_f] = \left[ \exp \left[ \theta \left[ \frac{1}{L} \right] \right] \right] [Y_l]$$

the hard formulas get activation cores by unit vectors

$$\epsilon_1[Y_f] = \begin{bmatrix} 0 \\ 1 \end{bmatrix} [Y_l] .$$

As shown in Sect. 11.2.1, the soft activation cores converge to these hard activation cores in the limit of large parameters, when imposing a local normalization. We further notice, that the probabilistic activation cores are trivial tensors if and only if the corresponding parameter coordinate vanishes.

For an example see Figure 32.

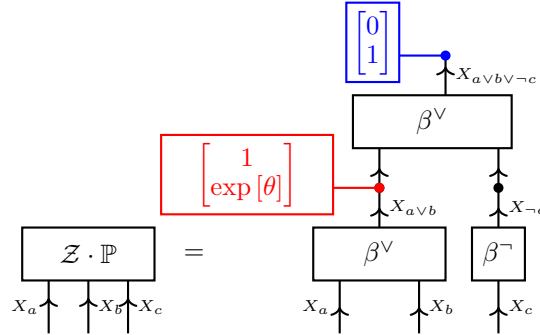


Figure 32: Diagram of a formula tensor with activated heads, containing **hard constraint cores** and **probabilistic weight cores**.

**Remark 13.** *Probability interpretation using the Partition function* The tensor networks here represent unnormalized probability distributions. The probability distribution can be normed by the quotient with the naive contraction of the network, the partition function.

### 11.3.2 Reasoning Properties

Deciding probabilistic entailment (see Def. 47) with respect to Hybrid Logic Networks can be reduced to the hard logic parts of the network.

**Theorem 66.** *Let  $(\mathcal{F}, \theta, \mathcal{KB})$  define a Hybrid Logic Network. Given a query formula  $f$  we have that*

$$\mathbb{P}^{(\mathcal{F}, \theta, \mathcal{KB})} \models f$$

*if and only if*

$$\mathcal{KB} \models f .$$

*Proof.* This follows from Theorem 51 on the representation of Hybrid Logic Networks as Markov Networks in Theorem 65. □

Formulas in  $\mathcal{F}$ , which are entailed or contradicted by  $\mathcal{KB}$  are redundant, as we show next.

**Theorem 67.** *If for a formula  $f$  and  $\mathcal{KB}$  we have*

$$\mathcal{KB} \models f \quad \text{or} \quad \mathcal{KB} \models \neg f$$

*then for any  $(\mathcal{F}, \theta, \mathcal{KB})$*

$$\mathbb{P}^{(\mathcal{F}/\{f\}, \tilde{\theta}, \mathcal{KB})} [X_{[d]}] = \mathbb{P}^{(\mathcal{F}, \theta, \mathcal{KB})} [X_{[d]}] ,$$

*where  $\tilde{\theta}$  denotes the tensor  $\theta$ , where the coordinate to  $f$  is dropped, if  $f \in \mathcal{F}$ .*

*Proof.* Isolate the factor to the hard formula, which is constant for all situations. □

A similar statement holds for the hard formulas itself, as shown in Theorem 46. However, notice that if  $\mathcal{KB}/\{f\} \models \neg f$ , then  $\mathcal{KB} \cup \{f\}$  is not satisfiable and a hybrid logic network cannot be defined for  $\mathcal{KB} \cup \{f\}$  as hard logic formulas.

These results are especially interesting for the efficient implementation of Algorithm 7, which has been introduced in Chapter 8. By Theorem 66 only the hard logic parts of a Hybrid Logic Network are required in the ASK operation.

### 11.3.3 Expressivity

Hybrid Logic Networks extend the expressivity result of Theorem 60 to arbitrary probability tensors, dropping the positivity constraints for Markov Logic Networks.

**Theorem 68.** *Let  $\mathbb{P} [X_{[d]}]$  a possibly not positive probability tensor we build a base measure*

$$\nu^{\mathcal{KB}} = \mathbb{I}_{\neq 0} (\mathbb{P} [X_{[d]}])$$

*and a parameter tensor*

$$\theta [L_{[d]} = x_{[d]}] = \begin{cases} 0 & \text{if } \mathbb{P} [X_{[d]} = x_{[d]}] = 0 \\ \ln [\mathbb{P} [X_{[d]} = x_{[d]}]] & \text{else} \end{cases} .$$

*Then the probability tensor is the member of the minterm exponential family with base measure  $\mathcal{KB}$  and parameter  $\theta$ , that is*

$$\mathbb{P} [(\mathcal{F}_{\wedge}, \theta, \nu^{\mathcal{KB}})]$$

*Proof.* It suffices to show that

$$\langle \nu^{\mathcal{KB}}, \exp [\langle \sigma^{\mathcal{F}_{\wedge}} \theta \rangle [X_{[d]}]] \rangle [X_{[d]}] = \mathbb{P} [X_{[d]}] .$$

For indices  $x_{[d]}$  with  $\mathbb{P} [X_{[d]} = x_{[d]}] = 0$  we have  $\nu^{\mathcal{KB}} [X_{[d]} = x_{[d]}] = 0$  and thus also

$$\langle \nu^{\mathcal{KB}}, \exp [\langle \sigma^{\mathcal{F}_{\wedge}} \theta \rangle [X_{[d]}]] \rangle [X_{[d]} = x_{[d]}] = 0 .$$

For indices  $x_{[d]}$  with  $\mathbb{P} [X_{[d]} = x_{[d]}] > 0$  we have  $\nu^{\mathcal{KB}} [X_{[d]} = x_{[d]}] = 1$  and

$$\begin{aligned} \langle \nu^{\mathcal{KB}}, \exp [\langle \sigma^{\mathcal{F}_{\wedge}} \theta \rangle [X_{[d]}]] \rangle [X_{[d]} = x_{[d]}] &= \prod_{l_{[d]}} \exp [\theta [L_{[d]} = l_{[d]}] \cdot Z_{l_{[d}}}^{\wedge} [X_{[d]} = x_{[d]}]] \\ &= \exp [\theta [L_{[d]} = x_{[d]}]] \\ &= \mathbb{P} [X_{[d]} = x_{[d]}] . \end{aligned}$$

□

### 11.4 Polynomial Representation

We now sparse representation formats for the introduced logic networks, namely the basis+ CP format introduced in Chapter 18 which are understood as polynomial decompositions. First of all, we establish a sparsity result for terms and clauses (see Def. 45).

**Lemma 17.** *Any term is representable by a single monomial and any clause is representable by a sum of at most two monomials.*

*Proof.* Let  $\mathcal{V}_0$  and  $\mathcal{V}_1$  be disjoint subsets of  $\mathcal{V}$ , then we have

$$Z_{\mathcal{V}_0, \mathcal{V}_1}^\wedge = \epsilon_{\{x_k=0:k \in \mathcal{V}_0\} \cup \{x_k=1:k \in \mathcal{V}_1\}} [X_{\mathcal{V}_0 \cup \mathcal{V}_1}] \otimes \mathbb{I} [X_{\mathcal{V}/(\mathcal{V}_0 \cup \mathcal{V}_1)}]$$

and

$$Z_{\mathcal{V}_0, \mathcal{V}_1}^\vee = \mathbb{I} [X_{\mathcal{V}}] - \epsilon_{\{x_k=0:k \in \mathcal{V}_0\} \cup \{x_k=1:k \in \mathcal{V}_1\}} [X_{\mathcal{V}_0 \cup \mathcal{V}_1}] \otimes \mathbb{I} [X_{\mathcal{V}/(\mathcal{V}_0 \cup \mathcal{V}_1)}] .$$

We notice, that any tensors  $\mathbb{I}$  and  $\epsilon_x \otimes \mathbb{I}$  have basis+-rank of 1 and therefore  $Z_{\mathcal{V}_0, \mathcal{V}_1}^\wedge$  of 1 and  $Z_{\mathcal{V}_0, \mathcal{V}_1}^\vee$  of at most 2.  $\square$

A formula in conjunctive normal form is a conjunction of clauses, where clauses are disjunctions of literals being atoms (positive literals) or negated atoms (negative literals). Based on these normal forms, we show representations of formulas as sparse polynomials. We apply Lem. 17 to show the following sparsity bound.

**Theorem 69.** *Any formula  $f$  with a conjunctive normal form of  $n$  clauses satisfies*

$$\text{rank}^d(f) \leq 2^n .$$

*Proof.* Let  $f$  have a conjunctive normal form with clauses indexed by  $l \in [n]$  and each clause represented by subsets  $\mathcal{V}_0^l, \mathcal{V}_1^l$ , that is

$$f = \bigwedge_{l \in [n]} Z_{\mathcal{V}_0^l, \mathcal{V}_1^l}^\vee .$$

We now use the rank bound of The. 124 and Lem. 17 to get

$$\text{rank}^d(f) \leq \prod_{l \in [n]} \text{rank}^d(Z_{\mathcal{V}_0^l, \mathcal{V}_1^l}^\vee) \leq 2^n .$$

$\square$

We apply this result on the sparse representation of a single formula to derive sparse representations for Hard Logic Networks and the energy tensor of Hybrid Logic Networks. Both results use in addition to The. 69 sparsity bounds, which are shown by explicit representation construction in Chapter 18.

**Corollary 6.** *Any Hard Logic Network  $\mathcal{KB}$  obeys*

$$\text{rank}^d(\mathcal{KB}) \leq \prod_{f \in \mathcal{KB}} 2^{n_f}$$

*Proof.* We apply the contraction bound The. 124 for the decomposition

$$\mathcal{KB} [X_{[d]}] = \langle \{f [X_{[d]}] : f \in \mathcal{KB}\} \rangle [X_{[d]}]$$

and get

$$\text{rank}^d(\mathcal{KB}) \leq \prod_{f \in \mathcal{KB}} \text{rank}^d(f) .$$

The claimed bound follows with The. 69.  $\square$

**Corollary 7.** *The energy tensor of a Hybrid Logic Network with statistic  $\mathcal{F}$*

$$\text{rank}^d(\langle \sigma^\mathcal{F} [X_{[d]}, L], \theta [L] \rangle [X_{[d]}]) \leq \sum_{l \in [p] : \theta [L=l] \neq 0} 2^{n_{f_l}} .$$

where  $n_{f_l}$  denotes the number of clauses in a conjunctive normal form of  $f_l$ .

*Proof.* We decompse the energy into the sum

$$\langle \sigma^\mathcal{F} [X_{[d]}, L], \theta [L] \rangle [X_{[d]}] = \sum_{l \in [p] : \theta [L=l] \neq 0} \theta [L=l] \cdot f_l [X_{[d]}]$$

and apply The. 123 to get

$$\text{rank}^d(\langle \sigma^\mathcal{F} [X_{[d]}, L], \theta [L] \rangle [X_{[d]}]) \leq \sum_{l \in [p] : \theta [L=l] \neq 0} \text{rank}^d(f_l [X_{[d]}]) \leq \sum_{l \in [p] : \theta [L=l] \neq 0} 2^{n_{f_l}} .$$

$\square$

## 11.5 Applications

Hybrid Logic Networks as neuro-symbolic architectures:

- Neural Paradigm here by decompositions of logical formulas into their connectives. In more generality by decompositions of sufficient statistics into composed functions, using Basis Calculus. Deeper nodes as carrying correlations of lower nodes.
- Symbolic Paradigm by interpretability of propositional logics.

Hybrid Logic Networks as trainable Machine Learning models:

- Expressivity: Can represent any positive distribution, as shown by Theorem 60, with  $2^d$  formulas.
- Efficiency: Can only handle small subsets of possible formulas, since their possible number is huge. Tensor networks provide means to efficiently represent formulas depending on many variables and reason based on contractions.
- Differentiability: Distributions are differentiable functions of their weights, see Parameter Estimation Chapter. The log-likelihood of data is therefore also differentiable function of their weights and we can exploit first-order methods in their optimization.
- Structure Learning: We need to find differentiable parametrizations of logical formulas respecting a chosen architecture. In Chapter 10 such representations are described based on Selector Tensor Networks.

Differentiability and structure learning will be investigated in more detail in the next chapter.

When understanding atoms as observed variables, and the computed as hidden, Hybrid Logic Networks are deep higher-order boltzmann machines: More generic correlations can be captured by a logical connective, calculated by a basis encoding and activated by an activation core.

Hybrid Logic Networks as bridging soft and hard logics within the formalism of exponential families.

A more general class of problems, which have natural representations by Hard Logic Networks are Constraint Satisfaction Problems (see Chapter 5 in Russell and Norvig (2021)). Solving such problems is then equivalent to sampling from the worlds in a logical interpretation, and can be approached by the methods of Chapter 8. Among these classed, we have only discussed the Sudoku game in Example 7. Extensions by Hybrid Logic Networks can be interpreted as implementations of preferences among possible solutions by probabilities.

## 12 Logical Network Inference

In this chapter we investigate the inference properties of Hybrid Logic Networks starting with characterizations of its mean parameter polytopes. We investigate unconstrained parameter estimated for Markov Logic Networks and Hybrid Logic Networks, which are special cases of the backward maps introduced in Chapter 5. We then motivate structure learning based on sparsity constraints on the parameters on the minterm exponential family and present heuristic strategies leading to efficient structure learning algorithms.

### 12.1 Mean parameters of Hybrid Logic Networks

While mean parameter polytopes  $\mathcal{M}_{\mathcal{S},\nu}$  to generic exponential families have been subject to Chapter 6, we in this section restrict to the mean polytopes of hybrid logic networks, which we characterize using propositional logics. Hybrid Logic Networks are exponential families, which statistic  $\mathcal{S}$  consists of coordinates with  $\text{im}(\mathcal{S}_l) \subset \{0, 1\}$  and are therefore propositional formulas. The convex polytope of mean parameters (see Def. 32) is for a statistic  $\mathcal{F}$  of propositional formulas and a base measure  $\nu$

$$\mathcal{M}_{\mathcal{F},\nu} = \{ \langle \mathbb{P}, \sigma^{\mathcal{F}} \rangle [L] : \mathbb{P} \in \Gamma^{\delta,\nu} \} ,$$

where by  $\Gamma^{\delta,\nu}$  we denote the set of all by  $\nu$  representable probability distributions. By The. 16 the convex polytope has a characterization as a convex hull

$$\mathcal{M}_{\mathcal{F},\nu} = \text{conv} \left( \sigma^{\mathcal{F}} [X_{[d]} = x_{[d]}, L] : x_{[d]} \in \bigtimes_{k \in [d]} [2], \nu [X_{[d]} = x_{[d]}] = 1 \right) . \quad (25)$$

We notice, that all  $\sigma^{\mathcal{F}} [X_{[d]} = x_{[d]}, L]$  are boolean vectors in  $\mathbb{R}^p$ . The mean parameter polytopes are thus of 0/1-polytopes Ziegler (2000); Gillmann (2007), from which a few obvious properties follow. Since those are convex subsets of the cube  $[0, 1]^p$ , which vertices are all binary vectors, also each  $\sigma^{\mathcal{F}} [X_{[d]} = x_{[d]}, L]$  (with  $\nu [X_{[d]} = x_{[d]}] = 1$ ) is an extreme point. Further, if for any  $l \in [p]$  we have  $\mu [L = l] \in \{0, 1\}$ , then  $\mu [L]$  is in the boundary of the cube and thus also of  $\mathcal{M}_{\mathcal{F}, \nu}$ .

In the following, we characterize the faces of the mean parameter

### 12.1.1 Vertices by hard logic networks

$\nu^{\mathcal{F}, \mu} [X_{[d]}]$  are the subset encodings of preimages of the statistics encoding.

We exploit this characterization to show, that the vertices of  $\mathcal{M}_{\mathcal{F}, \nu}$  are exactly those reproducible by Hard Logic Networks.

**Theorem 70.** Any set  $\{\mu [L]\} \subset \mathcal{M}_{\mathcal{F}, \nu}$  of cardinality 1 is a vertex of  $\mathcal{M}_{\mathcal{F}, \nu}$ , if and only if its unique element  $\mu [L]$  is boolean and the formula

$$\nu^{\mathcal{F}, \mu} [X_{[d]}] := \bigwedge_{l \in [p]} \neg^{(1-\mu[L=l])} f_l [X_{[d]}]$$

is satisfiable. In that case,  $\mu$  is the mean parameter of the Hard Logic Network with formulas

$$\mathcal{KB} = \{\neg^{(1-\mu[L=l])} f_l : l \in [p]\},$$

where we denote by  $\neg^0$  the identity connective and by  $\neg^1 = \neg$  the logical negation.

*Proof.* " $\Rightarrow$ ": Let  $\mu$  be an extreme point of  $\mathcal{M}_{\mathcal{F}, \nu}$ . Since by (25)  $\mathcal{M}_{\mathcal{F}, \nu}$  is a convex hull of vectors, there exists a  $x_{[d]} \in \times_{k \in [d]} [2]$  such that

$$\mu [L] = \sigma^{\mathcal{F}} [X_{[d]} = x_{[d]}, L].$$

By definition of  $\sigma^{\mathcal{F}}$ ,  $\mu [L]$  is a boolean vector and for any  $l \in [p]$  we have

$$f_l [X_{[d]} = x_{[d]}] = \mu [L = l]$$

and thus

$$\neg^{(1-\mu[L=l])} f_l [X_{[d]} = x_{[d]}] = 1.$$

It follows that  $x_{[d]}$  is also a model of  $\nu^{\mathcal{F}, \mu}$  and  $\nu^{\mathcal{F}, \mu}$  is satisfiable.

" $\Leftarrow$ ": To show the converse direction, let  $\mu [L]$  be a boolean vector such that  $\nu^{\mathcal{F}, \mu}$  is satisfiable. Then there exists a model  $x_{[d]}$  of  $\nu^{\mathcal{F}, \mu}$ . We have for any  $l \in [p]$

$$\sigma^{\mathcal{F}} [X_{[d]} = x_{[d]}, L = l] = f_l [X_{[d]} = x_{[d]}] = \mu [L = l]$$

and thus  $\mu [L] = \sigma^{\mathcal{F}} [X_{[d]} = x_{[d]}, L]$ . With the characterization (25) this establishes in particular  $\mu [L] \in \mathcal{M}_{\mathcal{F}, \nu}$ . Since  $\mu [L]$  is boolean and therefore an extreme point of the cube  $[0, 1]^p$ , it is also an extreme point of the subset  $\mathcal{M}_{\mathcal{F}, \nu} \subset [0, 1]^p$ .  $\square$

### 12.1.2 Faces of larger rank

Since the by inclusion partially ordered set of faces is a graded lattice (see Theorem 2.7 in Ziegler (2013)), the faces are ranked by dimension of th

The face base measure is thus

$$\nu^{\mathcal{I}} = \bigcup_{\mu \in Q_{\mathcal{S}, \nu}^{\mathcal{I}} \cup \{0, 1\}^p} f^{\mu} = \bigcup_{\mu \in Q_{\mathcal{S}, \nu}^{\mathcal{I}} \cup \{0, 1\}^p} \bigwedge_{l \in [p]} \neg^{(1-\mu[L=l])} f_l.$$

### 12.1.3 Mean parameters in the interior

By The. 24 the interior points are those realizable by a Hybrid Logic Network with statistics  $\mathcal{F}$  and base measure  $\nu$ , as we state in the following Corollary.

**Corollary 8.** If  $\mu [L] \in \mathcal{M}_{\mathcal{F}, \nu}^{\circ}$ , or equivalently the statistic is minimal and  $\mu [L]$  is reproducible by a distribution positive with respect to  $\nu$ , then there is  $\theta [L]$  such that  $\mathbb{P}^{\mathcal{F}, \theta, \nu}$  reproduces  $\mu [L]$ .

### 12.1.4 Mean parameters outside the interior

By The. 18 mean parameter vectors outside the interior of  $\mathcal{M}_{\mathcal{F},\nu}$  are not realizable by distributions, which are positive with respect to the base measure  $\nu$ . Instead, we in this section construct refined base measures  $\tilde{\nu}$  (refinement in the sense of  $\tilde{\nu} \prec \nu$ ), such that there are distributions, which are positive with respect to  $\tilde{\nu}$  and represent such mean parameters.

First of all, we can use the criterion of mean parameter coordinates in  $\{0, 1\}$  as a sufficient condition for  $\mu[L] \notin \mathcal{M}_{\mathcal{F},\nu}^\circ$ . In this case, the next theorem provides us with a procedure to refine the base measure in these cases.

**Theorem 71** (Base measure refinement for mean coordinates in  $\{0, 1\}$ ). *To any  $\mu[L] \in \mathcal{M}_{\mathcal{F},\nu}$  we build the base measure*

$$\nu^{\mathcal{F},\mu}[X_{[d]}] := \bigwedge_{l \in [p] : \mu[L=l] \in \{0,1\}} \neg^{(1-\mu[L=l])} f_l[X_{[d]}] .$$

*Then any probability distribution reproducing  $\mu[L]$  has a representation with respect to the base measure*

$$\tilde{\nu}[X_{[d]}] = \langle \nu, \nu^{\mathcal{F},\mu} \rangle [X_{[d]}] .$$

*Proof.* Let  $\mathbb{P}$  be a distribution representable with respect to  $\nu$  and reproducing  $\mu$ , that is

$$\langle \mathbb{P}, \sigma^{\mathcal{F}} \rangle [L] = \mu[L] .$$

For any  $l \in [p]$ , if  $\mu[L=l] = 1$  we have

$$\langle \mathbb{P}[X_{[d]}], f_l[X_{[d]}] \rangle [\emptyset] = 1$$

and with The. 47 probabilistic entailment  $\mathbb{P} \models f_l$  (see Def. 47). On the other hand for any  $l \in [p]$ , if  $\mu[L=l] = 0$  we have

$$\langle \mathbb{P}[X_{[d]}], \neg f_l[X_{[d]}] \rangle [\emptyset] = 1$$

and with the same argumentation  $\mathbb{P} \models \neg f_l$ .

Together it holds that  $\mathbb{P} \models \nu^{\mathcal{F},\mu}$  and

$$\mathbb{P}[X_{[d]}] = \langle \mathbb{P}[X_{[d]}], \nu^{\mathcal{F},\mu}[X_{[d]}] \rangle [ ] .$$

Thus,  $\mathbb{P}$  is representable with respect to the base measure  $\nu^{\mathcal{F},\mu}$  and also with respect to the base measure  $\tilde{\nu}$ .  $\square$

By The. 71 we can reduce the statistics to the set  $\mathcal{F}^\mu = \{f_l \in \mathcal{F} : \mu[L=l] \notin \{0, 1\}\}$  and continue searching for a reproducing distribution, now for the reduced mean parameter with coordinates not in  $\{0, 1\}$ . The criterion of mean parameter coordinates in  $\{0, 1\}$  is, however, not a necessary condition for  $\mu[L] \notin \mathcal{M}_{\mathcal{F},\nu}^\circ$ , see Example 13 for minimal representations where mean parameters outside the interior exist with coordinates not in  $\{0, 1\}$ .

Following a different approach, any mean parameter outside the interior of the mean parameter polytope is on a face of the polytope. We can use this insight, to refine base measures by face base measures (see Def. 34), which has been shown in more generality in The. 32. For hybrid logic network we characterize the face base measures in the next theorem.

**Theorem 72** (Base measures by formula satisfaction). *Let  $\theta[L=l]$  be a normal to a face of  $\mathcal{M}$ . If and only if the face formula*

$$f_{\mathcal{F},\theta} = \bigwedge_{l \in [p] : \theta[L=l] \neq 0} \neg^{(1-\mathbb{I}_{>0}(\theta[L=l]))} f_l$$

*is satisfiable, then it is the base measure to the face with normal  $\theta$ . Here  $\mathbb{I}_{>0}(z)$  denotes the indicator of  $z > 0$ .*

*If the above is not satisfiable, then the base measure is*

$$\bigvee_{v[L] : \langle \theta, v \rangle [\emptyset] \in \max_{\mu} \langle \theta, \mu \rangle [\emptyset]} f_{\mathcal{F}, \langle \theta, v \rangle [L]}$$

*where  $v$  are boolean vectors.*

*Proof.* We show this lemma based on characterizations of

$$\operatorname{argmax}_{x_{[d]}} \langle \theta[L], \sigma^{\mathcal{F}}[X_{[d]} = x_{[d]}, L] \rangle [\emptyset] .$$

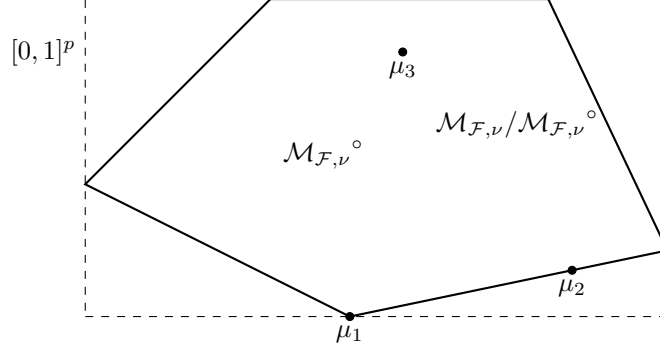


Figure 33: Sketch of the mean polytope  $\mathcal{M}_{\mathcal{F}, \nu}$  to a statistic  $\mathcal{F}$  which is minimal with respect to  $\nu$ , as a special case of the more generic sketch Figure 11. The mean polytope is a subset of the  $p$ -dimensional cube  $[0, 1]^p$  (here sketched as a 2-dimensional projection), where each mean parameter in one of the three cases  $\mu_1, \mu_2$  or  $\mu_3$ . Extreme points  $\mu_1 \in \mathcal{M}_{\mathcal{F}, \nu} \cap \{0, 1\}^p$  are those reproducible by a Hard Logic Network given  $\mathcal{F}$ . Non-extreme points outside the interior  $\mu_2 \in \mathcal{M}_{\mathcal{F}, \nu} / (\mathcal{M}_{\mathcal{F}, \nu}^\circ \cup \{0, 1\}^p)$  are reproducible by Hybrid Logic Networks with statistic  $\mathcal{F}$  and refined base measure  $\tilde{\nu}$ . Interior points  $\mu_3 \in \mathcal{M}_{\mathcal{F}, \nu}^\circ$  are reproducible by a Hybrid Logic Network with statistic  $\mathcal{F}$  and refined base measure  $\tilde{\nu}$ .

For the first claim: We have since  $\sigma^{\mathcal{F}}$  is boolean that

$$\|\theta[L]\|_1 \geq \max_{x_{[d]}} \langle \theta[L], \sigma^{\mathcal{F}}[X_{[d]} = x_{[d]}, L] \rangle [\emptyset] .$$

If and only if the formula  $f_{\mathcal{F}, \theta}$  is satisfiable, then we find a model  $x_{[d]}$  and the inequality is straight. The maximum is taken at any  $x_{[d]}$  if and only if for any  $l$  in the support of  $\theta$  we have  $\sigma^{\mathcal{F}}[X_{[d]} = x_{[d]}, L = l] = \mathbb{I}_{>0}(\theta[L = l])$ . This condition is equal to  $x_{[d]}$  being a model of  $\neg(1 - \mathbb{I}_{>0}(\theta[L = l])) f_l$  for any  $l$  in the support of  $\theta$  and thus to being a model of  $f_{\mathcal{F}, \theta}$ .

For the second claim: The auxiliary boolean vector  $v$  serves as an indicator for the formulas  $\neg(1 - \mathbb{I}_{>0}(\theta[L = l])) f_l$  being satisfied. In the worlds, for which  $\max_{x_{[d]}} \langle \theta[L], \sigma^{\mathcal{F}}[X_{[d]} = x_{[d]}, L] \rangle [\emptyset]$  is attained, we have

$$\langle \theta, v \rangle [\emptyset] \in \max_{\mu} \langle \theta, \mu \rangle [\emptyset] .$$

They are the models of  $f_{\mathcal{F}, \langle \theta, v \rangle [L]}$ . Note that then, all formulas  $f_l$  with  $v[L = l] = 0$  need to be contradicted from  $f_{\mathcal{F}, \langle \theta, v \rangle [L]}$ , since otherwise the maximum would be larger. The mean parameters solving the optimization problem are convex combinations of these extreme points, which correspond with distributions being convex combinations of the one-hot encodings of these models. Such distributions are therefore representable by the base measure being the indicator of these models, which is

$$\bigvee_{v[L] : \langle \theta, v \rangle [\emptyset] \in \max_{\mu} \langle \theta, \mu \rangle [\emptyset]} f_{\mathcal{F}, \langle \theta, v \rangle [L]} .$$

□

### 12.1.5 Expressivity of Hybrid Logic Networks

To summarize the above, we have characterized all the mean parameter vectors in  $\mathcal{M}_{\mathcal{F}, \nu}$  by their reproducing distributions, as sketched in Figure 33. Extreme points in  $\mathcal{M}_{\mathcal{F}, \nu}$  are boolean mean parameter vectors and by The. 70 exactly those reproduced by hard logic networks. Points, which are neither extreme nor in the interior of  $\mathcal{M}_{\mathcal{F}, \nu}$ , are reproduced by hybrid logic networks with a refined base measure  $\tilde{\nu}$ . If the statistic  $\mathcal{F}$  is minimal with respect to  $\nu$ , then  $\tilde{\nu}$  does not coincide with  $\nu$ . Finally, interior points are by Cor. 8 those reproduced by a hybrid logic network with base measure  $\nu$ .

Let us recall, that the set  $\Lambda^{\mathcal{F}, \text{EL}}$  contains all Hybrid Logic Networks, which can be realized as tensor networks with the same structure of computation and activation cores. We now investigate, whether we can reduce the set of probability distributions in the definition of the convex polytope of mean parameters to the set  $\Lambda^{\mathcal{F}, \text{EL}}$  (see Def. 60), that is

$$\mathcal{M}_{\mathcal{F}, \nu} |_{\Lambda^{\mathcal{F}, \text{EL}}} = \{ \langle \mathbb{P}, \sigma^{\mathcal{F}} \rangle [L] : \mathbb{P} \in \Lambda^{\mathcal{F}, \text{EL}} \} .$$

While  $\mathcal{M}_{\mathcal{F},\nu}|_{\Lambda^{\mathcal{F},\text{EL}}} \subset \mathcal{M}_{\mathcal{F},\nu}$  is obvious, we pose the question, for which  $\mathcal{F}$  there is an equivalence. We will refer to the equality of both sets as sufficient expressivity of  $\Lambda^{\mathcal{F},\text{EL}}$ . In the next example we provide a class of formulas, for which  $\Lambda^{\mathcal{F},\text{EL}}$  does not have sufficient expressivity.

**Example 13** (Insufficient expressivity of  $\Lambda^{\mathcal{F},\text{EL}}$  in cases of disjoint models). *To provide an example, where the set of hybrid logic networks does not suffice to reproduce all possible mean parameters, consider the formulas*

$$f_0 = X_0 \wedge X_1 \quad , \quad f_1 = \neg X_0 \wedge \neg X_1 .$$

*The probability distributions on the facet with normal  $\theta = [1 \ 1]$  are those with support on the models of  $f_0 \vee f_1$ . The Hybrid Logic Networks can only reproduce those which are supported on the model of  $f_0$  or the model of  $f_1$ , but not their convex combinations.*

*More generally, we can construct similar examples by arbitrary sets of formulas with pairwise disjoint model sets. If they do not sum to  $\mathbb{I}$ , i.e. there is a world which is not a model to any formula, the statistic is minimal. The vector  $\theta[L] = \mathbb{I}[L]$  is then the normal of the facet with all probabilities supported on the models of the formulas have mean parameters on this facet.*

Before presenting the example class of atomic formulas as a case, where  $\Lambda^{\mathcal{F},\text{EL}}$  has sufficient expressivity, let us first prove a generic criterion.

**Theorem 73.** *The set of mean parameters realizable by  $\Lambda^{\mathcal{F},\text{EL}}$  coincides with the set of mean parameters realizable by any distribution, if for any boolean vector  $\theta[L]$  the formula  $f_{\mathcal{F},\theta}$  is satisfiable.*

*Proof.* If these formulas are satisfiable, the base measures to the facets coincides with those realizable by  $\Lambda^{\mathcal{F},\text{EL}}$ .  $\square$

When the assumptions of The. 73 are not satisfied, there are mean parameters, which can not be reproduced by a distribution in  $\Lambda^{\mathcal{F},\text{EL}}$ . In that case, we can flexibilize the distribution, to also represent the base measures used for refinement in The. 72. This can be done by adding activation cores with multiple variables, or further computation cores calculating the disjunctions of formulas.

To characterize all the pre-images of faces, which are in  $\Lambda^{\mathcal{F},\text{EL}}$ , we exploit the faces of the  $p$ -dimensional cube. The cube  $[0, 1]^p$  is a polytope in  $\mathbb{R}^p$ , namely the convex hull of all boolean vectors

$$[0, 1]^p = \text{conv}(\mu[L] : \forall l \in [p] : \mu[L = l] \in \{0, 1\}) .$$

We have for the statistic of atoms  $X_{[d]}$ , that  $\mathcal{M}_{X_{[d]},\mathbb{I}} = [0, 1]^p$ . Its faces are enumerated by tuples  $(A, x_A)$ , where  $A \subset [p]$  and  $x_A \in \{0, 1\}^{|A|}$ .

$$Q_{S,\nu}^{(A,x_A)} = \{\mu[L] : \mu \in [0, 1]^p, \forall l \in A : \mu[L = l] = x_l\} .$$

**Theorem 74.** *The pre-image of faces of  $\mathcal{M}_{\mathcal{F},\nu}$ , which subset encodings are in  $\Lambda^{\mathcal{F},\text{EL}}$ , are those faces, which are intersections of  $\mathcal{M}_{\mathcal{F},\nu}$  with faces of the cube  $[0, 1]^p$ .*

*Proof.* Any subset encoding is a boolean tensor and any boolean tensor in  $\Lambda^{\mathcal{F},\text{EL}}$  has a representation with boolean activation core. If the activation core is elementary, it has a basis+ decomposition with rank 1, since we have leg dimensions of 2 for boolean statistics. Let  $A \subset [p]$  be the legs, which vector is a basis vector, and let  $x_A \in \{0, 1\}^{|A|}$  be the tuple storing to  $l \in A$  the number of the basis vector by  $x_l$ . Therefore, we have a parametrization of the boolean tensors in  $\Lambda^{\mathcal{F},\text{EL}}$  by

$$\Lambda^{\mathcal{F},\text{EL}} \cup \bigotimes_{k \in [d]} \{0, 1\}^2 = \bigcup_{A \subset [p]} \bigcup_{x_A \in \{0, 1\}^{|A|}} \langle \beta^{\mathcal{F}}[Y_{[p]}, X_{[d]}], \epsilon_{x_A}[Y_A] \rangle [X_{[d]}]$$

For any tuple  $(A, x_A)$  and any  $x_{[d]}$  we then have

$$\langle \beta^{\mathcal{F}}[Y_{[p]}, X_{[d]}], \epsilon_{x_A}[Y_A] \rangle [X_{[d]} = x_{[d]}] = 1 \quad \Leftrightarrow \quad \forall l \in A : \sigma^S[X_{[d]} = x_{[d]}, L = l] = x_l$$

And thus

$$\{x_{[d]} : \langle \beta^{\mathcal{F}}[Y_{[p]}, X_{[d]}], \epsilon_{x_A}[Y_A] \rangle [X_{[d]} = x_{[d]}] = 1\} = \mathcal{M}_{\mathcal{F},\nu} \cup Q_{S,\nu}^{(A,x_A)} .$$

Since to each boolean tensors in  $\Lambda^{\mathcal{F},\text{EL}}$  and to each face in Note, that there might be multiple pairs  $(A, x_A)$  to a boolean tensor in  $\Lambda^{\mathcal{F},\text{EL}}$ , corresponding to situations where multiple faces of the cube have identical intersections with  $\mathcal{M}_{\mathcal{F},\nu}$ .  $\square$



The vertices of  $\mathcal{M}_{\mathcal{F},\nu}$  consist of  $\{\mu\}$  where  $\mu[L] = \sigma^S[X_{[d]} = x_{[d]}, L]$  for an  $x_{[d]}$ . As a corollary of The. 74, their pre-image is always in  $\Lambda^{\mathcal{F},\text{EL}}$ .

**Corollary 9.** *All encoded pre-images of vertices are in  $\Lambda^{\mathcal{F},\text{EL}}$ .*

*Proof.* Since vertices in  $\mathcal{M}_{\mathcal{F},\nu}$  are also vertices of the cube. □

Pre-images of vertices have a representation with basis tensors as activation tensors.

**Corollary 10.** *When  $\mathcal{F}$  is the set of atomic formulas, all encoded pre-images of faces are in  $\Lambda^{\mathcal{F},\text{EL}}$ .*

*Proof.* Since for atomic formulas, the mean parameter polytope is the cube  $[0, 1]^p$ . □

For non-boolean features, we can derive limited versions of the above statement. Any face, which is the intersection with faces of the distorted cube

$$\bigotimes_{l \in [p]} [\min_x \mathcal{S}_l[X = x], \max_x \mathcal{S}_l[X = x]]$$

has a representation of its encoded pre-image with an activation tensor of basis+ rank 1. Since for leg dimensions larger than 2 the boolean tensors with elementary decomposition contain more tensors than the basis+ tensors of rank 1, we can represent further faces.

**Theorem 75.** *The encoded pre-image of any face of rank  $r$  is represented by a activation core of basis CP rank  $r$ .*

*Proof.* We build for each vertex  $\{\mu\}$  a basis vector  $\bigotimes_{l \in p} \epsilon_{\mu[L=l]} [Y_l]$ , which contraction with  $\beta^{\mathcal{F}}$ . Note, that the vertex base measures have disjoint support and their disjunction is thus a summation. □

### 12.1.6 Case of tree computation networks

In this case, the mean polytope can be embedded into a markov network and characterized by local consistency of the mean parameters of the markov network.

### 12.1.7 Examples

We can relate our two standard examples of the atomic and the minterm formula sets to well-studied polytopes, namely the  $d$ -dimensional hypercube and the standard simplex (see Lecture 0 in Ziegler (2013) )

**Example 14** (Atomic formulas). *The assumption of The. 73 is satisfied in the case for atomic formulas, where the formulas  $f_{\mathcal{F},\theta}$  are the minterms, which are always satisfiable in exactly one situation. The mean polytope in this case is the  $d$ -dimensional hypercube,*

$$\mathcal{M}_{\mathcal{F}_{[d]},\mathbb{I}} = [0, 1]^d$$

*which is called a simple polytope, since each vertex is contained in the minimal number of  $d$  facets.*

**Example 15** (Minterm formulas). *The mean polytope is in the case of the minterm exponential family is the  $2^d - 1$ -dimensional standard simplex. In this case,  $\Lambda^{\mathcal{F}_{\wedge},\text{EL}}$  contains any distribution and therefore trivially realizes any mean parameter in  $\mathcal{M}_{\mathcal{F}_{\wedge},\mathbb{I}}$ .*

## 12.2 Entropic Motivation of unconstrained Parameter Estimation

Markov Logic Networks are exponential families with statistics by a set  $\mathcal{F}$  of propositional formulas. We furthermore allow for propositional formulas as base measures, to also include the discussion of Hybrid Logic Networks. Based on this, we apply the theory of probabilistic inference, developed in Chapter 6 for the generic exponential families.

### 12.2.1 Maximum Likelihood in Hybrid Logic Networks

The Maximum Likelihood Problem on a hybrid logic network family is the moment projection

$$\operatorname{argmax}_{\theta[L] \in \mathbb{R}^p} \quad \mathbb{H} \left[ \mathbb{P}, \mathbb{P}^{(\mathcal{S}, \theta, \nu)} \right]$$

in the case  $\mathbb{P} = \mathbb{P}^D$  for a sample selector map  $D$ .

The moment projection coincides, after dropping constant terms in case of non-trivial base measure, with the backward map

$$\operatorname{argmax}_{\theta[L] \in \mathbb{R}^p} \quad \langle \theta[L], \mu[L] \rangle [\emptyset] - A^{(\mathcal{S}, \nu)}(\theta[L])$$

where

$$\mu[L] = \langle \sigma^{\mathcal{F}}, \mathbb{P} \rangle [L] \quad \text{and} \quad A^{(\mathcal{S}, \nu)}(\theta[L]) = \langle \exp [\langle \sigma^{\mathcal{F}} [X_{[d]}, L], \theta[L] \rangle [X_{[d]}]], \nu \rangle [\emptyset] .$$

We now extend to Hybrid Logic Networks

$$\operatorname{argmax}_{\tilde{\mathbb{P}} \in \Lambda^{\mathcal{F}, \text{EL}}} \quad \mathbb{H} [\tilde{\mathbb{P}}, \tilde{\mathbb{P}}]$$

**Corollary 11.** *Let  $\mu[L] = \langle \mathbb{P}, \sigma^{\mathcal{F}} \rangle [L]$  and*

$$\tilde{\mathcal{F}} = \{f_l : \mu[L = l] \in \{0, 1\}\} \quad , \quad \nu^{\tilde{\mathcal{F}}, \mu} [X_{[d]}] = \bigwedge_{f_l \in \tilde{\mathcal{F}}} \neg^{(1 - \mu[L=l])} f_l [X_{[d]}] .$$

*If  $\mu[L]$  is reproduceable by a positive distribution with respect to  $\nu^{\tilde{\mathcal{F}}, \mu} [X_{[d]}]$ , then the solution of the M-projection of  $\mathbb{P}$  onto the set of hybrid logic networks is representable by  $\mathcal{F}$  then coincides with the projection of  $\mathbb{P}$  onto  $\Gamma^{\mathcal{F}/\tilde{\mathcal{F}}, \nu^{\tilde{\mathcal{F}}, \mu}}$ .*

### 12.2.2 Maximum Entropy in Hybrid Logic Networks

The Maximum Entropy Problem for a statistic  $\mathcal{F}$  is

$$\operatorname{argmax}_{\mathbb{P}} \quad \mathbb{H} [\mathbb{P}] \quad \text{subject to} \quad \langle \mathbb{P}, \sigma^{\mathcal{F}} \rangle [L] = \mu[L] \quad (26)$$

**Corollary 12.** *Let  $\mathbb{P}^D$  be a distribution such that there is a positive distribution  $\mathbb{P}$  with  $\langle \mathbb{P}, \sigma^{\mathcal{F}} \rangle [L] = \langle \mathbb{P}^D, \sigma^{\mathcal{F}} \rangle [L]$ . Among all positive distributions  $\mathbb{P}$  of  $\times_{k \in [d]} [2]$  satisfying this moment matching condition, the Markov Logic Network with formulas  $\mathcal{F}$  and weights  $\theta$  being the solution of the maximum likelihood problem has minimal entropy.*

We notice, that the solution of the maximum entropy problem is thus a Markov Logic Network. This is remarkable, because this motivates our restriction to Markov Logic Networks as those distributions with maximal entropy given satisfaction rates of formulas in  $\mathcal{F}$ .

When now extend to the situations  $\mu[L = l] \in \{0, 1\}$  can appear. In that case the formula is entailed or contradicted by the facts, and dropping should be considered in both cases.

The max entropy - max likelihood duality still holds for hybrid logic networks as we show in the next theorem.

**Theorem 76.** *Given a set of formulas  $\tilde{\mathcal{F}}$  and  $\tilde{\mu}$ , with coordinates  $\tilde{\mu}_l \in [0, 1]$  in the closed interval  $[0, 1]$ . If the corresponding maximum entropy problem is feasible, its solution is a hybrid logic network with*

- $\mathcal{KB} = \{f_l : l \in [p], \mu[L = l] = 1\} \cup \{\neg f_l : l \in [p], \mu[L = l] = 0\}$
- $\mathcal{F} = \{f_l : l \in [p], \mu[L = l] \in (0, 1)\}$
- $\theta$  being the backward map evaluated at the vector  $\mu$  consisting of the coordinates of  $\tilde{\mu}$  not in  $\{0, 1\}$

*Proof.* Feasible distributions have a density with base measure by  $\mathcal{KB}$ , we therefore reduce the set of distributions in the argmax to those with density to the base measure. The max entropy is a max entropy problem with respect to that base measure, where we only keep the constraints to the mean parameters different from  $\{0, 1\}$  (those are trivially satisfied). The statement then follows from the generic property (see Sec3.1 in Wainwright and Jordan (2008)).  $\square$

### 12.3 Alternating Algorithms to Approximate the Backward Map

Let us now introduce an implementation of the Alternating Moment Matching Algorithm 5 in case of Markov Logic Networks. To solve the moment matching condition at a formula  $f_l$  we refine Lem. 8 in the following.

**Lemma 18.** *Let there be a base measure  $\nu$ , a formula selecting map  $\mathcal{F} = \{f_l : l \in [p]\}$  and weights  $\theta$ , and choose  $l \in [p]$  such that  $f_l \notin \{\mathbb{I} [X_{[d]}], 0 [X_{[d]}]\}$ . The moment matching condition relative to  $\theta$ ,  $l \in [p]$  and  $\mu_D [L = l] \in (0, 1)$  is then satisfied, if*

$$\theta[L = l] = \ln \left[ \frac{\mu_D [L = l]}{(1 - \mu_D [L = l])} \cdot \frac{\tau [X_{f_l} = 0]}{\tau [X_{f_l} = 1]} \right] \quad (27)$$

where by  $\tau[X_{f_l}]$  we denote the contraction

$$\tau[X_{f_l}] = \left\langle \{\beta^{f_l} : l \in [p]\} \cup \{\alpha^{\tilde{l}} : \tilde{l} \in [p], \tilde{l} \neq l\} \cup \{\nu\} \right\rangle [X_{f_l}] .$$

*Proof.* Since  $\text{im}(f_l) \subset [2]$  we have

$$\text{Id}|_{\text{im}(f_l)} = \epsilon_1 [X_{f_l}]$$

and the moment matching condition is by Lem. 8 satisfied if

$$\langle \alpha^l, \epsilon_1, \tau \rangle [\emptyset] = \langle \alpha^l, \tau \rangle [\emptyset] \cdot \mu_D [L = l] .$$

This is equal to

$$\exp[\theta [L = l]] \cdot \tau[X_{f_l} = 1] = (\exp[\theta [L = l]] \cdot \tau[X_{f_l} = 1] + \tau[X_{f_l} = 0]) \cdot \mu_D [L = l] .$$

Rearranging the equations this is equal to

$$\tau[X_{f_l}] = \left\langle \{\beta^{f_l}\} \cup \{\alpha^{\tilde{l}} : \tilde{l} \in [p], \tilde{l} \neq l\} \cup \{\nu\} \right\rangle [L] .$$

We notice that the right side is well defined, since we have by assumption  $\mu_D [L = l], (1 - \mu_D [L = l]) \neq 0$  and  $\tau[X_{f_l} = 0], \tau[X_{f_l} = 1] \neq 0$  since Markov Logic networks are positive distributions and  $f_l \notin \{\mathbb{I}[X_{[d]}], 0[X_{[d]}]\}$ .  $\square$

In the case  $\mu_D [L = l] \in \{0, 1\}$  the moment matching conditions are not satisfiable for  $\theta [L = l] \in \mathbb{R}$ . But, we notice, that in the limit  $\theta [L = l] \rightarrow \infty$  (respectively  $-\infty$ ) we have

$$\mu [L = l] \rightarrow 1 \quad (\text{respectively } 0) ,$$

and the moment matching can be satisfied up to arbitrary precision. In Sect. 11.2 we will allow infinite weights and interpret the corresponding factors by logical formulas. As a consequence, we will be able to fit graphical models, which we will call hybrid networks on arbitrary satisfiable mean parameters.

The cases  $\tau[X_{f_l} = 1] = 0$ , respectively  $\tau[X_{f_l} = 0] = 0$  only appear for nontrivial formulas when the distribution is not positive. This is not the case for Markov Logic Networks, but will happen when formulas are added as cores of a Markov Network. This situation has been investigated in Sect. 11.2.

Since the likelihood is concave (see Koller and Friedman (2009)), there are not local maxima the coordinate descent could run into and coordinate descent will give a monotonic improvement of the likelihood.

We suggest an alternating optimization by Algorithm 9, solving the moment matching equation iteratively for all formulas  $f \in \mathcal{F}$  and repeat the optimization until a convergence criterion is met. This is an coordinate ascent algorithm, when interpreted the loss  $\mathcal{L}_D(\mathbb{P}^{(\mathcal{S}, \theta, \nu)})$  as an objective depending on the vector  $\theta$ .

In the initialization phase of Algorithm 9, each parameters is initialized relative to a uniform distribution. The algorithm would be finished, if the variables  $X_f$  are independent. This would be the case, if the Markov Logic Network consists of atomic formulas only. When they fail to be independent, the adjustment of the weights influence the marginal distribution of other formulas and we need an alternating optimization. This situation corresponds with couplings of the weights by a partition contraction, which does not factorize into terms to each formula.

**Solving Equation 27 requires inference of a current model by answering a query.** This can be a bottleneck and circumvented by approximative inference, see e.g. CAMEL Ganapathi et al. (2008).

**Remark 14** (Grouping of coordinates with trivial sum). *When having a set of coordinates, such that the coordinate functions are binary and sum to the trivial tensor, one can find simultaneous updates to the canonical parameters, such that the partition function is staying invariant. Given a parameter  $\theta^t$  we compute*

$$\mu^t = \left\langle \mathbb{P}^{(\mathcal{S}, \theta^t)}, \mathcal{S} \right\rangle [L]$$

and build the update

$$\theta^{t+1} = \theta^t + \ln[\mu^D] \mu^t .$$

Then,  $\theta^{t+1}$  satisfies the moment matching equations for all coordinates in the set.

The assumptions are met when taking all features to any hyperedge in a Markov Network seen as an exponential family. In that case, the update algorithm is referred to as Iterative Proportional Fitting Wainwright and Jordan (2008). Further, when activating both  $f$  and  $\neg f$ .

**Algorithm 9** Alternating Weight Optimization (AWO)**Require:** Empirical distribution  $\mathbb{P}^D$ , boolean features  $\mathcal{F}$ **Ensure:** Canonical parameter  $\theta[L]$ , such that  $\mathbb{P}^{(\mathcal{S}, \theta, \nu)}$  is the (approximative) moment projection of  $\mathbb{P}^D$  onto  $\Gamma^{\mathcal{S}, \nu}$ Compute  $\mu_D[L] = \langle \mathbb{P}^D, \sigma^{\mathcal{S}} \rangle [L]$  $\mathcal{KB} = \mathbb{I}, \tilde{\mathcal{V}} = \emptyset$ **for**  $l \in [p]$  **do**    **if**  $\mu[L = l] = 1$  **then**         $\mathcal{KB} \leftarrow \mathcal{KB} \cup \{f_l\}$     **else if**  $\mu[L = l] = 0$  **then**         $\mathcal{KB} \leftarrow \mathcal{KB} \cup \{\neg f_l\}$     **else**         $\tilde{\mathcal{V}} \leftarrow \tilde{\mathcal{V}} \cup \{l\}$     **end if****end for****for**  $l \in \tilde{\mathcal{V}}$  **do**

Compute

$$\tau[X_{f_l}] \leftarrow \langle \beta^{f_l} \rangle [X_{f_l}]$$

Set

$$\theta[L = l] \leftarrow \ln \left[ \frac{\mu_D[L = l]}{(1 - \mu_D[L = l])} \cdot \frac{\tau[X_{f_l} = 0]}{\tau[X_{f_l} = 1]} \right]$$

**end for****if**  $\langle \mathcal{KB} \rangle [\emptyset] = 0$  **then**    **raise** "Inconsistent Knowledge Base"**end if****while** Convergence criterion is not met **do**    **for**  $l \in \tilde{\mathcal{V}}$  **do**

Compute

$$\tau[X_{f_l}] = \left\langle \{\beta^{f_l} : l \in [p]\} \cup \{\alpha^{\tilde{l}} : \tilde{l} \in [p], \tilde{l} \neq l\} \cup \{\nu\} \right\rangle [X_{f_l}]$$

Set

$$\theta[L = l] = \ln \left[ \frac{\mu_D[L = l]}{(1 - \mu_D[L = l])} \cdot \frac{\tau[X_{f_l} = 0]}{\tau[X_{f_l} = 1]} \right]$$

**end for****end while****return**  $\theta[L]$ **12.4 Forward and backward mappings in closed form**

We recall from Chapter 6, that while forward mappings are always in closed form by contractions, backward mapping in general do not have a closed form representation. Instead, the backward map is in general implicitly characterized by a maximum entropy problem constrained to matching expected sufficient statistics. We investigate in this section specific examples, where closed forms are available for both. In these cases, parameter estimation can thus be solved by application of the inverse on the expected sufficient statistics with respect to the empirical distribution, and iterative algorithms can be avoided.

**12.4.1 Maxterms and Minterms**

Minterms (respectively maxterms) are ways in propositional logics to get a syntactical formula representation based on a formula to each world which is a model (respectively fails to be a model). We have already studied in Sect. 11.1.3 how to represent any distribution as a MLN of maxterms (respectively minterms), see The. 60.

We use the tuple enumeration of the maxterms and minterms by  $\times_{k \in [d]} [2]$  introduced in Sect. 7.3.3. With respect to this enumeration the canonical parameters and mean parameters are tensors in  $\otimes_{k \in [d]} \mathbb{R}^2$ . Since the statistic of the

minterm family is the identity, the mean parameters for the minterm family are

$$\mu [L_{[d]} = x_{[d]}] = \mathbb{P} [x_{[d]}]$$

and therefore after a relabeling of categorical variables to selection variables  $\mu = \mathbb{P}$ . For maxterms we have analogously

$$\mu [L_{[d]} = x_{[d]}] = 1 - \mathbb{P} [x_{[d]}]$$

and  $\mu = \mathbb{I} - \mathbb{P}$ . We can use these insights to provide a characterization of the forward and backward maps of the minterm and maxterm family.

**Theorem 77.** *Given the Markov Logic Networks to the formula sets*

$$\mathcal{F}_{\wedge} := \{Z_{x_0, \dots, x_{d-1}}^{\wedge} : x_0, \dots, x_{d-1} \in \bigtimes_{k \in [d]} [2]\} \quad \text{and} \quad \mathcal{F}_{\vee} := \{Z_{x_0, \dots, x_{d-1}}^{\vee} : x_0, \dots, x_{d-1} \in \bigtimes_{k \in [d]} [2]\}$$

of all minterms, respectively of all maxterms, the forward mapping are

$$F^{\wedge}(\theta) = \langle \exp [\theta] \rangle [X_{[d]} | \emptyset] \quad \text{and} \quad F^{\vee}(\theta) = \langle \exp [-\theta] \rangle [X_{[d]} | \emptyset] ,$$

where in a slight abuse of notation we assigned the variables  $X_{[d]}$  to the canonical parameters  $\theta$ .

Possible choices of the backward mappings are

$$B^{\wedge}(\mu) = \ln [\mu] \quad \text{and} \quad B^{\vee}(\mu) = -\ln [\mu] .$$

*Proof.* For the minterms we use that

$$\mathcal{F}_{\wedge}[X_{[d]}, X_{\mathcal{F}_{\wedge}}] = \delta [X_{[d]}, X_{\mathcal{F}_{\vee}}]$$

and get

$$F^{\wedge}(\theta) = \langle \exp [\langle \{\mathcal{F}_{\wedge}, \theta\} \rangle [X_{[d]}]] \rangle [X_{[d]} | \emptyset] = \langle \exp [\theta] \rangle [X_{[d]} | \emptyset] .$$

We notice that for any  $\mu$  in the image of the forward map we have

$$F^{\wedge}(B^{\wedge}(\mu)) = \mu$$

Therefore,  $B^{\mathcal{F}_{\wedge}}$  is indeed a backward mapping to the exponential family of minterms.

For the maxterms we use that

$$\mathcal{F}_{\vee}[X_{[d]}, X_{\mathcal{F}_{\vee}}] = \mathbb{I} [X_{[d]}, X_{\mathcal{F}_{\vee}}] - \delta [X_{[d]}, X_{\mathcal{F}_{\vee}}]$$

and get

$$\begin{aligned} F^{\vee}(\theta) &= \langle \exp [\langle \{\mathcal{F}_{\vee}, \theta\} \rangle [X_{[d]}]] \rangle [X_{[d]} | \emptyset] \\ &= \langle \{ \exp [\langle \{\mathbb{I}, \theta\} \rangle [X_{[d]}]] , \exp [-\langle \theta \rangle [X_{[d]}]] \} \rangle [X_{[d]} | \emptyset] \\ &= \langle \exp [-\theta] \rangle [X_{[d]} | \emptyset] \end{aligned}$$

where we used, that  $\exp [\langle \{\mathbb{I}, \theta\} \rangle [X_{[d]}]]$  is a multiple of  $\mathbb{I} [X_{[d]}]$  and is thus eliminated in the normalization. For any  $\mu \in \text{im} (F^{\vee})$  we have

$$F^{\vee}(B^{\vee}(\mu)) = \mu$$

and  $B^{\wedge}$  is thus a backward map for the exponential family of maxterms.  $\square$

Any positive probability distribution can thus be fitted by minterms when we choose  $\theta = \ln [\mathbb{P}]$ , respectively by maxterms when we choose  $\theta = \mathbb{I} - \ln [\mathbb{P}]$ . Thus, we have identified a subset of  $2^d$  formulas, which is rich enough to fit any distribution.

## 12.4.2 Atomic formulas

Let us now derive a closed form backward mapping for the statistic

$$\mathcal{F}_{[d]} := \{X_k : k \in [d]\} .$$

The mean parameters coincide with the queries on the atomic formulas, that is the marginal

$$\mu [L = k] = \mathbb{P} [X_k = 1] .$$

**Theorem 78.** *Given a Markov Logic Network with the statistic  $\mathcal{F}_{[d]}$  of atomic formulas, the forward mapping from canonical parameters to mean parameters is the coordinatewise sigmoid, that is*

$$F^{[d]}(\theta[L]) = \frac{\exp[\theta[L]]}{\mathbb{I}[L] + \exp[\theta[L]}}$$

where the quotient is performed coordinatewise.

A backward mapping is the coordinatewise logit, that is

$$B^{[d]}(\mu[L]) = \ln \left[ \frac{\mu[L]}{\mathbb{I}[L] - \mu[L]} \right].$$

*Proof.* We have for any  $\theta[L] \in \mathbb{R}^d$

$$\mathbb{P}^{(\mathcal{F}_{[d]}, \theta)}[X_{[d]}] = \bigotimes_{k \in [d]} \langle \exp[\theta[L = k]] \cdot X_k \rangle [X_k | \emptyset].$$

For any  $k \in [d]$  it therefore holds, that

$$\begin{aligned} F^{[d]}(\theta[L])[L = k] &= \left\langle X_k, \mathbb{P}^{(\mathcal{F}_{[d]}, \theta)}[X_{[d]}] \right\rangle [\emptyset] \\ &= \langle X_k, \langle \exp[\theta[L = k]] \cdot X_k \rangle [X_k | \emptyset] \rangle [\emptyset] \\ &= \frac{\exp[\theta[L = k]]}{1 + \exp[\theta[L = k]]}. \end{aligned}$$

Since the coordinatewise logit is the inverse function of the coordinatewise sigmoid the map

$$B^{[d]}(\mu[L])[L = k] = \ln \left[ \frac{\mu[L = k]}{1 - \mu[L = k]} \right]$$

satisfies for any  $\mu$  in the image of the forward map

$$F^{[d]}(B^{[d]}(\mu)) = \mu$$

and is therefore a backward map. □

In a selection tensor networks they are represented by a single neuron with identity connective and variable selection to all atoms. We will investigate such examples in more detail in Chapter 18, where atomic formulas Markov Logic Networks are specific cases of monomial decomposition of order 1.

The maximum likelihood estimator of a positive probability distribution by the MLN of atomic formulas is therefore the tensor product of the marginal distributions. The Kullback-Leibler divergence between the distribution and its projection is the mutual information of the atoms, see for example Chapter 8 in MacKay (2003).

**Remark 15** (Decomposition into systems of atomic networks). *By Independence Decomposition we reduce to a system of atomic MLN. The minterms of such MLNs are the literals. By redundancy (literals sum up to  $\mathbb{I}$ ), it suffices to take only the positive or the negative literal.*

## 12.5 Constrained parameter estimation in the minterm family

We approach structure learning as constrained parameter estimation in the naive exponential family (see Sect. 5.3.6), which coincides with the minterm family  $\mathcal{F}_{\wedge}$ . The minterm family is defined by the statistic  $\mathcal{S} = \delta[X_{[d]}, L_{[d]}]$  and has energy tensors coinciding with the canonical parameters.

For the minterm family, we have as mean parameter set the convex hull of one-hot encodings. Each basis vector is an extreme point is an extreme point.

By The. 60 all positive distributions are member of the minterm markov logic network family. This expressivity result was generalized to arbitrary distributions, when allowing for formulas as basemeasures by The. 68.

Finding the distribution maximizing the likelihood of data would then be the empirical distribution. In this case we would have  $\mu_D[L_{[d]} = x_{[d]}] = \mathbb{P}^D[X_{[d]} = x_{[d]}]$  and the maximum likelihood distribution is found by the problem

$$\operatorname{argmax}_{\theta \in \bigotimes_{k \in [d]} \mathbb{R}^{m_k}} \langle \theta, \mathbb{P}^D \rangle [\emptyset] - A^{(\mathcal{S}, \nu)}(\theta)$$

which is solved at  $\theta = \ln [\mathbb{P}^D]$  with  $\mathbb{P}^{(\delta, \ln[\mathbb{P}^D])} = \mathbb{P}^D$ . This follows from  $\mathcal{L}_D(\mathbb{P}^{(\delta, \theta)}) = D_{\text{KL}}[\mathbb{P}^D || \mathbb{P}^{(\delta, \theta)}]$ , which is by Gibbs inequality minimized at  $\mathbb{P}^{(\delta, \theta)} = \mathbb{P}^D$ , which is the case for  $\theta = \ln [\mathbb{P}^D]$ .

We here allow for  $\ln [0] = -\infty$ , with the convention of  $\exp [-\infty] = 0$ , to handle datasets where specific worlds are not represented. **Better: Use The. 68 with basemeasure dropping non appearing data.**

To avoid this overfitting situation, we regularize by restricting the parameter to be a set  $\Theta \subset \bigotimes_{k \in [d]} \mathbb{R}^{m_k}$  and state

$$\operatorname{argmax}_{\theta \in \Theta} \langle \theta, \mathbb{P}^D \rangle [\emptyset] - A^{(S, \nu)}(\theta). \quad (\text{P}_{\Theta, \mathbb{P}^D})$$

Problem  $\text{P}_{\Theta, \mathbb{P}^D}$  has two important types of instantiation, which we discuss in the next sections.

### 12.5.1 Parameter Estimation

**Projecting onto the markov logic family to the statistic  $\mathcal{F}$  is the instance of Problem ?? with the hypothesis choice**

$$\Theta^{\mathcal{F}} = \operatorname{span}(\{f : f \in \mathcal{F}\}).$$

Then, the problem is the parameter estimation problem studied in Sect. 12.2. To see this, we reparametrize by the coefficient vectors of the elements in the span, which are then understood as the canonical parameter of the respective distribution in the markov logic family to  $\mathcal{F}$ .

**Remark 16** (Overparametrization). *Taking  $\mathcal{F}$  to consist of all propositional formulas, we get a massive overparametrization: The essential statistics maps to a  $2^{(2^d)}$  dimensional real vector space. All possible distributions of the  $d$  atomic variables are mapped to an  $2^d - 1$  dimensional submanifold, where also the essential statistics maps to.*

*Thus, to identify probabilistic knowledge bases, we need to drastically restrict the shape of formulas allowed. It is in principle impossible to decide which formulas to be activated, based only on statistics and not on prior assumptions.*

*When having  $d$  atoms, there are  $2^d$  states in the factored system. Since each state can either be a model of a formula or not, there are*

$$|\mathcal{F}| = 2^{(2^d)}$$

*formulas. Having, for example,  $d = 10$ , then  $|\mathcal{F}| > 10^{308}$ .*

*One regularization is by allowing only a small number of formulas to be active. This corresponds with regularization with  $\ell_0(\theta)$ . The problem is then non-convex.*

*A further regularization strategy is the restriction of the size of the possible formulas to maintain interpretability. Thus, we choose small formula selection networks.*

### 12.5.2 Structure Learning

The problem of structure learning arises, when the set of parameters in Problem ?? is chosen as

$$\Theta^{\mathcal{H}} = \bigcup_{\mathcal{F} \in \mathcal{H}} \operatorname{span}(\mathcal{F}).$$

In this case, the problem in general fails to be convex.

Each formula set  $\mathcal{F}$  represents a subspace in the parameters of the minterm family, which is spanned by the propositional formulas  $f \in \mathcal{F}$ .

## 12.6 Greedy Structure Learning

It can be impracticable to learn all formulas at once, since the set  $\mathcal{H}$  often grows combinatorically, for example when choosing as a powerset of formulas. **Further, we need to avoid overfitting and carefully choose a hypothesis.** To avoid intractabilities and overfitting, one can choose a greedy approach and learn in addition formulas  $f$  when already having learned a set  $\mathcal{F}$  of formulas. We in this section assume a current model  $\tilde{\mathbb{P}}$ , which is a generic positive distribution not necessarily a Markov Logic Network.

We will use the effective selection tensor network representation of exponentially many formulas described in Chapter 10 and select from them a small subset.

### 12.6.1 Greedy formula inclusions

Having a current set of formulas  $\mathcal{F}$  we want to choose the best  $f \in \mathcal{H}$  to extend the set of formulas to  $\mathcal{F} \cup \{f\}$  in a way minimizing the cross entropy. Given this, add each step we solve the greedy cross entropy minimization

$$\operatorname{argmin}_{f \in \mathcal{H}} \operatorname{argmin}_{\theta \in \mathbb{R}^{|\mathcal{F}|+1}} \mathbb{H} \left[ \mathbb{P}^D, \mathbb{P}^{(\mathcal{F} \cup \{f\}, \theta, \nu)} \right]. \quad (\mathbf{P}_{D, \mathcal{F}, \mathcal{H}})$$

A brute force solution would require parameter estimation for each candidate in  $\mathcal{H}$ . We provide two more efficient approximative heuristics in the following (see Chapter 20 in Koller and Friedman (2009)).

### 12.6.2 Gain Heuristic

In the gain heuristic, only the parameters of the new formula are optimized and the others left unchanged. This amounts to

$$\operatorname{argmin}_{f \in \mathcal{H}} \left( \min_{\theta \in \mathbb{R}^{|\mathcal{F}|+1}} \mathbb{H} \left[ \mathbb{P}^D, \mathbb{P}^{(\mathcal{F} \cup \{f\}, \theta, \nu)} \right] \right). \quad (\mathbf{P}_{D, \mathcal{F}, \mathcal{H}}^{\text{gain}})$$

Here we denote by  $\theta$  the first  $|\mathcal{F}|$  coordinates of the M-projection  $\tilde{\mathbb{P}}$  of  $\mathbb{P}^D$  onto  $\mathcal{F}$  and the variable new coordinate at position  $\theta \in \mathbb{R}^{|\mathcal{F}|}$ .

**Lemma 19.** *The gain heuristic objective is an upper bound on the true greedy objective.*

*Proof.* Since

$$\begin{aligned} & \operatorname{argmin}_{f \in \mathcal{H}} \left( \operatorname{argmin}_{\theta \in \mathbb{R}^{|\mathcal{F}|+1}} \mathbb{H} \left[ \mathbb{P}^D, \mathbb{P}^{(\mathcal{F} \cup \{f\}, \theta, \nu)} \right] \right) \\ & \leq \operatorname{argmin}_{f \in \mathcal{H}} \left( \operatorname{argmin}_{\theta \in \mathbb{R}^{|\mathcal{F}|}} \mathbb{H} \left[ \mathbb{P}^D, \mathbb{P}^{(\mathcal{F} \cup \{f\}, \theta, \nu)} \right] \right). \end{aligned}$$

□

Further, this is Problem  $\mathbf{P}_{\Theta, \mathbb{P}^D}$  in the case

$$\Theta = \ln \left[ \tilde{\mathbb{P}} \right] + \cup_{f \in \mathcal{F}} \operatorname{span}(f).$$

Let us choose a formula  $f \in \mathcal{F}$  and consider Problem  $\mathbf{P}_{\Theta, \mathbb{P}^D}$  in the case

$$\Theta^f = \ln \left[ \tilde{\mathbb{P}} \right] + \operatorname{span}(f).$$

This is parameter estimation on the exponential family with the single feature  $f$  and the base measure  $\tilde{\mathbb{P}}$ . Therefore we can apply the theory of Chapter 6 and characterize the solution by the  $\theta$  satisfying the moment matching condition

$$\left\langle \tilde{\mathbb{P}}, \langle \exp[\theta] \rangle [X_{[d]}|\emptyset] \right\rangle [\emptyset] = \langle \mathbb{P}^D, f \rangle [\emptyset].$$

We state the solution of this condition in the next theorem.

**Theorem 79.** *Problem  $(\mathbf{P}_{D, \mathcal{F}, \mathcal{H}}^{\text{gain}})$  is solved at any*

$$\hat{\theta} = \theta_{\hat{f}} \cdot \hat{f}$$

where the formula  $\hat{f}$  is in

$$\hat{f} \in \operatorname{argmax}_{f \in \mathcal{F}} \mathbf{D}_{\text{KL}} \left[ \langle \mathbb{P}^D, f \rangle [\emptyset] \parallel \langle \tilde{\mathbb{P}}, f \rangle [\emptyset] \right]$$

and  $\theta_{\hat{f}}$  is the weight of  $\hat{f}$  in the solution of Problem  $\mathbf{P}_{\Theta, \mathbb{P}^D}$  with  $\Gamma = \tilde{\mathbb{P}} + \operatorname{span}(f)$ . Here we denote by  $\mathbf{D}_{\text{KL}}[p_1||p_2]$  the Kullback-Leibler divergence between Bernoulli distributions with parameters  $p_1, p_2 \in [0, 1]$ , that is

$$\mathbf{D}_{\text{KL}}[p_1||p_2] = p_1 \cdot \ln \left[ \frac{p_1}{p_2} \right] + (1 - p_1) \cdot \ln \left[ \frac{(1 - p_1)}{(1 - p_2)} \right]$$



*Proof.* For any formula  $f$ , the inner minimum of Problem  $(\mathbb{P}_{D,\mathcal{F},\mathcal{H}}^{\text{gain}})$  is by Lem. 18 taken at

$$\theta_f = \ln \left[ \frac{\mu_D}{(1 - \mu_D)} \cdot \frac{(1 - \tilde{\mu})}{\tilde{\mu}} \right]$$

where

$$\tilde{\mu} = \langle \tilde{\mathbb{P}}, f \rangle [\emptyset]$$

and

$$\mu_D = \langle \mathbb{P}^D, f \rangle [\emptyset] .$$

The difference of the likelihood at the current distribution and the optimum is

$$\mathbb{H} [\mathbb{P}^D, \tilde{\mathbb{P}}] - \mathbb{H} [\mathbb{P}^D, \mathbb{P}^{\mathcal{F} \cup \{f\}, \tilde{\theta} \cup \{\theta_f\}, \nu}] = \mu_D \cdot \theta_f - A^{\mathcal{F} \cup \{f\}, \nu} (\tilde{\theta} \cup \{\theta_f\}) .$$

We use the representation scheme of Theorem 65 and get

$$\begin{aligned} \langle \tilde{\mathbb{P}}, \exp [\theta_f \cdot f] \rangle [\emptyset] &= \langle \tilde{\mathbb{P}}, \beta^f [X_f], \alpha^f [X_f] \rangle [\emptyset] \\ &= (1 - \tilde{\mu}) + \tilde{\mu} \cdot \exp [\theta_f] \\ &= (1 - \tilde{\mu}) + \frac{\mu_D \cdot (1 - \tilde{\mu})}{(1 - \mu_D)} \\ &= (1 - \tilde{\mu}) \cdot \frac{1}{(1 - \mu_D)} . \end{aligned}$$

It follows, that

$$\begin{aligned} A^{\mathcal{F} \cup \{f\}, \nu} (\tilde{\theta} \cup \{\theta_f\}) &= \ln \left[ \langle \tilde{\mathbb{P}}, \exp [\theta_f \cdot f] \rangle [\emptyset] \right] \\ &= \ln [1 - \tilde{\mu}] - \ln [1 - \mu_D] . \end{aligned}$$

We further have

$$\mu_D \cdot \theta_f = \mu_D \cdot \left[ \ln \left[ \frac{\mu_D}{(1 - \mu_D)} \cdot \frac{(1 - \tilde{\mu})}{\tilde{\mu}} \right] \right] = \mu_D \ln [\mu_D] - \mu_D \ln [1 - \mu_D] + \mu_D \ln [1 - \tilde{\mu}] - \mu_D \ln [\tilde{\mu}]$$

and arrive at

$$\begin{aligned} \mathbb{H} [\mathbb{P}^D, \tilde{\mathbb{P}}] - \mathbb{H} [\mathbb{P}^D, \mathbb{P}^{(f, \theta_f, \tilde{\mathbb{P}})}] &= \mu_D \ln [\mu_D] - \mu_D \ln [1 - \mu_D] + \mu_D \ln [1 - \tilde{\mu}] - \mu_D \ln [\tilde{\mu}] - \ln [1 - \tilde{\mu}] - \ln [1 - \mu_D] \\ &= (-\mu_D \ln [\tilde{\mu}] - (1 - \mu_D) \ln [1 - \tilde{\mu}]) - (-\mu_D \ln [\mu_D] - (1 - \mu_D) \ln [1 - \mu_D]) . \end{aligned}$$

By definition, this is the Kullback-Leibler divergence between Bernoulli distributions with parameters  $\mu_D$  and  $\tilde{\mu}$ . Since the gain in the likelihood loss when restricting to  $\Theta = \text{span}(f)$  is thus given by  $D_{\text{KL}} [\langle \mathbb{P}^D, f \rangle [\emptyset] \parallel \langle \tilde{\mathbb{P}}, f \rangle [\emptyset]]$ , we have that Problem ?? in the case  $\Theta = \bigcup_{f \in \mathcal{F}} \text{span}(f)$  is solved at  $\hat{\theta} = \theta_{\hat{f}} \cdot \hat{f}$  where

$$\hat{f} = D_{\text{KL}} [\langle \mathbb{P}^D, f \rangle [\emptyset] \parallel \langle \tilde{\mathbb{P}}, f \rangle [\emptyset]] .$$

□

Thus, we solve the grain heuristic with a coordinatewise transform of the mean parameter tensors to  $\mathbb{P}^D$  and  $\tilde{\mathbb{P}}$ , using the Bernoulli Kullback-Leibler divergence as transform function.

One therefore takes the formula, which marginal distribution in the current model and the targeted distribution are differing at most, measured in the KL divergence.

One optimization method would thus be the computation of the mean parameters to both distribution, building the coordinatewise KL divergence and choosing the maximum. Since we need to evaluate each coordinate, this can be intractable for large sets of formulas.

Further improvement of the model can be achieved by iteratively optimizing the other weights as well, since their corresponding moment matching conditions might be violated after the integration of a new formula. This would require the computation of backward mappings for each candidate formula, for which we only have an alternating approach in general.

### 12.6.3 Gradient heuristic and the proposal distribution

**Advantage:** Might avoid formulawise calculus, when sampling from proposal distribution. Brute force solution of gain heuristic require formulawise approach.

We now derive a heuristic of choosing features based on the maximal coordinate of the gradient when differentiating the canonical parameter in the minterm family. To prepare for this, we build the gradient of the loss

$$\mathcal{L}_D \left( \mathbb{P}^{(\delta, \tilde{\theta})} \right) = \left\langle \mathbb{P}^D, \sigma^\delta, \tilde{\theta} \right\rangle [\emptyset] - \ln \left[ \left\langle \exp \left[ \left\langle \sigma^\delta, \tilde{\theta} \right\rangle [X_{[d]}] \right] \right\rangle [\emptyset] \right]$$

as

$$\begin{aligned} \nabla_{\tilde{\theta}[L]} \mathcal{L}_D \left( \mathbb{P}^{(\delta, \tilde{\theta})} \right) &= \left\langle \sigma^\delta, \mathbb{P}^D \right\rangle [L] - \left\langle \sigma^\delta, \mathbb{P}^{(\delta, \tilde{\theta})} \right\rangle [L] \\ &= \mathbb{P}^D - \mathbb{P}^{(\delta, \tilde{\theta})}. \end{aligned}$$

The gradient shows the typical decomposition into a positive and a negative phase. While the positive phase comes from the data term and prefers directions of large data support, the negative phase originates in the partition function and draws the gradient away from directions already supported by the current model  $\mathbb{P}^{(\delta, \tilde{\theta})}$ . The negative phase is a regularization, by comparing with what has already been learned. When nothing has been learned so far, we can take the current model to be the uniform distribution, which is the naive exponential family with vanishing canonical parameters.

Given a set  $\mathcal{H}$  of features we vary  $\tilde{\theta}$  by the function

$$q(\theta) = \tilde{\theta} + \langle \theta, \sigma^{\mathcal{H}} \rangle [X_{[d]}] .$$

At  $\theta = 0$  we have the gradient of the loss of the parametrized formula by

$$\begin{aligned} \nabla_{\theta|0} \mathcal{L}_D \left( \mathbb{P}^{(\delta, q(\theta), \nu)} \right) &= \left\langle \nabla_{q(\theta)|\tilde{\theta}} \mathcal{L}_D \left( \mathbb{P}^{(\delta, q(\theta), \nu)} \right), \nabla_{\theta|0} q(\theta) \right\rangle [\emptyset] \\ &= \left\langle \mathbb{P}^D, \sigma^{\mathcal{S}} \right\rangle [L] - \left\langle \mathbb{P}^{(\delta, \tilde{\theta}, \nu)}, \sigma^{\mathcal{S}} \right\rangle [L] . \end{aligned}$$

We want to choose the formula, which is best aligned with the gradient of the log-likelihood, that is using a formula selecting map  $\mathcal{H}$

$$\operatorname{argmax}_{l \in [p]} \left\langle \mathbb{P}^D, \mathcal{H} \right\rangle [L = l] - \left\langle \mathbb{P}^{(\delta, \tilde{\theta}, \nu)}, \mathcal{H} \right\rangle [L = l] . \quad (\mathbb{P}_{D, \mathcal{F}, \mathcal{H}}^{\text{grad}})$$

This method is known as the gradient heuristic or grafting. The objective of Problem  $(\mathbb{P}_{D, \mathcal{F}, \mathcal{H}}^{\text{grad}})$  has another interpretation by the difference of the mean parameter  $\mu_D$  and  $\tilde{\mu}$  of the projections of the empirical and current distributions on the family to  $\mathcal{H}$ .

Problem  $(\mathbb{P}_{D, \mathcal{F}, \mathcal{H}}^{\text{grad}})$  is further equivalent to the formula alignment

$$\operatorname{argmax}_{f \in \mathcal{H}} \left\langle f, \mathbb{P}^D - \tilde{\mathbb{P}} \right\rangle [\emptyset] .$$

The objective can be interpreted as the difference of the satisfaction probability of the formula with respect to the empirical distribution and the current distribution.

### 12.6.4 Iterations

Let us now iterate the search for a best formula at a current model with the optimization of weights after each step. The result is Algorithm 10, which is a greedy algorithm adding iteratively the currently best feature.

When having used the same learning architecture multiple times, the energy of the corresponding formulas are all representable by a formula selecting architecture. Their energy term is therefore a contraction of the selecting tensor with a parameter tensor  $\theta$  in a basis CP decomposition with rank by the number of learned formulas. When multiple selection architectures have been used, the energy is a sum of such contractions. Let us note, that this representation is useful after learning, when performing energy-based inference algorithms on the result. During learning, one needs to instantiate the proposal distribution, which requires instantiation of the probability tensor. **However, one could alternate data energy-based and use this as a particle-based proxy for the probability tensor.**

**Remark 17** (Sparsification by Thresholding). *To maintain a small set of active formulas, one could combine greedy learning approaches with thresholding on the coordinates of  $\theta$ . This is a standard procedure in Iterative Hard Thresholding algorithms of Compressed Sensing, but note that here we do not have a linear in  $\theta$  objective.*

**Algorithm 10** Greedy Structure Learning**Require:** Empirical distribution  $\mathbb{P}^D$ , hypothesis  $\mathcal{H}$  of formulas**Ensure:** Distribution  $\mathbb{P}^{(S, \theta, \nu)}$  approximating  $\mathbb{P}^D$ 

Initialize

$$\tilde{\mathbb{P}} \leftarrow \frac{1}{\prod_{k \in [d]} m_k} \cdot \mathbb{I}[X_{[d]}] \quad , \quad \mathcal{F} = \emptyset$$

**while** Stopping criterion is not met **do****Structure Learning:** Compute a (approximative) solution  $\hat{f}$  to Problem  $P_{\Theta, \mathbb{P}^D}$  and add the formula to  $\mathcal{F}$ , i.e.

$$\mathcal{F} \leftarrow \mathcal{F} \cup \{\hat{f}\}$$

Extend dimension of  $L$  by one, by  $f_p = \hat{f}$  and  $\theta[p] = 0$ **Weight Estimation:** Estimate the best weights for the added formula and recalibrate the weights of the previous formulas, by calling Algorithm 9.

$$\tilde{\mathbb{P}} \leftarrow \mathbb{P}^{\mathcal{F}, \theta}$$

**end while****return**  $\mathcal{F}, \theta$ **12.7 Proposal distribution**

Let us now understand the likelihood gradient as the energy tensor of a probability distribution, which we call the proposal distribution.

**Definition 61** (Proposal Distribution). *Let there be a base distribution  $\tilde{\mathbb{P}}$ , a targeted distribution  $\mathbb{P}^D$  and a formula selecting map  $\mathcal{H}[X_{[d]}, L]$ . The proposal distribution at inverse temperature  $\beta > 0$  is the distribution of  $L$  defined by*

$$\left\langle \exp \left[ \left\langle \beta \cdot (\mathbb{P}^D - \tilde{\mathbb{P}}), \mathcal{H} \right\rangle [L] \right] \right\rangle [L|\emptyset] .$$

The proposal distribution is the member of the exponential family with statistics  $\mathcal{H}$  and parameter  $\beta \cdot (\mathbb{P}^D - \tilde{\mathbb{P}})$ .

The proposal distribution is in the exponential family with sufficient statistic by the formula selecting map  $\mathcal{H}$ , namely the member with the canonical parameters  $\theta = \mathbb{P}^D - \tilde{\mathbb{P}}$ . Of further interest are tempered proposal distributions, which are in the same exponential family with canonical parameters  $\beta \cdot (\mathbb{P}^D - \tilde{\mathbb{P}})$  where  $\beta > 0$  is the inverse temperature parameter.

As Markov Logic Networks, the proposal distributions are in exponential families with the sufficient statistic defined in terms of formula selecting maps. While Markov Logic Networks contract the maps on the selection variables  $L$ , the proposal distributions contract them along the categorical variables  $X$  to define energy tensors.

The grafting Problem ( $P_{D, \mathcal{F}, \mathcal{H}}^{\text{grad}}$ ) is the search for the mode of the proposal distribution. To solve grafting, we thus need to answer a mode query, for which we can apply the methods introduced in Chapter 6, such as Gibbs Sampling or Mean Field Approximations in combination with annealing.

**12.7.1 Mean parameter polytope**

The mean parameter polytope of the proposal distribution with statistic  $\mathcal{H}^T$  is the convex hull of the formulas in  $\mathcal{F}$ , that is

$$\mathcal{M}_{\mathcal{H}^T} = \text{conv} \left( \sigma^{\mathcal{H}^T} L = l, X_{[d]} : l \in [p] \right) = \text{conv} \left( f[X_{[d]}] : f \in \mathcal{H} \right)$$

As it was the case for Markov Logic Networks, the mean parameter polytopes are instances of a 0/1-polytopes Ziegler (2000); Gillmann (2007).

The extreme points are the formulas selectable by the formula selecting map  $\mathcal{H}$ .

**12.8 Discussion**

**Remark 18** (Bayesian approach). *We only treated the estimation of a single resulting distribution by the data, while in a Bayesian approach one typically considers an uncertainty over possible distributions. When treating  $\theta$  as a*

*random tensor, which prior distribution is given and posteriori distribution wanted, we have a more involved Bayesian approach. When having a prior  $\mathbb{P}[\mathcal{F}, \theta]$  over the Markov Logic Networks we alternatively want to find the parameters  $\mathcal{F}, \theta$  solving the maximum a posteriori problem*

$$\operatorname{argmax}_{\mathcal{F}, \theta} \mathbb{P}^{\mathcal{F}, \theta}[\{D(j)\}_{j \in [m]}] \cdot \mathbb{P}[\mathcal{F}, \theta] . \quad (28)$$

To summarize some insights on the mean polytopes  $\mathcal{M}_{\mathcal{F}, \nu}$ :

- If and only if all coordinates are in  $\{0, 1\}$  then an extreme points, then  $\mu$  is reproduced by a hard logic network.
- If some mean params in  $\{0, 1\}$ , then not in the interior, and not reproduced by a markov logic network. Back direction not correct: There are interior points where no coordinate in  $\{0, 1\}$ .
- If not in the interior, we can identify with base measure refinement a base measure, such that reproducible by a distribution representable by the base measure.

The polytopes of mean parameters to hybrid logic networks and proposal distributions are an interesting connection between the fields of combinatorial optimization and the study of expressivity of tensor networks. This is of special interest, when the computation cores of a hybrid logic network are minimally connected, the mean parameters are captured by local consistencies. Similar investigations have been made in the field of tensor networks, where minimal connected tensor networks are referred to by Hierarchical Tucker formats (HT). Minimal connection is exploited in the tensor network community to show numerical properties of the format, such as closedness and existence of best approximators.

### 13 Probabilistic Guarantees

When drawing data independently from a random distribution, we are limited by random effects. We in this chapter derive guarantees, that the learning methods introduced in Chapter 6 and Chapter 12 are robust against such effects.

#### 13.1 Fluctuations of random data

A random tensor is a random element of a tensor space  $\bigotimes_{k \in [d]} \mathbb{R}^{m_k}$ , drawn from a probability distribution on  $\bigotimes_{k \in [d]} \mathbb{R}^{m_k}$ . In contrast to the discrete distributions investigated previously in this work, the random tensors are in most generality continuous distributions.

##### 13.1.1 Fluctuation of the empirical distribution

When drawing random states  $D(j) \in \times_{k \in [d]} [m_k]$  by a distribution  $\mathbb{P}^*$ , we use the one-hot encoding to forward each random state to the random tensor

$$\epsilon_{D(j)} [X_{[d]}] .$$

The expectation of this random tensor is

$$\mathbb{E} [\epsilon_{D(j)}] = \sum_{x_{[d]} \in \times_{k \in [d]} [m_k]} \mathbb{P}^* [X_{[d]} = x_{[d]}] \epsilon_{x_{[d]}} [X_{[d]}] = \mathbb{P}^* [X_{[d]}] .$$

The empirical distribution is then the average of independent random one-hot encodings, namely the random tensor

$$\mathbb{P}^D = \frac{1}{m} \sum_{j \in [m]} \epsilon_{D(j)} [X_{[d]}] .$$

To avoid confusion let us strengthen, that in this chapter we interpret  $\mathbb{P}^D$  as a random tensor taking values in  $\bigotimes_{k \in [d]} \mathbb{R}^{m_k}$ , whereas each supported value of  $\mathbb{P}^D$  is an empirical distribution taking values in  $\times_{k \in [d]} [m_k]$ . The forwarding of  $\times_{k \in [d]} [m_k]$  under the one-hot encoding is a multinomial random variable, see The. 86.

When the marginal of each datapoint is  $\mathbb{P}^*$ , the expectation of the empirical distribution is

$$\mathbb{E} [\mathbb{P}^D] = \frac{1}{m} \sum_{j \in [m]} \mathbb{E} [\epsilon_{D(j)}] = \mathbb{P}^* .$$

From the law of large numbers it follows, that in the limit of  $m \rightarrow \infty$  at any coordinate  $x \in \times_{k \in [d]} [m_k]$  almost everywhere

$$\mathbb{P}^D [X_{[d]} = x_{[d]}] \rightarrow \mathbb{E} [\mathbb{P}^D [X_{[d]} = x_{[d]}]] = \mathbb{P}^* [X_{[d]} = x_{[d]}] .$$

At finite  $m$  the empirical distribution differs from the by the difference

$$\mathbb{P}^D - \mathbb{P}^*$$

which we call a fluctuation tensor.

### 13.1.2 Mean parameter of the empirical distribution

We now investigate the empirical mean parameter

$$\mu_D [L] = \langle \sigma^S [X_{[d]}, L], \mathbb{P}^D [X_{[d]}] \rangle [L] .$$

Each coordinate of  $\mu_D$  is decomposed as

$$\mu_D [L = l] = \frac{1}{m} \sum_{j \in [m]} \mathcal{S}_l [D(j)]$$

and thus stores the empirical average of the feature  $\mathcal{S}_l$  on the dataset  $\{D(j)\}_{j \in [m]}$ .

Since the mean parameter depends linearly on the corresponding distribution, we can show the following correspondence between the empirical and the expected mean parameter.

**Theorem 80.** *When drawing data independently from  $\mathbb{P}^*$ , we have  $\mathbb{E} [\mu_D [L]] = \mu^* [L]$ , where we call*

$$\mu^* [L] = \langle \sigma^S [X_{[d]}, L], \mathbb{P}^D [X_{[d]}] \rangle [L]$$

*the expected mean parameter.*

*Proof.* Since the expectation commutes with linear functions. □

For each  $l \in [p]$  the law of large numbers guarantees that  $\mu^* [L = l]$  converges almost surely against  $\mu^* [L = l]$  when  $m \rightarrow \infty$ . To utilize these we need to approach the following issues:

- We need non-asymptotic convergence bounds, since one has access to finite data when learning
- The convergence has to happen uniformly for all  $l \in [p]$
- Guarantees on the result of an estimated model are more accessible when provided for quantities like the canonical parameter and KL-divergences of the learning result. Those, however, depend nonlinearly on  $\mu_D [L]$  and therefore require further investigation.

### 13.1.3 Noise tensor and its width

Motivated by The. 80, we build our derivation of probabilistic guarantees on non-asymptotic and uniform convergence bounds for  $\mu_D [L]$ . Let us first define the fluctuations of the empirical mean parameter, when drawing the data independently from a random distribution, as the noise tensor.

**Definition 62.** *Given a statistic  $\mathcal{S}$ ,  $m \in \mathbb{N}$  and a distribution  $\mathbb{P}^*$ , we call*

$$\eta^{S, \mathbb{P}^*, D} = \langle (\mathbb{P}^D - \mathbb{P}^*), \sigma^S \rangle [L]$$

*the noise tensor, where  $D$  is a collection of  $m$  independent samples of  $\mathbb{P}^*$ .*

The fluctuation of the empirical distribution around the generating distribution corresponds in this notation with the minterm exponential family, taking the identity as statistics. Besides this, fluctuation tensors appears in Markov Logic Networks as fluctuations of random mean parameters and in proposal distributions as fluctuation of random energy tensor. We will discuss these examples in the following sections.

We notice, that the fluctuation tensor  $\eta^{S, \mathbb{P}^*, D}$  is the centered mean parameter to the empirical distribution, that is

$$\mu_D - \mathbb{E} [\mu_D] = \langle \sigma^S, \mathbb{P}^D - \mathbb{P}^* \rangle [L] .$$

In the following we will use the supremum of contractions with random tensors in the derivation of success guarantees for learning problems. Such quantities are called widths.

**Definition 63.** Given a set  $\Gamma \subset \bigotimes_{k \in [d]} \mathbb{R}^{m_k}$  and  $\eta^{\mathcal{S}, \mathbb{P}^*, D}$  a random tensor taking values in  $\bigotimes_{k \in [d]} \mathbb{R}^{m_k}$  we define the width as the random variable

$$\Omega_\Gamma \left( \eta^{\mathcal{S}, \mathbb{P}^*, D} \right) = \sup_{\theta \in \Gamma} \left| \left\langle \theta, \eta^{\mathcal{S}, \mathbb{P}^*, D} \right\rangle [\emptyset] \right|.$$

Bounds on the widths are also called uniform concentration bounds Goeßmann (2021) and generic probabilistic bounds will be provided in Sect. 13.4.

### 13.2 Error bounds based on the noise width

We now derive error bounds for parameter estimation and structure learning, as introduced in Chapter 12. When combined with probabilistic bounds on the noise width, they are probabilistic success guarantees.

#### 13.2.1 Parameter Estimation

We in this section always assume, that  $\mathbb{P}^D$  is representable by the base measure  $\nu$  of the respective exponential families.

Parameter Estimation is the M-projection of the empirical distribution onto an exponential family. In Chapter 6 we have characterized those by the backward map acting on the mean parameter. Thus, while we are interested in the expected canonical parameter

$$\theta_* [L] = B^{(\mathcal{S}, \nu)}(\mu^* [L])$$

we get an estimation by the empirical canonical parameter

$$\theta_D [L] = B^{(\mathcal{S}, \nu)}(\mu_D [L]).$$

Unfortunately, since the backward map is not linear, we in general do not have that  $\mathbb{E} [B^{(\mathcal{S}, \nu)}(\mu_D)]$  coincides with  $B^{(\mathcal{S}, \nu)}(\mu^*)$ . To build intuition on the concentration we recall the expression of the backward map as

$$B^{(\mathcal{S}, \nu)}(\mu) = \operatorname{argmax}_\theta - \mathbb{H} [\mathbb{P}^\mu, \mathbb{P}^{(\mathcal{S}, \theta, \nu)}]$$

where  $\mathbb{P}^\mu$  is any distribution reproducing the mean parameter. We want to compare the solutions  $B^{(\mathcal{S}, \nu)}(\mu_D)$  and  $B^{(\mathcal{S}, \nu)}(\mu^*)$ , in which case  $\mathbb{P}^\mu$  can be chosen as  $\mathbb{P}^D$  and  $\mathbb{P}^*$ . It is common to call the objectives  $\mathbb{H} [\mathbb{P}^D, \mathbb{P}^{(\mathcal{S}, \theta, \nu)}]$  and  $\mathbb{H} [\mathbb{P}^*, \mathbb{P}^{(\mathcal{S}, \theta, \nu)}]$  empirical and expected risk Shalev-Schwartz, Shai and Ben-David, Shai (2014) Since the empirical risk has a linear dependence on  $\mu_D$ , we have at each  $\theta$

$$\begin{aligned} \mathbb{E} [\mathbb{H} [\mathbb{P}^D, \mathbb{P}^{(\mathcal{S}, \theta, \nu)}]] &= \mathbb{E} [\langle \mu_D, \theta \rangle [\emptyset] - A^{(\mathcal{S}, \nu)}(\theta)] \\ &= \langle \mathbb{E} [\mu_D], \theta \rangle [\emptyset] - A^{(\mathcal{S}, \nu)}(\theta) \\ &= \mathbb{H} [\mathbb{P}^*, \mathbb{P}^{(\mathcal{S}, \theta, \nu)}] \end{aligned}$$

By the law of large numbers, in the limit  $m \rightarrow \infty$  we thus have at each  $\theta$  a convergence of the empirical risk to the expected risk. However, since the backward map is defined by the minima of these risks, we need a uniform and non-asymptotical concentration guarantee to get more useful bounds. To this end, we now consider constrained parameter estimation and relate the supremum on the differences between expected and empirical risks with the width of the noise tensor.

**Lemma 20.** For any  $\Gamma$  and  $D$  we have

$$\Omega_\Gamma \left( \eta^{\mathcal{S}, \mathbb{P}^*, D} \right) = \sup_{\theta \in \Gamma} \left| \mathbb{H} [\mathbb{P}^D, \mathbb{P}^{(\mathcal{S}, \theta, \nu)}] - \mathbb{H} [\mathbb{P}^*, \mathbb{P}^{(\mathcal{S}, \theta, \nu)}] \right|.$$

*Proof.* For any  $\theta \in \Gamma$  and by  $\mathbb{P}^\mu$  realizable mean parameter  $\mu$  we have

$$\mathbb{H} [\mathbb{P}^\mu, \mathbb{P}^{(\mathcal{S}, \theta, \nu)}] = -\langle \mu, \theta \rangle [\emptyset] + A^{(\mathcal{S}, \nu)}(\theta).$$

It follows that

$$\mathbb{H} [\mathbb{P}^D, \mathbb{P}^{(\mathcal{S}, \theta, \nu)}] - \mathbb{H} [\mathbb{P}^*, \mathbb{P}^{(\mathcal{S}, \theta, \nu)}] = -\langle (\mu_D - \mu^*), \theta \rangle [\emptyset]$$

and the claim follows from comparison with Def. 62 and Def. 63.  $\square$

As a direct consequence, we have at any  $\theta \in \Gamma$

$$\left| \mathbb{H} \left[ \mathbb{P}^D, \mathbb{P}^{(S, \theta, \nu)} \right] - \mathbb{H} \left[ \mathbb{P}^*, \mathbb{P}^{(S, \theta, \nu)} \right] \right| \leq \Omega_\Gamma \left( \eta^{S, \mathbb{P}^*, D} \right).$$

Thus, the absolute difference of the expected risk and the empirical risk is bounded by the width of the noise tensor. This is especially useful for the solution  $\mu_D$  of the empirical risk minimization, where we can state

$$\mathbb{H} \left[ \mathbb{P}^*, \mathbb{P}^{(S, \theta_D, \nu)} \right] \leq \mathbb{H} \left[ \mathbb{P}^D, \mathbb{P}^{(S, \theta_D, \nu)} \right] + \Omega_\Gamma \left( \eta^{S, \mathbb{P}^*, D} \right).$$

At the solution of a empirical risk minimization problem over  $\Gamma$ , the expected risk exceeds the empirical risk at most by the noise tensor width.

When the generating distribution is in the hypothesis, we can further show the following KL-divergence bound for the estimated distribution.

**Theorem 81.** *Let us assume that for  $\theta_* \in \Gamma$  we have  $\mathbb{P}^* = \mathbb{P}^{(S, \theta_*, \nu)}$ . Then for any solution  $\theta_D$  of the empirical problem we have*

$$D_{\text{KL}} \left[ \mathbb{P}^{(S, \theta_*, \nu)} || \mathbb{P}^{(S, \theta_D, \nu)} \right] \leq 2\Omega_\Gamma \left( \eta^{S, \mathbb{P}^*, D} \right). \quad (29)$$

*Proof.* For the solution  $\theta_D$  of the empirical risk minimization on  $\Gamma$  we have since  $\theta_* \in \Gamma$  that

$$\mathbb{H} \left[ \mathbb{P}^D, \mathbb{P}^{(S, \theta_D, \nu)} \right] \leq \mathbb{H} \left[ \mathbb{P}^D, \mathbb{P}^{(S, \theta_*, \nu)} \right].$$

It follows that

$$\begin{aligned} D_{\text{KL}} \left[ \mathbb{P}^{(S, \theta_*, \nu)} || \mathbb{P}^{(S, \theta_D, \nu)} \right] &\leq D_{\text{KL}} \left[ \mathbb{P}^{(S, \theta_*, \nu)} || \mathbb{P}^{(S, \theta_D, \nu)} \right] + \mathbb{H} \left[ \mathbb{P}^D, \mathbb{P}^{(S, \theta_*, \nu)} \right] - \mathbb{H} \left[ \mathbb{P}^D, \mathbb{P}^{(S, \theta_D, \nu)} \right] \\ &= \left( \mathbb{H} \left[ \mathbb{P}^{(S, \theta_*, \nu)}, \mathbb{P}^{(S, \theta_D, \nu)} \right] - \mathbb{H} \left[ \mathbb{P}^D, \mathbb{P}^{(S, \theta_D, \nu)} \right] \right) \\ &\quad - \left( \mathbb{H} \left[ \mathbb{P}^{(S, \theta_*, \nu)}, \mathbb{P}^{(S, \theta_*, \nu)} \right] - \mathbb{H} \left[ \mathbb{P}^D, \mathbb{P}^{(S, \theta_*, \nu)} \right] \right), \end{aligned}$$

where we expanded the KL-divergence as a difference of cross entropies. We apply Lem. 20 to estimate the terms in brackets and get

$$\begin{aligned} &\left( \mathbb{H} \left[ \mathbb{P}^{(S, \theta_*, \nu)}, \mathbb{P}^{(S, \theta_D, \nu)} \right] - \mathbb{H} \left[ \mathbb{P}^D, \mathbb{P}^{(S, \theta_D, \nu)} \right] \right) - \left( \mathbb{H} \left[ \mathbb{P}^{(S, \theta_*, \nu)}, \mathbb{P}^{(S, \theta_*, \nu)} \right] - \mathbb{H} \left[ \mathbb{P}^D, \mathbb{P}^{(S, \theta_*, \nu)} \right] \right) \\ &\leq 2\Omega_\Gamma \left( \eta^{S, \mathbb{P}^*, D} \right). \end{aligned}$$

Combined with the above inequality we arrive at

$$D_{\text{KL}} \left[ \mathbb{P}^{(S, \theta_*, \nu)} || \mathbb{P}^{(S, \theta_D, \nu)} \right] \leq 2\Omega_\Gamma \left( \eta^{S, \mathbb{P}^*, D} \right) \square$$

One technical issue arises from the fact, that when we allow for  $\Gamma = \mathbb{R}^p$ , then  $\Omega_\Gamma \left( \eta^{S, \mathbb{P}^*, D} \right)$  vanishes or is infinity. To apply the result on the unconstrained parameter estimation, we therefore need to argue on bounded sets for the canonical parameter. When restricting to the sphere  $\mathbb{S} \subset \mathbb{R}^p$  we have

$$\|\mu_D - \mu^*\|_2 = \Omega_{\mathbb{S}} \left( \eta^{\mathcal{F}, \mathbb{P}^*, D} \right),$$

We apply this insight to state the following guarantee for unconstrained parameter estimation.

**Theorem 82.** *Let  $\theta$*

$$\left| \mathbb{H} \left[ \mathbb{P}^D, \mathbb{P}^{(S, \theta_D, \nu)} \right] - \mathbb{H} \left[ \mathbb{P}^*, \mathbb{P}^{(S, \theta_D, \nu)} \right] \right| \leq \Omega_{\mathbb{S}} \left( \eta^{\mathcal{F}, \mathbb{P}^*, D} \right) \cdot \|\theta_D\|_2.$$

*Proof.* As in the proof of Lem. 20 we use that

$$\mathbb{H} \left[ \mathbb{P}^D, \mathbb{P}^{(S, \theta_D, \nu)} \right] - \mathbb{H} \left[ \mathbb{P}^*, \mathbb{P}^{(S, \theta_D, \nu)} \right] = \langle \mu_D - \mu^*, \theta_D \rangle [\emptyset].$$

By Cauchy-Schwartz we further have

$$|\langle \mu_D - \mu^*, \theta_D \rangle [\emptyset]| \leq \|\mu_D - \mu^*\|_2 \cdot \|\theta_D\|_2.$$

Using that  $\|\mu_D - \mu^*\|_2 = \Omega_{\mathbb{S}} \left( \eta^{\mathcal{F}, \mathbb{P}^*, D} \right)$  we arrive at the claim.  $\square$

### 13.2.2 Structure Learning

In the gradient heuristic of structure learning, one selects the statistic to the maximal coordinate of the energy tensor of the proposal distribution. This tensor coincides with the mean parameter of a markov logic network and has thus a fluctuation by the noise tensor. We now use these insights to show a guarantee, that the formula chosen by grafting with respect to the empirical proposal distribution coincides with the formula chosen with respect to the expected proposal distribution. To this end, we need to define the max gap, which is the difference between the maximal coordinate of a tensor to the second maximal coordinate.

**Definition 64.** *The max gap of a tensor  $\tau [X_{[d]}]$  is the quantity*

$$\Delta(\tau) = \left( \max_{x_{[d]}} \tau [X_{[d]} = x_{[d]}] \right) - \left( \max_{x_{[d]} \notin \operatorname{argmax}_{x_{[d]}} \tau [X_{[d]} = x_{[d]}]} \tau [X_{[d]} = x_{[d]}] \right).$$

When comparing the gap with the noise width, we get the following guarantee.

**Theorem 83.** *Whenever*

$$\Delta(\mu^*) > 2 \cdot \Omega_{\{\epsilon_{x_{[d]}} : x_{[d]} \in \times_{k \in [d]} [m_k]\}} \left( \eta^{S, \mathbb{P}^*, D} \right),$$

*then any mode  $x_{[d]}$  of the empirical proposal distribution is a mode of the expected proposal distribution.*

*Proof.* Let us assume that for a mode  $l^D \in \operatorname{argmax}_{l \in [p]} \mu_D [L = l]$  of the empirical mean parameter we have

$$l^D \notin \operatorname{argmax}_{l \in [p]} \mu^* [L = l].$$

For a mode  $l^* \in \operatorname{argmax}_{l \in [p]} \mu^* [L = l]$  of the expected mean parameter we then have

$$\mu^* [L = l^D] \leq \mu^* [L = l^*] - \Delta(\mu^*)$$

and

$$\mu_D [L = l^D] \geq \mu_D [L = l^*].$$

Comparing both inequalities we get

$$(\mu_D [L = l^D] - \mu^* [L = l^D]) + (-\mu_D [L = l^*] + \mu^* [L = l^*]) \geq \Delta(\mu^*).$$

Estimating the terms in the bracket by the width of the noise tensor with respect to basis vectors, we get

$$2 \cdot \Omega_{\{\epsilon_{x_{[d]}} : x_{[d]} \in \times_{k \in [d]} [m_k]\}} \left( \eta^{S, \mathbb{P}^*, D} \right) \geq \Delta(\mu^*),$$

which is a contradiction to the assumption. Thus, any mode of the empirical mean parameter is also a model of the expected mean parameter.  $\square$

### 13.2.3 Mode recovery

Let us now consider a more general problem than in the section above, namely the estimation of the modes of a distribution. Let  $\hat{\theta}$  be the estimator of the canonical parameter  $\theta_*$ , then the mode set of both coincide, if and only if they are elements in the same max cone, i.e.

$$C_{S, \nu}^{\hat{\theta}} = C_{S, \nu}^{\theta_*}.$$

To ensure, that this is the case, we generalize the gap at  $\theta_*$  as the minimal distance to other cones and bound uniform concentration events implying that the distance between  $\hat{\theta}$  and  $\theta_*$  is smaller than the gap.

**Definition 65.** *Let  $d(\cdot, \cdot)$  be a metric on  $\mathbb{R}^p$ , then the generalized gap of  $\theta \in \mathbb{R}^p$  is defined as*

$$\Delta_d(\theta) = \inf_{\tilde{\theta} \notin C_{S, \nu}^{\theta}} d(\tilde{\theta}, \theta).$$

Atomic norms induce metrics, which are widths (see ? and Chapter 5 in Goeßmann (2021)).



**Definition 66.** Given a set  $\Gamma \subset \mathbb{R}^p$  such that an open neighborhood of the origin is contained in  $\text{conv}(\Gamma)$ . Then

$$\|\theta\|_\Gamma = \Omega_\Gamma(\theta)$$

is the dual atomic norm and

$$d_\Gamma(\theta, \hat{\theta}) = \Omega_\Gamma(\theta - \hat{\theta})$$

is the dual atomic distance to  $\Gamma$ .

Examples of atomic norms are:

- Euclidean distance  $\ell_2$ , when  $\Gamma = \mathbb{S}$
- Supremum distance  $\ell_\infty$ , when  $\Gamma = \{\lambda \cdot \epsilon_l[L] : \lambda \in \{-1, +1\}, l \in [p]\}$

**Theorem 84.** Let  $\Gamma \subset \mathbb{R}^p$  induce an atomic norm. If

$$\Omega_\Gamma(\hat{\theta} - \theta_*) < \Delta_{d_\Gamma}(\theta_*)$$

then the modes of  $\mathbb{P}^{(S, \hat{\theta}, \nu)}$  and  $\mathbb{P}^{(S, \theta_*, \nu)}$  coincide.

*Proof.* If  $\hat{\theta} \notin C_{S, \nu}^{\theta_*}$ , then

$$\Delta_{d_\Gamma}(\theta_*) \geq d_\Gamma(\hat{\theta}, \theta_*) = \Omega_\Gamma(\hat{\theta} - \theta_*) ,$$

which contradicts the assumption. Therefore if the assumption holds, then  $\hat{\theta} \in C_{S, \nu}^{\theta_*}$  and the modes of  $\mathbb{P}^{(S, \hat{\theta}, \nu)}$  and  $\mathbb{P}^{(S, \theta_*, \nu)}$  coincide.  $\square$

The guarantee on structure learning is the special case, where  $\Gamma = \{\lambda \cdot \epsilon_l[L] : \lambda \in \{-1, +1\}, l \in [p]\}$  and

$$\Delta_{d_\Gamma}(\theta) = \frac{1}{2} \max_{l \notin \arg\max_l \theta[L=l]} \left| \theta[L=l] - \max_l \theta[L=l] \right| .$$

### 13.3 Fluctuations in Logic Networks

In case of logical formulas being statistics, the coordinates of the mean parameter are satisfaction rates to the formulas.

For Logic Networks we have statistics consistent of boolean statistics  $f_l$ , which are logical formulas. In this case the marginal distributions of the coordinates of  $\eta^{S, \mathbb{P}^*, D}$  are scaled and centered binomials, which we show now.

**Theorem 85.** For any  $\mathcal{F}$  the marginal distribution of the coordinate  $\eta^{\mathcal{F}, \mathbb{P}^*, D}[L=l]$  is the scaled and centered binomial distribution

$$\frac{1}{m} (B(m, \mu[L=l]) - \mu[L=l])$$

with parameters  $m$  and  $\mu[L=l]$ .

*Proof.* We notice that when forwarding a random sample  $D(j)$  of  $\mathbb{P}^*$  is the random tensor

$$\epsilon_{D(j)}[X_{[d]}]$$

and since  $\text{im}(\mathcal{S}_l) \subset \{0, 1\}$  the contraction

$$\langle \mathcal{S}_l, \epsilon_{D(j)}[X_{[d]}] \rangle [\emptyset]$$

is a random variable taking values in  $\{0, 1\}$ . The variable therefore follows a Bernoulli distribution with mean parameter

$$\mu[L=l] = \mathbb{E}[\langle \mathcal{S}_l, \epsilon_{D(j)}[X_{[d]}] \rangle [\emptyset]] = \langle \mathcal{S}_l, \mathbb{P}^* \rangle [\emptyset] \quad \square$$

The mean parameter of the M-projection of the empirical distribution on the family of Markov Logic Networks with statistic  $\mathcal{H}$  is the random tensor

$$\mu_D[L] = \langle \sigma^{\mathcal{F}}, \mathbb{P}^D \rangle [L] .$$

The expectation of this random tensor is

$$\mathbb{E}[\mu_D] = \langle \sigma^{\mathcal{F}}, \mathbb{E}[\mathbb{P}^D] \rangle [L] = \langle \sigma^{\mathcal{F}}, \mathbb{P}^* \rangle [L] = \mu^* ,$$

where we used that the expectation and contraction operation can be commuted due to the multilinearity of contractions.

### 13.3.1 Energy tensor in proposal distributions

The fluctuation tensor appears as a fluctuation of the energy of the proposal distribution. The expectation of the energy of the proposal distribution is

$$\begin{aligned}\mathbb{E} \left[ \phi^{\mathcal{H}^T, \mathbb{P}^D - \tilde{\mathbb{P}}} \right] &= \mathbb{E} \left[ \left\langle \sigma^{\mathcal{H}^T}, \mathbb{P}^D - \tilde{\mathbb{P}} \right\rangle [L] \right] = \left\langle \sigma^{\mathcal{H}^T}, \mathbb{E} \left[ \mathbb{P}^D - \tilde{\mathbb{P}} \right] \right\rangle [L] = \left\langle \sigma^{\mathcal{H}^T}, \mathbb{P}^* - \tilde{\mathbb{P}} \right\rangle [L] \\ &= \mathbb{E} \left[ \phi^{\mathcal{H}^T, \mathbb{P}^* - \tilde{\mathbb{P}}} \right].\end{aligned}$$

The fluctuation of this random tensor is

$$\mathbb{E} \left[ \phi^{\mathcal{H}^T, \mathbb{P}^D - \tilde{\mathbb{P}}} \right] - \mathbb{E} \left[ \phi^{\mathcal{H}^T, \mathbb{P}^* - \tilde{\mathbb{P}}} \right] = \mathbb{E} \left[ \phi^{\mathcal{H}^T, \mathbb{P}^D - \mathbb{P}^*} \right]$$

and coincides with  $\eta^{\mathcal{F}, \mathbb{P}^*, D}$ .

### 13.3.2 Minterm Exponential Family

In case of the minterm exponential family, we have  $\mathcal{S} = \delta [X_{[d]}, L]$  and the noise tensor is

$$\eta^{\delta, \mathbb{P}^*, D} = \mathbb{P}^D - \mathbb{P}^*.$$

This noise tensor follows a multinomial distribution as we show next. To this end, we notice that a multinomial distribution can be defined as the average of one-hot encodings of independently and identically distributed datapoints. When drawing  $\{D(j)\}_{j \in [m]}$  independently from  $\mathbb{P}^*$  we denote

$$\sum_{j \in [m]} \epsilon_{D(j)} [X_{[d]}] \sim \underline{B}(m, \mathbb{P}^*).$$

**Theorem 86.** *The noise tensor  $\eta^{\delta, \mathbb{P}^*, D}$  is a by  $\frac{1}{m}$  rescaled centered multinomial random tensor with parameters  $\mathbb{P}^*$  and  $m$ , that is*

$$\eta^{\delta, \mathbb{P}^*, D} \sim \frac{1}{m} (\underline{B}(m, \mathbb{P}^*) - \mathbb{P}^*).$$

*Proof.* By the above construction we have

$$\mathbb{P}^D - \mathbb{P}^* = \frac{1}{m} \sum_{j \in [m]} (\epsilon_{D(j)} [X_{[d]}] - \mathbb{E} [\epsilon_{D(j)} [X_{[d]}]])$$

We further have

$$\mathbb{E} [\epsilon_{D(j)} [X_{[d]}]] = \mathbb{P}^* [X_{[d]}].$$

□

The noise tensor characterization by multinomial distributions, which holds for minterm statistics, is a more detailed characterization compared to the characterization of its marginals by binomial distribution in The. 85, which holds for generic statistics  $\mathcal{F}$ .

### 13.3.3 Guarantees for Mode of the Proposal Distribution

Let us now derive probabilistic guarantees, that the mode of the proposal distribution at the empirical and the generating distribution are equal.

**Theorem 87.** *Whenever the energy tensor of the expected proposal distribution has a gap of  $\Delta$ , then for every  $p > 0$  any mode of the empirical proposal distribution coincides is also a mode of the expected proposal distribution with probability at least  $1 - \exp \left[ -\frac{1}{p^2} \right]$ , provided that*

$$m > C \frac{(1 + \ln [p])}{\Delta^2}$$

where  $C$  is a universal constant.

*Proof.* To proof the theorem we combine the deterministic guarantee The. 83 with the width bound of The. 89, which we show in the next section. Given the assumed bound, the sub-gaussian norm of the width is upper bounded by  $C_2 \cdot \Delta$ , thus for any  $p > 0$  we have

$$\Omega_{\{\epsilon_l[L] : l \in [p]\}} \left( \eta^{\mathcal{F}, \mathbb{P}^*, D} \right) < 2\Delta$$

with probability at least  $1 - \exp \left[ -\frac{1}{p^2} \right]$ . The claim thus follows with The. 83.  $\square$

**Example 16** (Gap of a MLNs with single formulas). *Let there be the MLN of a maxterm  $f$  with  $d$  variables, and let  $\mathcal{F}$  be the maxterm selecting tensor, then*

$$\Delta \left( \phi^{\langle \mathcal{F}, \mathbb{P}^{\langle \{f\}, \theta \rangle} - \langle \mathbb{I} \rangle [X_{[d]} | \emptyset] \rangle} \right) = \frac{1}{2^d - 1 + \exp[-\theta]}$$

*If  $\theta > 0$  we have an exponentially small gap. Thus, for the above Lemma to apply, the width needs to be exponentially in  $d$  small.*

*Let there be the MLN of a minterm  $f$  with  $d$  variables, then*

$$\Delta \left( \phi^{\langle \mathcal{F}, \mathbb{P}^{\langle \{f\}, \theta \rangle} - \langle \mathbb{I} \rangle [X_{[d]} | \emptyset] \rangle} \right) = \frac{1}{1 + (2^d - 1) \cdot \exp[-\theta]}$$

*For large  $\theta$  and  $d$ , the gap tends to 1.*

### 13.3.4 Guarantees for Unconstrained Parameter Estimation

We here the sphere bounds and combine with The. 82.

**Theorem 88.** *For any  $p \in (0, 1)$  we have the following with probability at least  $1 - p$ . Let  $\hat{\theta}$  and  $t > 0$ , then*

$$\left| \mathbb{H} \left[ \mathbb{P}^*, \mathbb{P}^{\langle \mathcal{F}, \theta_D, \nu \rangle} \right] - \mathbb{H} \left[ \mathbb{P}^D, \mathbb{P}^{\langle \mathcal{F}, \theta_D, \nu \rangle} \right] \right| \leq \tau \cdot \|\theta_D\|_2$$

*provided that*

$$m \geq \frac{\langle \mu^* \rangle [\emptyset] - \langle (\mu^*)^2 \rangle [\emptyset]}{pt^2}.$$

*Proof.* The claim follows from the deterministic guarantee The. 82 with the probabilistic width bound The. 90 to be shown in the next section.  $\square$

## 13.4 Width bounds for the noise tensor

We here provide width bounds on the noise tensors  $\eta^{\mathcal{F}, \mathbb{P}^*, D}$  to logic networks, which coordinates have marginal distributions by Binomials, as shown in The. 85. All bounds hold for arbitrary statistics  $\mathcal{F}$  of propositional formulas and number  $m$  of data and the appearing constants are universal, that is independent of particular choices of  $\mathcal{F}$  and  $m$ .

### 13.4.1 Basis Vectors

We first introduce the sub-Gaussian Norm and show how we can exploit it to state concentration inequalities.

**Definition 67** (Sub-Gaussian Norm, see Def. 2.5.6 in Vershynin (2018)). *The sub-Gaussian norm of a random variable  $X$  is defined as*

$$\|X\|_{\psi_2} = \inf \left\{ C > 0 : \mathbb{E} \left[ \exp \left[ \frac{X^2}{C^2} \right] \right] \leq 2 \right\}.$$

The moment bound used to define the sub-Gaussian norm can then be combined with Markov's inequality to state concentration bounds. Before showing the utility of these norm, let us first connect with the contraction formalism of this work. When  $X$  is a random coordinate of  $\tau [X_{[d]}]$ , selected by a probability tensor  $\mathbb{P} [X_{[d]}]$  we have

$$\mathbb{E} \left[ \exp \left[ \frac{X^2}{C^2} \right] \right] = \left\langle \mathbb{P} [X_{[d]}], \exp \left[ \frac{1}{C^2} \cdot \tau [X_{[d]}] \right] \right\rangle [\emptyset]$$

and thus

$$\|X\|_{\psi_2} = \inf \left\{ C > 0 : \left\langle \mathbb{P} [X_{[d]}], \exp \left[ \frac{1}{C^2} \cdot \tau [X_{[d]}] \right] \right\rangle [\emptyset] \leq 2 \right\}.$$

We now show a sub-Gaussian norm bound on the coordinates of the noise tensor.

**Lemma 21.** *The marginal distribution of any coordinate of  $\eta^{\mathcal{F}, \mathbb{P}^*, D} [L]$  is sub-Gaussian with*

$$\left\| \eta^{\mathcal{F}, \mathbb{P}^*, D} [L = l] \right\|_{\psi_2} \leq C_0 \frac{1}{\sqrt{m}}$$

where  $C_0 > 0$  is a universal constant.

*Proof.* Any centered Bernoulli variable is bounded and therefore sub-Gaussian with

$$\left\| \langle f_l [X_{[d]}], \epsilon_{D(j)} [X_{[d]}] \rangle [\emptyset] - \langle f_l [X_{[d]}], \mathbb{P}^* \rangle [\emptyset] \right\|_{\psi_2} \leq \frac{1}{\sqrt{\ln [2]}}.$$

Binomial variables are sums of independent Bernoulli variables. We apply the sub-Gaussian norm bound for sums from Proposition 2.6.1 in Vershynin (2018), which states that for a universal constant  $C > 0$  we have

$$\left\| \left\langle f_l [X_{[d]}], \left( \sum_{j \in [m]} \epsilon_{D(j)} [X_{[d]}] - \mathbb{P}^* \right) \right\rangle [\emptyset] \right\|_{\psi_2} \leq \frac{C \cdot \sqrt{m}}{\sqrt{\ln [2]}}.$$

We therefore have

$$\left\| \eta^{\mathcal{F}, \mathbb{P}^*, D} [L = l] \right\|_{\psi_2} = \frac{1}{m} \left\| \left\langle f_l [X_{[d]}], \left( \sum_{j \in [m]} \epsilon_{D(j)} [X_{[d]}] - \mathbb{P}^* \right) \right\rangle [\emptyset] \right\|_{\psi_2} \leq \frac{C}{\sqrt{\ln [2]} \cdot m}.$$

We arrive at the claimed bound with a transform of the universal constant to  $C_0 = \frac{C}{\sqrt{\ln [2]}}$ .  $\square$

Based on this norm bound, we now show a bound on the sub-Gaussian norm of the width with respect to basis vectors.

**Theorem 89.** *For the set of basis vectors*

$$\Gamma = \{\epsilon_l [L] : l \in [p]\}$$

we have

$$\left\| \Omega_{\Gamma} \left( \eta^{\mathcal{F}, \mathbb{P}^*, D} \right) \right\|_{\psi_2} \leq C_1 \sqrt{\frac{1 + \ln [p]}{m}},$$

where  $C_1 > 0$  is a universal constant.

*Proof.* We first notice, that

$$\Omega_{\Gamma} (\eta) = \max_{l \in [p]} \left| \eta^{\mathcal{F}, \mathbb{P}^*, D} [L = l] \right|$$

By a generic bound on the supremum of sub-Gaussian variables (see Exercise 2.5.10 in Vershynin (2018)) we have for a universal constant  $C > 0$

$$\left\| \max_{l \in [p]} \left| \eta^{\mathcal{F}, \mathbb{P}^*, D} [L = l] \right| \right\|_{\psi_2} \leq C \left( \max_{l \in [p]} \left\| \eta^{\mathcal{F}, \mathbb{P}^*, D} [L = l] \right\|_{\psi_2} \right) \sqrt{1 + \ln [p]}.$$

We now apply Lem. 21 and get with  $C_1 = C \cdot C_0$  that

$$\left\| \Omega_{\Gamma} \left( \eta^{\mathcal{F}, \mathbb{P}^*, D} \right) \right\|_{\psi_2} \leq C_1 \sqrt{\frac{1 + \ln [p]}{m}}.$$

$\square$

The bound in The. 89 is furthermore sharp, see the construction of an identically scaling lower bound in Exercise 2.5.11 in Vershynin (2018). Note that the binomials used here tend to normal distributed variables used in the construction therein.

### 13.4.2 Sphere

For any tensor  $\eta [L]$  and the sphere  $\mathbb{S} \subset \mathbb{R}^p$  we have

$$\Omega_{\mathbb{S}}(\eta [L]) = \|\eta [L]\|_2 .$$

To show probabilistic width bounds with respect to the sphere, we therefore apply in the following Chebyshevs inequality on the norm of random tensors.

**Theorem 90.** *Let  $\mu [L]$  be a deterministic vector with coordinates in  $[0, 1]$  and  $\eta [L]$  a random vector, which coordiantes are for  $l \in [p]$  marginally distributed as*

$$\eta [L = l] \sim B(m, \mu [L = l]) .$$

Then we have for any  $p > 0$ ,  $t > 0$  and  $m \in \mathbb{N}$  with probability at least  $1 - p$

$$\left\| \frac{\eta - \mathbb{E}[\eta]}{m} \right\|_2 \leq t$$

provided that

$$m \geq \frac{\langle \mu [L], (\mathbb{I} [L] - \mu [L]) \rangle [\emptyset]}{p \cdot t^2} .$$

*Proof.* Since the squared norm of the noise is the sum of squared centered and averaged Binomials, we have

$$\mathbb{E} \left[ \|\eta [L] - \mathbb{E}[\eta [L]]\|_2^2 \right] = m \cdot \left( \sum_{l \in [p]} \mu [L = l] (1 - \mu [L = l]) \right)$$

Here we used that the variance of a variable distributed by  $B(m, \mu [L = l])$  is  $m \cdot \mu [L = l] (1 - \mu [L = l])$ .

It follows, that

$$\mathbb{E} \left[ \left( \left\| \frac{\eta - \mathbb{E}[\eta]}{m} \right\|_2 \right)^2 \right] = \frac{\langle \mu [L], (\mathbb{I} [L] - \mu [L]) \rangle [\emptyset]}{m} .$$

Then we apply a Chebyshev Bound to get for any  $t > 0$

$$\mathbb{P} \left[ \left\| \frac{\eta - \mathbb{E}[\eta]}{m} \right\|_2 > t \right] = \mathbb{P} \left[ \left( \left\| \frac{\eta - \mathbb{E}[\eta]}{m} \right\|_2 \right)^2 > t^2 \right] \leq \frac{\langle \mu [L], (\mathbb{I} [L] - \mu [L]) \rangle [\emptyset]}{m \cdot t^2} \quad (30)$$

For a  $p > 0$  we choose any  $m$  with

$$m \geq \frac{\langle \mu [L], (\mathbb{I} [L] - \mu [L]) \rangle [\emptyset]}{t^2 p}$$

and get

$$\mathbb{P} \left[ \left\| \frac{\eta - \mathbb{E}[\eta]}{m} \right\|_2 > t \right] \leq p . \quad (31)$$

Thus, we have

$$\mathbb{P} \left[ \left\| \frac{\eta - \mathbb{E}[\eta]}{m} \right\|_2 \leq t \right] = 1 - \mathbb{P} \left[ \left\| \frac{\eta - \mathbb{E}[\eta]}{m} \right\|_2 > t \right] \geq 1 - p . \quad (32)$$

□

For the minterm family where  $\mathcal{F} = \delta$  the noise tensor is a rescaled and centered multinomial. In that case, the bound of The. 90 can be simplified by

$$\langle \mu [L], (\mathbb{I} [L] - \mu [L]) \rangle [\emptyset] = 1 - \langle \mu [L]^2 \rangle [\emptyset] .$$

### 13.5 Discussion

We in this chapter only provided probabilistic width bounds for logic networks, that are exponential families with boolean statistics. Similar recovery bounds for parameter estimation and structure learning for more general exponential families would require width bounds in these generic cases. A general approach towards width bounds are chaining techniques on stochastic processes, see Talagrand (2014). While we showed bounds based on the sub-Gaussian norm, more general sub-exponential bounds could be used, see Wainwright (2019).

We further assumed that our random tensors to be projected are empirical distributions. More general random tensor networks and corresponding width bounds have been developed in Goeßmann (2021).

## 14 First Order Logic

We now extend the tensor representation from to structured representations, whereas we previously focused on factored representation of systems.

**We observe that the more expressive first-order logic bears another tensor structure: The representation of each world is a boolean tensor.**

### 14.1 World Tensors

Since first-order logic follows structured representations of a system, a first-order logic world consists in objects and relations between them. To each world there is a world domain  $\mathcal{U}$  of objects, which we assume to be finite (this is a restrictive assumption). We exploit the set-encoding formalism discussed in more detail in Chapter 17 and use bijective index interpretation maps

$$I : [r] \rightarrow \mathcal{U}.$$

A so-called term variable  $O$  takes states  $o \in [r]$ , which represent objects

$$I(o) \in \mathcal{U}.$$

The relations between objects are described by  $n$ -ary predicates  $g$ . Given a specific world  $x_W$  the truth of relations is represented by boolean tensors

$$g|_{x_W} : \bigtimes_{l \in [n]} [r] \rightarrow \{0, 1\}.$$

Given a tuple  $o_0, \dots, o_{n-1} \in \bigtimes_{l \in [n]} [r]$  the boolean

$$g|_{x_W} [O_0 = o_0, \dots, O_{n-1} = o_{n-1}] \in \{0, 1\}$$

is called a grounding and encodes, whether the relation  $g$  is satisfied in the world  $x_W$  for the objects  $I^{-1}(o_0), \dots, I^{-1}(o_{n-1})$ .

Let us assume, that we have a function-free theory with  $d$  predicates, where are predicates all of the same arity  $n$ . We then formalize a world in the following based on a selection variable  $L$  selecting a specific predicate and term variables  $O_{[n]} = O_0, \dots, O_{n-1}$  representing choices of objects from a given set  $\mathcal{U}$ .

**Definition 68 (FOL World).** *Given a set of objects  $\mathcal{U}$  enumerated by an index interpretation function  $I : [r] \rightarrow \mathcal{U}$  and a finite set  $\{g_0, \dots, g_{d-1}\}$  of  $n$ -ary predicates a world is a boolean tensor*

$$x_W[L, O_{[n]}] : [d] \times \left( \bigtimes_{l \in [n]} [r] \right) \rightarrow [2]. \quad (33)$$

*We interpret the world tensor as encoding in the coordinate  $x_W[L = k, O_{[n]} = o_{[n]}]$ , whether the  $k$ -th predicate is satisfied on the object tuple  $I^{-1}(o_0), \dots, I^{-1}(o_{n-1})$ .*

When the assumptions of function-free and constant variable order are not met, we can do the following tricks. Functions are turned to predicates by their relation interpretation. If there are predicates of different arity in the theory, we can trivially extend them to  $n$ -ary predicates by tensor products with the trivial tensor  $\mathbb{I}$ . This can be done by a tensor product with  $\epsilon_r [O]$ , where we add an auxiliary object  $I_r$  as a placeholder for predicates with smaller arity.

While in first order logics, depending on the chosen semantics, worlds can have infinite sets of objects, we here only treat worlds with finite objects.

### 14.1.1 Case of Propositional Logics

Before continuing with the one-hot encoding of first-order logic worlds, let us show that the previously discussed formalism of propositional logics (see Chapter 7) is a special case of first-order logics, namely when demanding  $n = 0$ . Consistent with Def. 68 we have a propositional logic world by

$$x_W : [d] \rightarrow [2],$$

which we have in Chapter 7 represented by the assignments  $x_k = x_W[L = k]$  to the categorical variables  $X_k$ .

To represent logical formulas as sets of possible worlds, and distributions of worlds, we applied in Part I one-hot encodings of possible worlds. For the case of propositional logics, this is

$$\epsilon_{x_W} [X_{[d]}] = \bigotimes_{k \in [d]} \epsilon_{x_W[L=k]} [X_k].$$

### 14.1.2 One-hot encoding of worlds

Let us now generalize the one-hot encodings of propositional logic worlds to worlds of first-order logic. To encode the boolean tensors  $x_W$  describing first order logics as basis elements of a tensor space, we take the one-hot encoding

$$\epsilon : \bigotimes_{k \in [d]} \bigotimes_{o_0 \in [r_0]} \cdots \bigotimes_{o_{n-1} \in [r_{n-1}]} [2] \rightarrow \bigotimes_{k \in [d]} \bigotimes_{o_0 \in [r_0]} \cdots \bigotimes_{o_{n-1} \in [r_{n-1}]} \mathbb{R}^2$$

defined by

$$\epsilon_{x_W} [X_{[d] \times [r]^n}] = \bigotimes_{k \in [d]} \bigotimes_{o_0 \in [r_0]} \cdots \bigotimes_{o_{n-1} \in [r_{n-1}]} \epsilon_{x_W[L=k, O_{[n]}=o_{[n]}]} [X_{k, o_{[n]}}].$$

This is a tensor of order  $d \cdot r^n$ , in a tensor space of dimension  $2^{(d \cdot r^n)}$ . Storage of such tensors in naive formats would not be possible. However, the basis CP format discussed in Chapter 18 still provides storage with demand linear in the order  $d \cdot r^n$ .

Another issue when comparing different first-order logic worlds arises in potentially different world domains. As we have explored, the cardinality of the domain influences the order of the one-hot encoding tensors. To avoid such issues we here enumerate worlds coinciding in their domains. This restriction is called database semantics (see e.g. Section 8.2.8 in Russell and Norvig (2021)), where only those worlds are considered, which domains have a one-to-one map to the constant symbols appearing in a respective knowledge base. When restricting to worlds coinciding in their domain, we still have a factored representation of the system, since we can enumerate the possible worlds by a cartesian product. However, the number of categorical variables representing the world is  $d \cdot r^n$  and tensor representations, even in sparse formats, are not feasible due to the large order required. These techniques to restrict to comparable factored representations are often referred to propositionalization of a first-order logic knowledge base.

### 14.1.3 Probability distributions

Having established the formalism of one-hot encodings also in the case of first-order logic worlds, we can now proceed with the definition of distributions and formulas, analogously to the development in Part I. Probability distributions over worlds coinciding on their domain are then non-negative and normed tensors

$$\mathbb{P} [X_{[d] \times [r]^n}] \in \bigotimes_{k \in [d], o_{[n]} \in [r]^n} \mathbb{R}^2.$$

where each coordinate of a world  $x_W$  is captured by a boolean random variable  $X_{k, o_{[n]}}$ , indicating whether the  $k$ -th predicate holds on the object tuple indexed by  $o_{[n]}$ .

We notice, that by definition these probability distributions are distributions of  $d \cdot r^n$  Booleans with  $2^{(d \cdot r^n)}$  many states. Unfortunately, it is not possible to design encoding spaces of smaller dimension, when our aim is to get any distribution over possible worlds by an element in the encoding space. This is due to the fact, that one-hot encodings provide a basis in the tensor space, as will be shown in Chapter 16. The reason for the large encoding space dimension is therefore rooted in the equal number of possible worlds and not in an overhead in the dimension of the one-hot encoding space. We will later in this chapter investigate methods to handle such high-dimensional distributions in the formalism of exponential families.

#### 14.1.4 Semantics of formulas

Following the development of Chapter 7, we can choose a semantic approach to the definition of formulas, under the assumption of database semantics. Since the semantic of a logical formula is the set of its models, we again have a one-to-one correspondence between logical formulas and the boolean tensors in the one-hot encoding space

$$\bigotimes_{k \in [d], o_{[n]} \in [r]^n} \mathbb{R}^2.$$

This correspondence between the semantics and boolean tensor is through a subset encoding (see Def. 80) of the respective formulas. However, due to the large state dimensions, we will in the following sections choose a syntactical approach to the construction of formulas, which will naturally provide efficient tensor network decompositions.

#### 14.1.5 Two levels of tensor representation

In comparison with propositional logics, first-order logic bears two levels of natural tensor representations. In the first level, which we call the structured level, each world (see Def. 68) has a natural structure by a tensor, since it encodes relations between objects chosen by assignments to term variables. This is different to the worlds of a propositional logic theory, which are represented by a boolean vector instead of a tensor. The second level arises as in propositional logics, by understanding each world as a uncertain state and studying distributions over states, which are understood themselves as a tensor (see Def. 16). We call this the factored level, since it arises in general in the discussion of factored representations. As argued above, the assumption of database semantics is central to exploit the tensor structure of the substitution level. Under this assumption, representation of an uncertain state, or a collection of possible states, is done in the tensor space

$$\bigotimes_{k \in [d], o_{[n]} \in [r]^n} \mathbb{R}^2$$

where the enumeration of the 2-dimensional axes contains the tensor structure of the substitution level.

### 14.2 Formulas in a fixed first-order logic world

Following the argumentation above, we in this section restrict to the exploitation of tensors in the structured level, namely a fixed world represented as a tensor  $x_W[L, O_{[n]}]$ , see Def. 68. We are specifically interested in the tensor network decomposition of first order formulas, which contain in full generality variables and therefore also have a tensor. The evaluation of a first-order formula on a specific world is therefore different to the case in propositional logics, where the evaluation was a boolean in  $\{0, 1\}$  indicating whether the world is a model.

#### 14.2.1 Grounding tensors

Given a first-order logic world  $x_W[L, O_{[n]}]$ , arbitrary formulas are interpreted in terms of the satisfactions of their groundings. We define their semantic first, and then relate their syntactical decomposition to tensor networks, similar to our approach to propositional logics in Chapter 7.

**Definition 69** (Grounding of a first-order formula given a world). *Given a specific world  $x_W$ , with an domain  $\mathcal{U}$  enumerated by  $[r]$ , the grounding of a formula  $q$  with variables  $O_q$  is the tensor*

$$q|_{x_W}[O_q] : \bigtimes_{l \in [O_q]} [r] \rightarrow \{0, 1\}.$$

*Each coordinate represents thereby the boolean, whether the substitution of the variables in the formula is satisfied given a world  $x_W$ , that is*

$$q|_{x_W}[O_q = o_q] = 1$$

*if and only if the substitution of  $q$  with the variables  $O_q$  replaced by the objects  $I(o_l)$  is satisfied on the world  $x_W$ .*

The grounding tensor formalism can be used to define formulas as a map

$$q : \left( \bigotimes_{k \in [d], o_{[n]} \in [r]^n} \mathbb{R}^2 \right) \rightarrow \left( \bigotimes_{k \in [d], o_q \in [r]^{|O_q|}} \mathbb{R}^2 \right)$$

where each world  $x_W$  is mapped to a grounding tensor

$$q(x_W) = q|_{x_W}.$$

This would involve the factored level of tensor interpretation, namely representation of all possible worlds.



### 14.2.2 Atomic Formulas

Atomic formulas in first-order logic are predicates, which are applied on terms. We restrict in this chapter to function-free logic, therefore terms are either constants or variables. If all arguments of a predicate are assigned by free variables, the corresponding grounding tensor is stored in the slices of the first axis of  $x_W$  and we have

$$g_k|_{x_W} = \langle x_W[L, O_{[n]}], \epsilon_k[L] \rangle [O_{[n]}] . \quad (34)$$

In contrast, when a constant object  $I_o$  is assigned to an argument of a predicate, the grounding tensor reduced to a slice of the grounding with exclusively free variables. We capture such slicings by contractions with one-hot encodings of the corresponding constant.

We formalize this approach by atom creating tensors, which contraction with the world tensor results in the grounding of the corresponding atomic formula.

**Definition 70.** Let there be an atomic formula  $q$ , which is constructed using the  $l$ -th predicate and has constants assigned on the arguments  $\mathcal{U}^C \subset [n]$  and free variables to the arguments  $\mathcal{U}^V = [n]/\mathcal{U}^C$ . Let the constant map  $C : \mathcal{U}^C \subset [n] \rightarrow [r]$  map to the specific objects represented by the constant and  $V : \mathcal{U}^V \subset [n] \rightarrow \mathcal{V}$  to free variables labeled by a set  $\mathcal{V}$ . Then the atom creating tensor to  $q$  is

$$\psi_q [O_{V(\mathcal{U}^V)}] = \epsilon_l[L] \otimes \left( \bigotimes_{l \in \mathcal{U}^C} \epsilon_{C(l)} [O_l] \right) \otimes \left( \bigotimes_{l \in \mathcal{U}^V} \delta [O_{V(l)}, O_l] \right) .$$

The ground of the atom is then the contraction of the atom creating tensor with the world tensor, that is

$$q|_{x_W} [O_{V(\mathcal{U}^V)}] = \langle x_W[L, O_{[n]}], \psi_q [O_{\mathcal{V}}] \rangle [O_{V(\mathcal{U}^V)}] .$$

What is more abstract, we can understand the predicate itself as an object, then take the first-order world as a grounding tensor of a more abstract formula. We will follow this thought in the ternary representation of Knowledge Graphs in Sect. 14.3.2.

### 14.2.3 Formula synthesis by connectives

In order to have a sound semantic, the grounding of FOL formulas is determined by the syntax of the formula, i.e. a decomposition of the formula into connectives and quantifiers acting on atomic formulas.

Quantifier-free formulas are connectives acting on atomic formulas. We can describe them as in the case of propositional logics in the  $\beta$ -formalism. While the atomic formulas where delta tensors copying states, they are more involved here.

**Theorem 91.** For any connective  $\circ$  and formulas  $q_1$  and  $q_2$  we have

$$(q_1 \circ q_2)|_{x_W} [O_{q_1 \cup q_2}] \quad (35)$$

$$= \left\langle \beta^{q_1}|_{x_W} [Y_{q_1}, O_{q_1}], \beta^{q_2}|_{x_W} [Y_{q_2}, O_{q_2}], \beta^\circ [Y_{q_1 \circ q_2}, Y_{q_1}, Y_{q_2}], \epsilon_1 [Y_{q_1 \circ q_2}] \right\rangle [O_{[n]}] . \quad (36)$$

*Proof.* This directly follows from The. 108. □

Here, variables can be shared by the connected formulas, therefore the variables in the combined formula are unions of the possible not disjoint variables of the connected formulas.

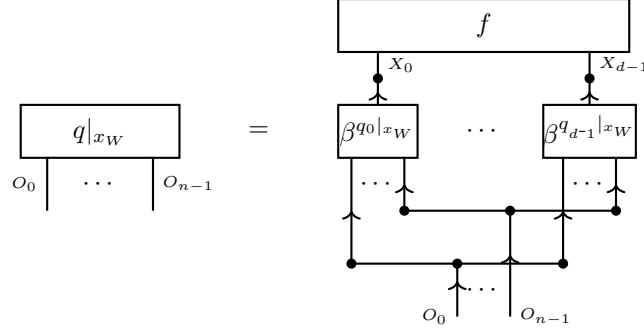
When interpreting the head variables of relational encoded atomic formulas as the atoms of a propositional theory, we find a propositional formula  $f$  associated with any decomposable first order logic formula.

**Definition 71.** Given a formula  $q$  in first order logic, we say that a propositional formula  $f [X_{[d]}]$  is the propositional equivalent to  $q$  given atomic formulas  $q_k$  in first order logic, when for any world  $x_W$  we have

$$q|_{x_W} [O_q] = \left\langle \{ \beta^{q_k}|_{x_W} [X_k, O_{q_k}] : k \in [d] \} \cup \{ f [X_{[d]}] \} \right\rangle [O_q] .$$

We here denote the head variables of the basis encoding to  $\beta^{q_k}|_{x_W}$  by  $X_k$  to highlight their interpretation as propositional atoms.

We depict the relation of a grounding tensor to a propositional formula as:



#### 14.2.4 Quantifiers

Existential and universal quantifiers appear in generic first order logic and are besides substitutions further means to reduce the number of variables in a formula.

The semantics of existential quantification consists in a formula being true, if at least one state of the quantified variable is true, as we define next.

**Definition 72.** Given a grounding tensor

$$q|_{x_W} [O_0, \dots, O_{n-1}]$$

the existential and universal quantification with respect to the first variable are the tensors

$$(\exists_{o_0} q)|_{x_W} [O_1, \dots, O_{n-1}] \quad \text{and} \quad (\forall_{o_0} q)|_{x_W} [O_1, \dots, O_{n-1}]$$

with coordinates as follows. For an assignment  $o_1, \dots, o$  to the non-quantified variables we have

$$(\exists_{o_0} q)|_{x_W} [O_1 = o_1, \dots, O_{n-1} = o_{n-1}] = 1$$

if and only if there is an assignment  $o_0 \in [r_0]$  such that

$$q|_{x_W} [O_0 = o_0, O_1 = o_1, \dots, O_{n-1} = o_{n-1}] = 1.$$

Conversely, we have for the universal quantification that

$$(\forall_{o_0} q)|_{x_W} [O_1 = o_1, \dots, O_{n-1} = o_{n-1}] = 1$$

if and only if for any assignment  $o_0 \in [r_0]$  we have

$$q|_{x_W} [O_0 = o_0, O_1 = o_1, \dots, O_{n-1} = o_{n-1}] = 1.$$

Let us now show, that existential and universal quantification are coordinatewise transforms (see Def. 78) of contracted grounding tensors. To this end, let us introduce the greater- $z$  indicator  $\mathbb{I}_{>z}$ , where  $z \in \mathbb{R}$ , as the function

$$\mathbb{I}_{>:\mathbb{R}} \rightarrow \{0, 1\} \quad , \quad \mathbb{I}_{>z}(x) = \begin{cases} 1 & \text{if } x > z \\ 0 & \text{else} \end{cases}.$$

**Theorem 92.** For any formula  $q$  with variables  $O_{[n]}$  we have

$$(\exists_{o_0} q)|_{x_W} [O_1, \dots, O_{n-1}] = \mathbb{I}_{>0}(\langle q|_{x_W} \rangle [O_1, \dots, O_{n-1}]) [O_1, \dots, O_{n-1}]$$

and

$$(\forall_{o_0} q)|_{x_W} [O_1, \dots, O_{n-1}] = \mathbb{I}_{>r-1}(\langle q|_{x_W} \rangle [O_1, \dots, O_{n-1}]) [O_1, \dots, O_{n-1}]$$

*Proof.* We proof the claimed equalities to arbitrary slices of the remaining variables, which amount to arbitrary substitutions of the formulas. For any indices  $o_1 \in [r_1], \dots, o_{n-1} \in [r_{n-1}]$  we notice, that

$$\begin{aligned} \langle q|_{x_W} \rangle [O_1 = o_1, \dots, O_{n-1} = o_{n-1}] &= \sum_{o_0 \in [r_0]} q|_{x_W} [O_0 = o_0, \dots, O_{n-1} = o_{n-1}] \\ &= |o_0 \in [r_0] : q|_{x_W} [O_0 = o_0, \dots, O_{n-1} = o_{n-1}] = 1|. \end{aligned}$$

We can thus understand the contracted grounding tensor as storing in its coordinates the count of the coordinate extensions to the zeroth variable, such that the grounding tensor is satisfied. This is analogous to our interpretation of contracted propositional formulas as world counts. From this it is obvious, that the existential quantification is satisfied, if the count is different from zero, which is captured by the coordinatewise transform with  $\mathbb{I}_{>0}$ . We therefore arrive at

$$(\exists_{o_0} q)|_{x_W} [O_1 = o_1, \dots, O_{n-1} = o_{n-1}] = \mathbb{I}_{>0} (\langle q|_{x_W} \rangle [O_1, \dots, O_{n-1}]) [O_1 = o_1, \dots, O_{n-1} = o_{n-1}] .$$

The first claim follows, since the assignment to the non-quantified variables was arbitrary. The universal quantification is satisfied, when all extensions are satisfied, and the count is  $r$ . Since  $r$  is the maximal count, this is captured by the coordinatewise transform with  $\mathbb{I}_{>r-1}$  and we get

$$(\forall_{o_0} q)|_{x_W} [O_1 = o_1, \dots, O_{n-1} = o_{n-1}] = \mathbb{I}_{>r-1} (\langle q|_{x_W} \rangle [O_1, \dots, O_{n-1}]) [O_1 = o_1, \dots, O_{n-1} = o_{n-1}] .$$

With the same argument, the second claim is established.  $\square$

We can extend this discussion towards more generic counting quantifiers, of which the existential and the universal quantifier are extreme cases. One can define quantifiers by demanding that at least  $z \in \mathbb{N}$  compatible groundings are satisfied, and show that they amount to coordinatewise transforms with  $\mathbb{I}_{>z}$ . What is more, quantifiers demanding that at most  $z \in \mathbb{N}$  are satisfied would be representable by transforms with an analogously defined function  $\mathbb{I}_{\leq z}$ . Such customized quantifiers appear for example in the OWL 2 standard of description logics (see Rudolph (2011) and Sect. 14.3).

As will be discussed in Chapter 17, any coordinatewise transform can be performed by a contraction of a basis encoding of the tensor with a head vector prepared by the transform function (see The. 109). In the case here, a direct implementation would require a dimension of these head variables by  $r$ , which can be infeasible when having large object sets.

To summarize, let us assume a formula is in its prenex normal form, that is a collection of quantifiers are acting on a quantifier free part. We can represent its grounding tensor by

- Instantiations of the tom groundings with the assigned variables, as contractions of the basis encoding of the world tensor with atom selecting tensors.
- Propositional formula acting on the head variables of the predicate instantiations, representing the connectives combining the formula.
- Quantifiers as a composition of contractions closing the quantified variable and coordinatewise transforms with the respective greater-than indicators.

#### 14.2.5 Storage in basis CP decomposition

In many situations, grounding cores are sparse and representations as single tensor cores comes with a drastic overhead. We often encounter sparse grounding tensors, where the number of non-zero coordinates (to be investigated by basis CP ranks in Chapter 18) satisfies

$$\ell_0(q|_{x_W}) << r^{|O_q|} .$$

In this case, since most coordinates vanish, the basis CP decomposition (see Sect. 18.1.2) enables a representation of the grounding with significantly lower storage demand, see The. 118. This is particularly useful for representing large relational databases, where each object has only a few relations with others, while the majority of possible relations remains unsatisfied. We depict such CP decomposition of a formula grounding in The. 34.

Most logical syntaxes exploit  $\ell_0$ -sparsity, explicitly storing only known assertions. The interpretation of unspecified assertions depends on the underlying assumptions. Under the Closed World Assumption, for example, all unspecified assertions are assumed to be false.

#### 14.2.6 Queries

A database is understood as a specific first order logic world, and are operations on such a single world. Queries are described by a formula  $p$ , which are asked against a specific world  $x_W$  to retrieve the grounding  $p|_{x_W}$ . The variables of such formulas are called projection variables. The answer  $p|_{x_W}$  of a query is most conveniently represented as a list of solution mappings from the projection variables to objects in the world, such that the query formula is satisfied. Answering a query by solution mappings corresponds with finding the basis CP Decomposition (see Sect. 18.1.2) of  $p|_{x_W}$ . We can understand these solution mappings as stored in the leg-matrices  $\rho^{a,l}$  (see Figure 34).

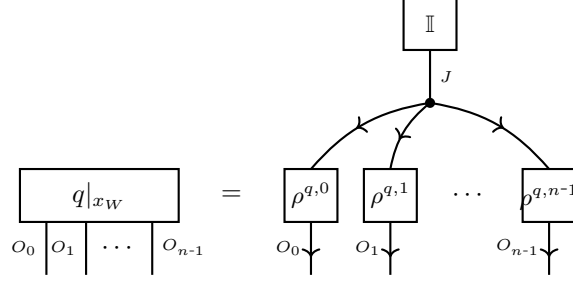


Figure 34: Basis CP Decomposition of the grounding of  $q$ , following the scheme of The. 118. Instead of direct storage of the grounding tensor  $q|_{x_W}$ , the non-zero coordinates are enumerated by a variable  $J$  and the corresponding coordinates stored in leg-matrices  $\rho^{q,l}$ .

Let us give with the outer join an example of a popular operation to define queries, which efficient execution and storage can be improved based on considerations in the tensor network formalism.

**Definition 73** (Outer join). *Let there be a world  $x_W$  and formulas  $q_l$  depending on variables  $O_{\mathcal{V}^l}$ , which have grounding tensors by*

$$q_l|_{x_W}[O_v] : \prod_{v \in \mathcal{V}^l} [r_v] \rightarrow \{0, 1\}.$$

*Then their (outer) JOIN is defined as the grounding of their conjunctions, as*

$$\text{JOIN}(q_0, \dots, q_{p-1})|_{x_W} \left[ \bigcup_{l \in [p]} O_{\mathcal{V}^l} \right] = \langle q_l|_{x_W}[O_{\mathcal{V}^l}] : l \in [p] \rangle \left[ \bigcup_{l \in [p]} O_{\mathcal{V}^l} \right].$$

We can understand the JOIN of groundings by a factor graph, where each grounding tensor decorates the hyperedge to the node set  $\mathcal{V}^l$ . The projection variable assignment to each formula combined in a JOIN operation provide a basic tensor network format to store the output of the operation. There are thus situations, in which the solution map storage corresponding with a CP Decomposition comes with unnecessary overheads compared with other formats.

We can also understand the JOIN operation as a coordinatewise transform (see Def. 78) with the product as transform function. To make this connection solid, one would need to extend each joined formula trivially to the variables appearing in other formulas.

The efficiency of evaluating the contraction to a JOIN operation might be improved by understanding it as an Constraint Satisfaction Problem (see Chapter 8). When applying efficient Message Passing algorithms such as Knowledge Propagation (see Algorithm 8), the groundings can be sparsified by local constraint propagation operations before turning to more global and more demanding contraction operations. Here the groundings  $q_l|_{x_W}$  would be used to initialize Knowledge Cores  $\kappa^e$  and sequentially sparsified during the algorithm.

### 14.3 Representation of Knowledge Graphs

Let us now represent a specific fragment of first-order logic, namely Description Logics which Knowledge Bases are often referred to as Knowledge Graphs. We here use the OWL 2 standard, which encodes the syntax of the description logic  $\mathcal{SROIQ}(\mathcal{D})$  Rudolph (2011).

#### 14.3.1 Representation as unary and binary predicates

Predicates in knowledge graphs are binary (owl:ObjectProperties) and unary (owl:Class). We enumerate the predicates by  $[d]$ , the objects in the domain  $\mathcal{U}$  by  $[r]$ , and extend the unary predicates to binaries by tensor product with  $\epsilon_0[O_1]$ . A Knowledge Graph on the set  $\mathcal{U}$  of constants (owl:NamedIndividuals) is then the tensor

$$\text{KG}|_{x_W}[L, O_0, O_1] : [d] \times [r] \times [r] \rightarrow \{0, 1\}.$$

#### 14.3.2 Representation as ternary predicate

It has been particularly convenient to represent a Knowledge Graph instead as a grounding of a single ternary predicate RDF. To this end, the predicates  $g_k$  and another object rdf:type are added to a domain  $\mathcal{U}$ , by extending the  $r$  and the index interpretation function accordingly.

Following our notation we understand a Knowledge Graph as a grounding of the rdf triple relation RDF (being a formula of order 3) on a specific world  $\text{KG}|_{x_W}$  with individuals  $\mathcal{U}$

We then construct a grounding tensor  $\text{RDF}|_{x_W}$  out of the world  $\text{KG}|_{x_W} [L, O_0, O_1]$  by

$$\text{RDF}|_{x_W} : [r] \times [r] \times [r] \rightarrow \{0, 1\}$$

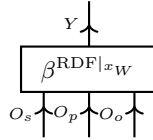
where

$$\begin{aligned} \text{RDF}|_{x_W} [O_s = o_s, O_p = o_p, O_o = o_o] \\ = \begin{cases} \text{KG}|_{x_W} [L = o_s, O_0 = o_o, O_1 = 0] & \text{if } o_p = I^{-1}(\text{rdftype}) \\ \text{KG}|_{x_W} [L = o_p, O_0 = o_s, O_1 = o_o] & \text{if } o_p = I^{-1}(g_k) \text{ for some } k \\ 0 & \text{else} \end{cases} \end{aligned}$$

Slicing the tensor  $\text{RDF}|_{x_W}$  along the predicate axis retrieves specific information about roles and can be efficiently be performed on these formats. The role `rdftype` has a specific meaning, since it contains from a DL perspective classifications (memberships of named concepts). Further slicing the tensor along object axis therefore results in membership lists for specific classes (concepts). One can thus regard `rdftype` as a placeholder for unitary formulas in a space of binary formulas.

Exploiting the  $\ell_0$ -sparsity now leads to a so-called triple store, where  $\text{RDF}|_{x_W}$  is stored by a listing of those triples  $o_s, o_p, o_o$  such that  $\text{RDF}|_{x_W} [O_s = o_s, O_p = o_p, O_o = o_o] = 1$ . A recent implementation of a triple store exploiting these intuitions is TETRIS, see Biggerl et al. (2020). In this work, such decompositions are generalized into more generic CP formats, see Chapter 18. Approximations of grounding tensors by decompositions leads to embeddings of the individuals such as Tucker, ComplEx and RESCAL (see Nickel et al. (2016)).

For our purposes of evaluating logical formulas such as SPARQL queries we use the basis encoding of the groundings, which are depicted by



### 14.3.3 SPARQL Queries

The SPARQL query language is a syntax to express first-order logic formulas  $q$  and intended to be evaluated given a Knowledge Graph. We here consider tensor network representations of the WHERE block. Given a specific knowledge graph  $\text{RDF}|_{x_W}$ , the execution of query is the interpretation  $q|_{x_W}$ , typical represented in a sparse basis CP format where each slice represents a solution mapping.

**Triple Patterns** Central to SPARQL queries are triple patterns, which we understand as slicings of the tensor  $\text{RDF}|_{x_W}$ . To each so-called triple pattern we build a corresponding atom creating tensor (see Def. 70). The triple pattern is then evaluated by contraction of the atom creating tensor with  $\text{RDF}|_{x_W}$ .

Let us now provide examples of such pattern tensors. A unary triple patterns contains a single projection variable, typically related with the subject variable  $O_s$  of  $\text{RDF}|_{x_W}$ . The corresponding pattern tensor is then

$$\psi_{\langle Z, \text{rdftype}, g_k \rangle} [O_s, O_p, O_o, Z] = \delta [O_s, Z] \otimes \epsilon_{I^{-1}(\text{rdftype})} [O_p] \otimes \epsilon_{I^{-1}(g_k)} [O_o] .$$

Binary triple patterns come with two projection variables, typically related with the subject and the object variables  $O_s$  and  $O_o$ . The pattern tensor to the  $k$ -th predicate is then

$$\psi_{\langle Z_0, g_k, Z_1 \rangle} [O_s, O_p, O_o, Z_0, Z_1] = \delta [O_s, Z_0] \otimes \epsilon_{I^{-1}(g_k)} [O_p] \otimes \delta [O_o, Z_1] .$$

Contraction with these pattern tensor evaluated the specific triple pattern, and outputs in a boolean tensor the indicator, which objects are members of a specific class (for unary patterns) or which pair of objects are related by a specific relation. Again, the output of such contractions is a subset encodings of the set of solutions (see Def. 80).

Examples of triple patterns, drawn in Figure 35 are

- Unary triple pattern with one variable, representing a formula with a single projection variable. For the example  $\langle Z, \text{rdftype}, C \rangle$  see Figure 35a.

$$\psi_{\langle Z, \text{rdftype}, g_k \rangle} [O_s, O_p, O_o, Z] = \delta [O_s, Z] \otimes \epsilon_{I^{-1}(\text{rdftype})} [O_p] \otimes \epsilon_{I^{-1}(C)} [O_o]$$

If and only if the output slice is  $\epsilon_1$ , then the corresponding object encoded by the input indices is of class  $C$ .

- Binary triple pattern with two variables, representing a formula with two projection variables. For the example  $\langle Z_0, R, Z_1 \rangle$  see Figure 35b. If and only if the output slice is  $\epsilon_1$ , then the corresponding object tuple encoded by the input indices has a relation  $R$ .

The composition  $\psi(\psi^T)$  of the matrifaction of the tensor  $\psi$  is an orthogonal projection. That means that applying  $\psi(\psi^T)$  is the same map as applying once.

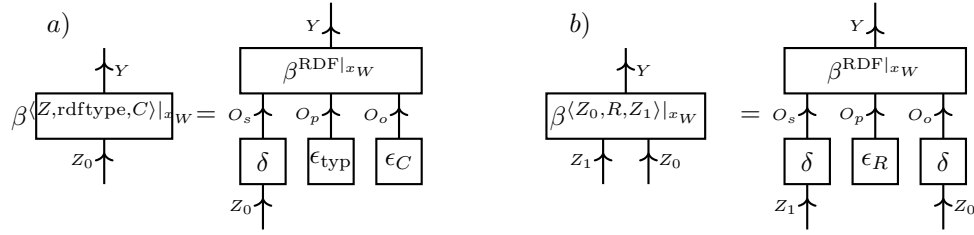


Figure 35: Triple patterns of SPARQL as tensor networks. a) Example of unary triple pattern  $\langle Z, \text{rdftype}, C \rangle$  specifying whether an individual  $I_{o_1}$  is a member of class  $C$ . b) Example of a binary triple pattern  $\langle Z_0, R, Z_1 \rangle$  specifying whether individuals  $I_{o_1}$  and  $I_{o_2}$  have a relation  $R$ . By  $\epsilon_{\text{typ}}, \epsilon_C, \epsilon_R$  we denote the one-hot encodings of the enumeration of the resources  $\text{rdf} : \text{type}, C$  and  $R$ .

**Basic Graph Patterns** Generic SPARQL queries are compositions of triple patterns by logical connectives. These triple patterns possibly share projection variables. Statements in SPARQL can be translated into Propositional Logics combining the triple patterns:

SPARQL	Propositional Logics	Tensor Representation
$\{f_1, f_2\}$	$f_1 \wedge f_2$	$\beta^\wedge [Y_{f_1 \wedge f_2}, Y_{f_1}, Y_{f_2}]$
UNION $\{f_1, f_2\}$	$f_1 \vee f_2$	$\beta^\vee [Y_{f_1 \vee f_2}, Y_{f_1}, Y_{f_2}]$
FILTER NOT EXISTS $\{f\}$	$\neg f$	$\beta^\neg [Y_{\neg f}, Y_f]$

If a SPARQL query consists of these keywords, we find a straight forward corresponding network of triple patterns and encoded logical connectives, by applying our findings of Sect. 14.2.3. To this end, we prepare for each appearing triple pattern the corresponding pattern tensor, and a copy of  $\text{RDF}|_{x_W}$ . Here we also copy the term variables  $O_s, O_p$  and  $O_o$ , to ensure that each copy of  $\text{RDF}|_{x_W}$  shares variables with a single pattern tensor. Projection variables are not copied, since we need to keep track of them shared among triple patterns. Then we prepare the basis encoding of logical connectives according to the hierarchy specified in the SPARQL query. Finally we add a  $\epsilon_1$ -vector to the final head variable representing the complete SPARQL query, to restrict the support to coordinates corresponding with solution mappings. We then contract the resulting tensor network, leaving all projection variables open.

If a projection variable is not appearing in the SELECT statement in front of the WHERE  $\{\cdot\}$ -block, we simply exclude it from the open variables of the described contraction. Note that in that case, the coordinates contain solution counts, i.e. how many assignments to the dropped variable have been a 1 coordinate. We can drop this additional information simply by performing a coordinatewise transform with the greater zero indicator  $\mathbb{I}_{>0}$ .

Here we represented a SPARQL query  $p$  consistent of multiple triple pattern by instantiating a head variables to each triple pattern. Alternatively, the more direct hybrid calculus developed in Sect. 17.6 can be applied and the additional head variables avoided. This is especially compelling, when the WHERE  $\{\cdot\}$ -block does not contain further keywords, i.e. it is the conjunction of all triple patterns. In that case, we avoid the instantiation of head variables (i.e. close the head variables separately by  $\epsilon_1$ -vectors) and represent the query by a contraction of all triple pattern tensors.

We further notice, that any propositional formula acting on the head variables of the triple patterns can be expressed by a hierarchical combination of the key words in the above table. To find the expression, one can transform a given formula into its conjunctive or disjunctive normal form and apply the statements according to the appearing operations  $\wedge, \vee$  and  $\neg$ .

#### 14.4 Probabilistic Relational Models

So far we have studied Markov Logic Networks in Propositional Logics as probability distributions over worlds. In FOL they define probability distributions over relations in worlds with a fixed set of objects. More generally, such models are probabilistic relational models (see for an overview Getoor and Taskar (2019)).

We in this section treat random worlds in first-order logics with fixed domains  $\mathcal{U}$ .

We in this section show, when and how we can interpret likelihoods of Markov Logic Networks in First Order Logic in terms of samples of a Markov Logic Network in Propositional Logics.

##### 14.4.1 Hybrid First-Order Logic Networks

Following Richardson and Domingos (2006) Markov Logic Networks in first-order logics are templates for distributions, which instantiate random worlds when choosing a set of objects  $\mathcal{U}$ . Given a fixed set of constants, they then define a distribution over the worlds, which objects correspond with the constants. This applies database semantics, where only those worlds are considered, where the unique name and domain closure assumptions given a set of constants are satisfied. **Here we directly define them as exponential families distributing  $X_W$  for a given set of objects  $\mathcal{U}$ . To avoid a similar discussion as in Chapter 11 we directly allow for boolean base measures and call the distributions Hybrid First-Order Logic Networks.**

**Definition 74** (Hybrid First-Order Logic Networks (HFLN)). *Let there be a set  $\mathcal{Q}$  of first-order logic formulas with maximal arity  $n$ , which is enumerated by a selection variable  $L$  of dimension  $p$ . Further, let there be a set of objects  $\mathcal{U}$  and a boolean base measure  $\nu [O_{[n]}]$ . The family of Hybrid First-Order Logic Networks  $\Gamma^{\mathcal{Q}|\mathcal{U},\nu}$  defined by the tuple  $(\mathcal{Q}, \mathcal{U}, \nu)$  is the exponential family of joint distributions to the variables  $X_W$  with the statistics*

$$\mathcal{S}_l^{\mathcal{Q}|\mathcal{U}} [X_W = x_W] = \langle q_l|_{x_W} \rangle [\emptyset]$$

and the base measure  $\nu$ .

Each element of the family  $\Gamma^{\mathcal{Q}|\mathcal{U},\nu}$  is represented by a canonical parameter  $\theta [L]$ .

The mean parameter polytope is the convex hull of the vectors

$$\sigma^{\mathcal{Q}} [X_W = x_W, L]$$

to the worlds  $x_W$  with  $\nu [X_W = x_W] = 1$ . These vectors store are the counts of satisfied groundings to each formula, that is

$$\sigma^{\mathcal{Q}} [X_W = x_W, L] = |o_{q_l} : q_l|_{x_W} [O_{q_l} = o_{q_l}] = 1| \ .$$

Each substitution of the variables in  $q_l$  by objects in  $\mathcal{U}$ , which satisfies the formula in the world  $x_W$ , therefore provides a factor of  $\exp [\theta [L = l]]$  to the probability of  $x_W$ .

Let us notice, that different to the case of Hybrid Logic Networks treated in Chapter 11, the statistic does not consist of boolean features, when formulas contain variables and we have multiple objects. One could, however, replace each  $q_l$  by the set of the possible groundings, i.e. substitutions of the formulas variables by any tuple of objects in  $\mathcal{U}$ . The resulting distribution would be an Hybrid Logic Network with boolean statistic, which coincides with the HFLN when posing certain weight sharing conditions on  $\theta$ . The downside of this construction is the increase in the number of features from  $p$  to  $\sum_{l \in [p]} |\mathcal{U}|^{|O_{q_l}|}$ . This polynomial in the cardinality of the domain set increase poses significant computational challenges, see Richardson and Domingos (2006). We will in the next sections explore an alternative way to apply the theory of Chapter 11 and Chapter 12, namely based on importance formulas.

##### 14.4.2 Base measures by importance formulas

The boolean base measure  $\nu$  of a Hybrid First-Order Logic Network is the subset encoding of the possible worlds which have a non-vanishing probability with respect to any member of the family. We now construct specific base measures based on a fixed grounding tensor of an importance formula. This will reduce the number of object tuples influencing the probability distribution in order to arrive at an interpretation of FOL MLNs as likelihoods to datasets of propositional MLNs.

To this end, we mark pairs of term indices relevant to the distributions by an auxiliary index  $j \in [m]$ . Given a set  $\{o_{[n]}^j : j \in [m]\}$  of indices to the important tuples we build a set encoding (see Def. 80)

$$\underline{p} = \sum_{j \in [m]} \left( \bigotimes_{l \in [n]} \epsilon_{o_l^j} [O_l] \right) \ .$$

We interpret the tensor  $\underline{p}$  as the grounding of a formula, which we call the importance formula.

To have a constant importance formula we define a syntactic representation and restrict the support of the HFLN to those world coinciding with groundings of the importance formula coinciding with  $\underline{p}$  by designing a base measure

$$\nu^{\underline{p}}[X_W] = \begin{cases} 1 & \text{if } p|_{x_W}[O_p] = \underline{p} \\ 0 & \text{else} \end{cases}.$$

The base measure restricts the HFLN to be those worlds, where  $p|_{x_W}$  is coincides with the fixed tensor  $\underline{p}$ . Intuitively,  $p|_{x_W}$  represents certain evidence about a first-order logic world, whereas other formulas are uncertain.

**Assumption 1.** Given a base measure  $\nu^{\underline{p}}[X_W]$ , we assume that there is an importance formula  $p[O_{[n]}]$  such that

$$\nu^{\underline{p}}[X_W] = \begin{cases} 1 & \text{if } p|_{x_W}[O_p] = \underline{p} \\ 0 & \text{else} \end{cases}.$$

#### 14.4.3 Decomposition of the log likelihood

To reduce the likelihood of a world to we make the assumption that all formulas in a HFLN are of the form

$$q_l[O_{q_l}] = (p[O_{[n]}] \Rightarrow h_l[O_{q_l}]) \quad (37)$$

that is a rule with the importance formula being the premise. In particular, we assume, that they depend on all term variables  $O_{[n]}$ . If this is not the case, we extend the formula trivially on the missing term variables. When this assumption holds, we can think of the importance formula as a conditions on individuals to satisfy a statistical relation given by  $h$ .

Towards connecting with propositional logics, we further make the assumption, that we can decompose the formula  $h_l$  in what we will call extraction formulas.

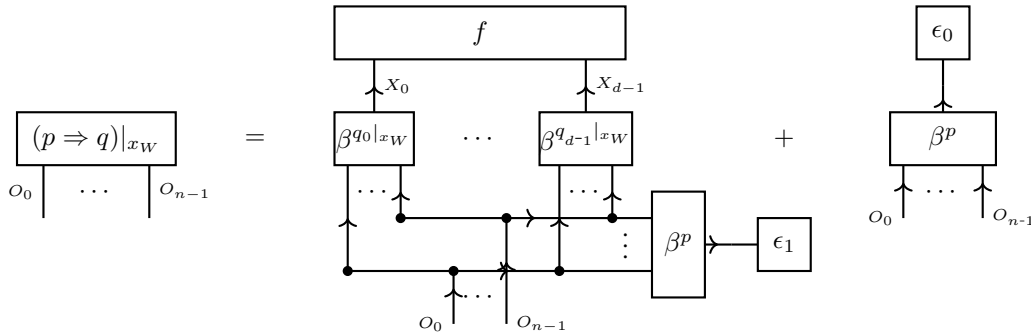
**Assumption 2.** We assume that there exist formulas  $\{q_k[O_{[n]}] : k \in [d]\}$ , which we refer to as atom extraction formulas, and an importance formula  $p[O_{[n]}]$  such that the following holds. To each first-order logic formula  $q_l$  there is another first-order logic formula  $h_l[O_{q_l}]$  and a propositional formula  $f_l[X_{[d]}]$  such that

$$q_l[O_{q_l}] = (p[O_{[n]}] \Rightarrow h_l[O_{q_l}])$$

and

$$h_l[O_{q_l}] = \langle \{f_l[X_{[d]}]\} \cup \{\beta^{q_k}[X_k, O_{[n]}] : k \in [d]\} \rangle [O_{f_l}].$$

We depict the assumption, that any formula is of the form (37) in the diagram



where the second summand depends only on the query  $p$  and therefore does not appear in the likelihood.

Let us now show, how to decompose the probability of a first-order logic world to a HFLN under the above assumptions. Given a HFLN  $\mathbb{P}^{\mathcal{Q}|\mathcal{U}, \theta, \nu^{\underline{p}}}$ , the probability of a world  $x_W$  with  $p|_{x_W} = \underline{p}$  is

$$\mathbb{P}^{\mathcal{Q}|\mathcal{U}, \theta, \nu^{\underline{p}}}[X_W = x_W] = \frac{1}{\mathcal{Z}(\mathcal{Q}|\mathcal{U}, \theta, \nu^{\underline{p}})} \exp \left[ \sum_{l \in [p]} \theta[L = l] \langle (p \Rightarrow h)|_{x_W} \rangle [\emptyset] \right]$$



where the partition function is

$$\mathcal{Z}(\mathcal{Q}|\mathcal{U}, \theta, \nu^{\underline{p}}) = \sum_{x_W : p|_{x_W}[O_{[n]}] = \underline{p}[O_p]} \exp \left[ \sum_{l \in [p]} \theta[L = l] \langle (p \Rightarrow h)|_{x_W} \rangle [\emptyset] \right].$$

Let us now decompose the statistics into constant and varying terms. We have

$$\langle (p \Rightarrow h)|_{x_W} \rangle [\emptyset] = \langle p \wedge h|_{x_W} \rangle [\emptyset] + \langle \neg p|_{x_W} \rangle [\emptyset],$$

where the second term is constant among the supported worlds and the first can be enumerated by the satisfied substitutions of  $p$ , that is

$$\langle p \wedge h|_{x_W} \rangle [\emptyset] = \sum_{j \in [m]} h|_{x_W} [O_{[n]} = o_{[n]}^j].$$

Using these insights we decompose a normalized log likelihood as

$$\frac{1}{m} \ln \left[ \mathbb{P}^{\mathcal{Q}|\mathcal{U}, \theta, \nu^{\underline{p}}} [X_W = x_W] \right] = \frac{1}{m} \sum_{j \in [m]} \sum_{l \in [p]} \theta[L = l] h|_{x_W} [o_{[n]} = o_{[n]}^j] \quad (38)$$

$$- \frac{1}{m} \ln \left[ \frac{\mathcal{Z}(\mathcal{Q}|\mathcal{U}, \theta, \nu^{\underline{p}})}{\exp[\langle \theta \rangle [\emptyset] \cdot \langle \neg p|_{x_W} \rangle [\emptyset]]} \right] \quad (39)$$

We notice a similarity with the likelihood in the case of MLNs in propositional logic. When we interpret each tuple  $o_{[n]} \in (\mathcal{U})^n$  satisfying  $p[O_{[n]} = o_{[n]}] = 1$  as a datapoint, and choose the formulas

$$\mathcal{F} = \{h_l : l \in [p]\}$$

from the propositional equivalents to the formulas  $\mathcal{Q}$ , the first term in (38) coincides with the first term of the likelihood

$$\mathbb{H} \left[ \mathbb{P}^D, \mathbb{P}^{(\mathcal{F}, \theta, \mathbb{I})} \right] = \frac{1}{m} \sum_{j \in [m]} \sum_{l \in [p]} \theta[L = l] f_l[D(j)] - \ln[\mathcal{Z}(\mathcal{F}, \theta)]$$

However, the partition function couples multiple samples, with possible couplings, and prevents a straight forward interpretation as an empirical dataset. We in the next section present assumptions on the tuples satisfying  $p$ , which lead to a factorization of the partition function.

#### 14.4.4 Decomposition of the Partition function

We now make additional assumptions to decompose the partition function of an HFLN as a product of HLN partition functions.

**Assumption 3.** Let  $\nu^{\underline{p}}[X_W]$  be a base measure of worlds such that the vectors

$$\left( q_0|_{x_W} [O_0 = o_0^j, \dots, O_{n-1} = o_{n-1}^j], \dots, q_{d-1}|_{x_W} [O_0 = o_0^j, \dots, O_{n-1} = o_{n-1}^j] \right) \quad (40)$$

for  $j \in [m]$  are independent and identical distributed by the normalization of a boolean base measure  $\nu$ , when drawing  $X_W$  by  $\nu^{\underline{p}}[X_W]$ .

When these assumption holds, we now show that the probability of a first-order logic world coincides with the likelihood of a propositional logic dataset.

**Theorem 93.** Let there be a set of formulas  $\mathcal{Q}$  and a base measure  $\nu^{\underline{p}}[X_W]$  such that Assumption 1, Assumption 2 and Assumption 3 hold. We then have for the likelihood of any by  $\nu^{\underline{p}}[X_W]$  supported world  $x_W$  that

$$\frac{1}{m} \ln \left[ \mathbb{P}^{\mathcal{Q}|\mathcal{U}, \theta, \nu^{\underline{p}}} [X_W = x_W] \right] = \mathbb{H} \left[ \mathbb{P}^D, \mathbb{P}^{\mathcal{F}, \theta} \right]$$

where  $\mathcal{F}$  is the set of propositional equivalents to  $\mathcal{Q}$  (see Assumption 2) and  $D$  the data map with evaluation at  $j \in [m]$  by the enumerated non-vanishing coordinates of  $\underline{p}[O_p]$

$$D(j) = \left( q_0|_{x_W} [O_0 = o_0^j, \dots, O_{n-1} = o_{n-1}^j], \dots, q_{d-1}|_{x_W} [O_0 = o_0^j, \dots, O_{n-1} = o_{n-1}^j] \right).$$

To show the theorem, we show first in the following lemma the factorization of the partition function of the HFLN.

**Lemma 22.** *Given the assumptions of The. 93, we have*

$$\frac{\mathcal{Z}(\mathcal{Q}|_{\mathcal{U}}, \theta, \nu^{\mathcal{L}})}{\exp[\langle \theta \rangle [\emptyset] \cdot \langle \neg p|_{x_W} \rangle [\emptyset]]} = (\mathcal{Z}(\mathcal{F}, \theta, \nu))^m.$$

*Proof.* We have

$$\begin{aligned} \mathcal{Z}(\mathcal{Q}|_{\mathcal{U}}, \theta, \nu^{\mathcal{L}}) &= \mathbb{E}_{x_W \sim \nu^{\mathcal{L}}[X_W]} \left[ \exp \left[ \sum_{q \in \mathcal{Q}} \theta_q \langle (p \Rightarrow h)|_{x_W} \rangle [\emptyset] \right] \right] \\ &= \exp[\langle \theta \rangle [\emptyset] \cdot \langle \neg p|_{x_W} \rangle [\emptyset]] \cdot \mathbb{E}_{x_W \sim \nu^{\mathcal{L}}[X_W]} \left[ \exp \left[ \sum_{q \in \mathcal{Q}} \theta_q \sum_{j \in [m]} h|_{x_W} [O_{[n]} = o_{[n]}^j] \right] \right] \\ &= \exp[\langle \theta \rangle [\emptyset] \cdot \langle \neg p|_{x_W} \rangle [\emptyset]] \cdot \mathbb{E}_{x_W \sim \nu^{\mathcal{L}}[X_W]} \left[ \prod_{j \in [m]} \exp \left[ \sum_{q \in \mathcal{Q}} \theta_q \cdot h|_{x_W} [O_{[n]} = o_{[n]}^j] \right] \right] \end{aligned}$$

Since the substitutions of the atom formulas at the respective object tuples are independent, also the variables

$$\exp \left[ \theta_q \cdot h|_{x_W} [O_{[n]} = o_{[n]}^j] \right]$$

for  $j \in [m]$  are independent. We therefore get

$$\mathcal{Z}(\mathcal{Q}|_{\mathcal{U}}, \theta, \nu^{\mathcal{L}}) = \exp[\langle \theta \rangle [\emptyset] \cdot \langle \neg p|_{x_W} \rangle [\emptyset]] \cdot \prod_{j \in [m]} \mathbb{E}_{x_W \sim \nu^{\mathcal{L}}[X_W]} \left[ \exp \left[ \sum_{q \in \mathcal{Q}} \theta_q \cdot h|_{x_W} [O_{[n]} = o_{[n]}^j] \right] \right] \quad (41)$$

Each  $h|_{x_W} [O_{[n]} = o_{[n]}^j]$  depends by Assumption 2 only on the random tuple  $\{q_k [O_{[n]} = o_{[n]}^j] : k \in [d]\}$ . We build the expectation over all possible values  $x_{[d]}$  of this tuple at any  $j \in [m]$  and get

$$\begin{aligned} &\mathbb{E}_{x_W \sim \nu^{\mathcal{L}}[X_W]} \left[ \exp \left[ \sum_{l \in [p]} \theta [L = l] \cdot h_l|_{x_W} [O_{[n]} = o_{[n]}^j] \right] \right] \\ &= \sum_{x_{[d]} \in \times_{k \in [d]} [2]} \mathbb{P}_{\forall k \in [d] : q_k [O_{[n]} = o_{[n]}^j] = x_k} [x_W \sim \nu^{\mathcal{L}}[X_W]] \cdot \exp \left[ \sum_{l \in [p]} \theta [L = l] \cdot f_l [X_{[d]} = x_{[d]}] \right] \\ &= \sum_{x_{[d]} \in \times_{k \in [d]} [2]} \nu [X_{[d]} = x_{[d]}] \cdot \exp \left[ \sum_{l \in [p]} \theta [L = l] \cdot f_l [X_{[d]} = x_{[d]}] \right] \\ &= \mathcal{Z}(\mathcal{F}, \theta, \nu). \end{aligned}$$

We arrive at the claim, when combining this equation with (41).  $\square$

With this lemma, we are now show The. 93.

*Proof of The. 93.* We have for the logarithm of the probability of a world  $x_W$  given the distribution  $\mathbb{P}^{\mathcal{Q}|_{\mathcal{U}}, \theta, \nu^{\mathcal{L}}}$ , that

$$\ln \left[ \mathbb{P}^{\mathcal{Q}|_{\mathcal{U}}, \theta, \nu^{\mathcal{L}}} [X_W = x_W] \right] = \sum_{l \in [p]} \theta [L = l] \langle q_l|_{x_W} \rangle [\emptyset] - \ln [\mathcal{Z}(\mathcal{Q}|_{\mathcal{U}}, \theta, \nu^{\mathcal{L}})]$$

The first term obeys with Assumption 2

$$\begin{aligned} \sum_{l \in [p]} \theta [L = l] \langle q_l|_{x_W} \rangle [\emptyset] &= \langle \theta \rangle [\emptyset] \cdot \langle \neg p|_{x_W} \rangle [\emptyset] + \sum_{l \in [p]} \sum_{j \in [m]} \theta [L = l] \cdot h_l [O_{[n]} = o_{[n]}^j] \\ &= \langle \theta \rangle [\emptyset] \cdot \langle \neg p|_{x_W} \rangle [\emptyset] + \sum_{l \in [p]} \sum_{j \in [m]} \theta [L = l] \cdot f_l [X_{[d]} = x_{[d]}^j]. \end{aligned}$$

With Lem. 22 we have under the given assumptions for the second term

$$\ln [\mathcal{Z} (\mathcal{Q}|_{\mathcal{U}}, \theta, \nu^p)] = m \cdot \ln [\mathcal{Z} (\mathcal{F}, \theta, \nu)] + \langle \theta \rangle [\emptyset] \cdot \langle \neg p|_{x_W} \rangle [\emptyset] .$$

Combining both, we have

$$\frac{1}{m} \ln [\mathbb{P}^{\mathcal{Q}|_{\mathcal{U}}, \theta, \nu^p} [X_W = x_W]] = \frac{1}{m} \sum_{j \in [m]} \sum_{l \in [p]} \theta [L = l] \cdot f_l [X_{[d]} = x_{[d]}^j] - \ln [\mathcal{Z} (\mathcal{F}, \theta, \nu)]$$

which coincides with  $\mathbb{H} [\mathbb{P}^D, \mathbb{P}^{\mathcal{F}, \theta}]$ .  $\square$

Let us now investigate, in which cases the Assumption 3 of independent data can be matched.

**Lemma 23.** *Let  $p$  and  $q_0, \dots, q_{d-1}$  be quantor and constant free and let the index tuples of the support of  $p$  be pairwise disjoint. Then the vectors (40) are pairwise independent.*

*Proof.* Then we can reduce each sample as dependent only on an independent random world with domain by the respective objects. Quantor and constant-free is needed that this reductions is possible.  $\square$

There are situations, where Assumption 3 is violated. For example, when two object tuples are not disjoint, then some formulas might always coincide on both datapoints, which would violate independence.

In further situations the atom base  $\nu$  are not the uniform  $\mathbb{I}$ :

- extraction formula being a) conjunctions of predicates: Probability that they are satisfied decreases b) disjunctions of predicates: Probability that they are satisfied increases
- extraction formula coinciding with importance formula: Always satisfied, in this case still boolean
- extraction formulas contradicting each other, more general not independent from each other

Let us notice, that non-boolean base measures could be treated in a same manner, but several developments in this work, such as cross-entropy decompositions in Chapter 6 would receive further terms.

**Remark 19** (Approximation by Independent Samples). *As observed above, we do not have independent samples in general. As a consequence, we cannot apply Lem. 22 to decompose the partition function term of the log-probability into factors to each solution map of  $p$ . In this case, it might be still beneficial to use the reduction to the likelihood of a HLN, but needs to understand it as a approximation to the true world probability.*

*If the expectations of each sample with respect to the marginalized distributions coincide, the average of empirical distribution also coincides with these (by linearity). When the creation of samples has sufficient mixing properties, the empirical distribution converges to this expectation in the asymptotic case of large numbers of samples.*

## 14.5 Sample extraction from first-order logic worlds

The decomposition of the likelihood suggests the following approach to generate samples from groundings:

- Define a query formula  $p$ , which we decompose in the basis CP decomposition and interpret each slice as the one-hot encoding of the datapoint.
- Define for  $k \in [d]$  queries  $q_k$  generating the the atoms  $X_k$ : Predicates along with assignment of variables / constants to its positions.
- Contract the groundings of each formula  $q_k$  with the grounding of  $p$  to build a data core

### 14.5.1 Representation by Tensor Networks

We model the extraction process as a relation between a tuple of individuals and the extracted world in the factored system of atoms  $X_k$ .

**Definition 75.** *Given a first-order logic world  $x_W$ , an importance formula  $p$  and extraction formulas  $q_k$  for  $k \in [d]$ , we define the extraction relation*

$$\mathcal{R} \subset \left( \times_{l \in [n]} [r] \right) \otimes \left( \times_{k \in [d]} [2] \right)$$

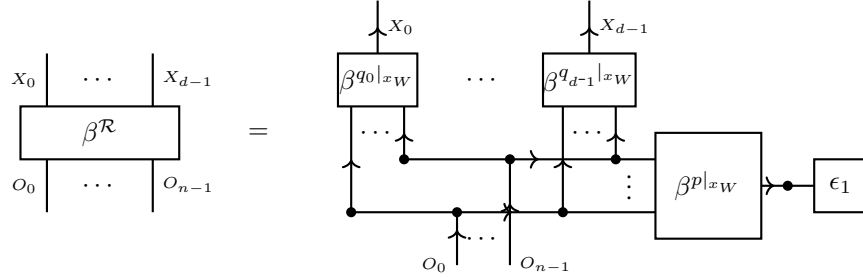
by

$$\mathcal{R} = \{ (o_{[n]}, x_{[d]}) : p|_{x_W} [O_{[n]} = o_{[n]}] = 1, \forall k \in [d] : x_k = q_k [O_{[n]} = o_{[n]}] \} .$$

The encoding of an extraction relation is

$$\beta^{\mathcal{R}} [O_{[n]}, X_{[d]}] \subset \left( \bigotimes_{l \in [n]} \mathbb{R}^r \right) \otimes \left( \bigotimes_{k \in [d]} \mathbb{R}^2 \right)$$

and drawn in a contraction diagram by



Here the contraction of  $\beta^p$  with the truth vector  $\epsilon_1$  represents the matching condition posed by  $p$  when extracting pairs of individuals.

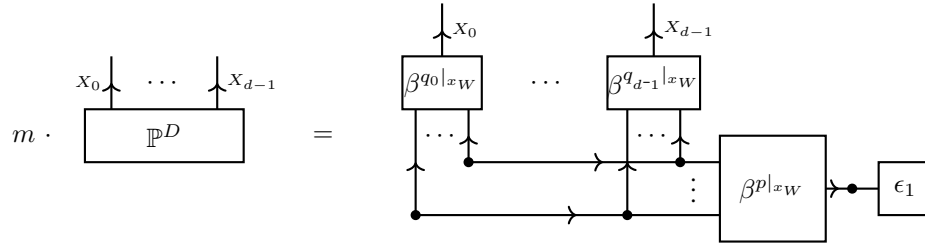
The empirical distribution is then the normalized contraction leaving only the legs to the extracted atomic formulas open, that is

$$\mathbb{P}^D = \frac{\langle \beta^{\mathcal{R}} \rangle [X_{[d]}]}{\langle \beta^{\mathcal{R}} \rangle [\emptyset]}.$$

Here the number of extracted data is the denominator

$$m = \langle \beta^{\mathcal{R}} \rangle [\emptyset] = \langle \beta^p [Y_p, O_{[n]}], \epsilon_1 [Y_p] \rangle [\emptyset].$$

We depict this by



#### 14.5.2 Basis CP Decomposition of extracted data

To connect with the empirical distribution introduced in Sect. 6.3.1 we now show how the empirical distribution extracted from the interpretations of the formulas  $p, q_0, \dots, q_{d-1}$  on a first-order logic world  $x_W$  can be represented by tensor networks.

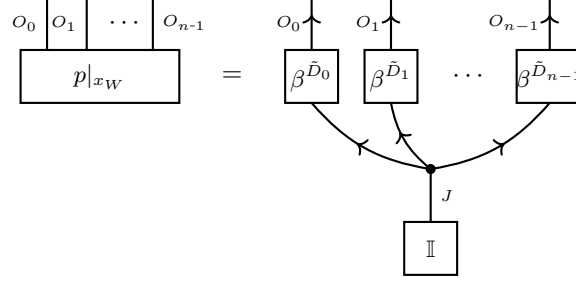
First of all, we decompose the importance formula into a basis CP format (see Chapter 18), that is a decomposition

$$p|_{x_W} [O_{[n]}] = \langle \{\rho^l [O_l, J] : l \in [n]\} \rangle [O_{[n]}]$$

such that all  $\rho^l [O_l, J]$  are directed and boolean tensors. Here an auxiliary variables  $J$  taking values in  $[m]$  is introduced, which we call the data variable, which enumerates the non-vanishing coordinates of  $p|_{x_W}$ . With this decomposition, we can understand the decomposition of  $p|_{x_W} [O_{[n]}]$  as a basis encoding of an term selection map  $\tilde{D}$  with coordinate maps defined such that

$$\beta^{\tilde{D}_l} [O_l, J] = \rho^l [O_l, J].$$

We depict this decomposition by:



Based on these construction, we now provide a tensor network decomposition of the extracted empirical distribution.

**Theorem 94.** *Given a first-order logic world  $x_W$ , an importance formula  $p$  and extraction formulas  $q_k$  for  $k \in [d]$ , we have*

$$\beta^{\mathcal{R}} [O_{[n]}, X_{[d]}] = \left\langle \{ \beta^{q_k|_{x_W}} [X_k, O_{[n]}] : k \in [d] \} \cup \{ \beta^{\tilde{D}_l} [O_l, J] : l \in [n] \} \right\rangle [O_{[n]}, X_{[d]}]$$

and thus

$$\mathbb{P}^D [X_{[d]}] = \frac{1}{m} \left\langle \{ \beta^{q_k|_{x_W}} [X_k, O_{[n]}] : k \in [d] \} \cup \{ \beta^{\tilde{D}_l} [O_l, J] : l \in [n] \} \right\rangle [X_{[d]}] .$$

*Proof.* To show the first claim, let us choose arbitrary state tuples  $o_{[n]}$  and  $x_{[d]}$ . We then have

$$\begin{aligned} & \left\langle \{ \beta^{q_k|_{x_W}} [X_k, O_{[n]}] : k \in [d] \} \cup \{ \beta^{\tilde{D}_l} [O_l, J] : l \in [n] \} \right\rangle [O_{[n]} = o_{[n]}, X_{[d]} = x_{[d]}] \\ &= \left\langle \{ \beta^{q_k|_{x_W}} [X_k = x_k, O_{[n]} = o_{[n]}] : k \in [d] \} \cup \{ \beta^{\tilde{D}_l} [O_l = o_l, J] : l \in [n] \} \right\rangle [\emptyset] . \end{aligned}$$

This contraction evaluates to 1, if and only if for all  $k \in [d]$  we have  $\beta^{q_k|_{x_W}} [X_k, O_{[n]}] = 1$  and

$$\left\langle \{ \beta^{\tilde{D}_l} [O_l = o_l, J] : l \in [n] \} \right\rangle [\emptyset] = 1 .$$

The first condition is equal to  $x_k = q_k [O_{[n]} = o_{[n]}]$  for all  $k \in [d]$  and the second to

$$p|_{x_W} [O_{[n]} = o_{[n]}] = 1 .$$

Comparing with the definition of the extraction relation (see Def. 75), we notice that these conditions are equal to  $(O_{[n]}, x_{[d]}) \in \mathcal{R}$  and therefore to

$$\beta^{\mathcal{R}} [O_{[n]} = o_{[n]}, X_{[d]} = x_{[d]}] .$$

The first claim follows, since  $\beta^{\mathcal{R}}$  is boolean, as is the contraction of the cores  $\beta^{q_k|_{x_W}}$  with the cores  $\beta^{\tilde{D}_l}$ , which leaves the outgoing variables  $X_{[d]}$  open. The second claim follows from the first using that  $\mathbb{P}^D [X_{[d]}] = \frac{1}{m} \langle \beta^{\mathcal{R}} \rangle [X_{[d]}]$ .  $\square$

To connect with the representation of empirical distributions based on data cores (see Sect. 6.3.1), we now form data cores by contractions with the grounding of extraction formulas with the cores  $\beta^{\tilde{D}_l}$  (see Figure 36),

$$\beta^{D_k} [X_k, J] = \left\langle \{ \beta^{q_k|_{x_W}} [X_k, O_{[n]}] \} \cup \{ \rho^l [O_l, J] : l \in [n] \} \right\rangle [X_k, J] .$$

The empirical distribution is then a tensor network of these tensors, as we show next.

**Theorem 95.** *We have*

$$\langle \beta^{\mathcal{R}} \rangle [X_{[d]}] = \langle \{ \beta^{D_k} [J, X_k] : k \in [d] \} \rangle [X_{[d]}]$$

and thus

$$\mathbb{P}^D [X_{[d]}] = \frac{1}{m} \langle \{ \beta^{D_k} [J, X_k] : k \in [d] \} \rangle [X_{[d]}] .$$

*Proof.* By The. ?? we have

$$\beta^{\mathcal{R}} [O_{[n]}, X_{[d]}] = \left\langle \{ \beta^{q_k|_{x_W}} [X_k, O_{[n]}] : k \in [d] \} \cup \{ \beta^{\tilde{D}_l} [O_l, J] : l \in [n] \} \right\rangle [O_{[n]}, X_{[d]}] .$$

Since  $\beta^{\tilde{D}_l} [O_l, J]$  are directed and boolean, they can be copied and separately contracted with each  $q_k|_{x_W}$ , without changing the contraction. We arrive at

$$\begin{aligned} \beta^{\mathcal{R}} [O_{[n]}, X_{[d]}] &= \left\langle \left\{ \{ \beta^{q_k|_{x_W}} [X_k, O_{[n]}] \} \cup \{ \beta^{\tilde{D}_l} [O_l, J] : l \in [n] \} \right\} [X_k, J] : k \in [d] \right\rangle [O_{[n]}, X_{[d]}] \\ &= \left\langle \{ \beta^{D_k} [J, X_k] : k \in [d] \} \right\rangle [X_{[d]}] , \end{aligned}$$

which established the claim.  $\square$

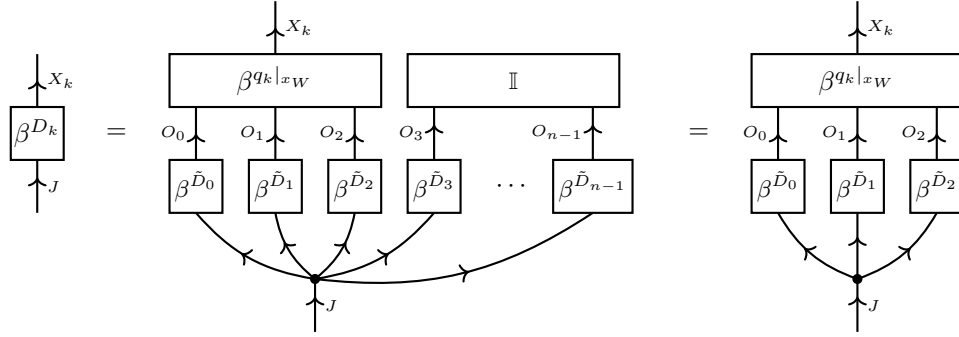


Figure 36: Generation of a data core for the variable  $X_k$  given an extraction formula  $q_k$  and an importance formula, which grounding is decomposed into a basis CP format with leg vectors  $\beta^{\tilde{D}_l} [O_l, J]$ . Term variables, which are appearing in the importance formula, but not in the extraction formula  $q_k$  can be treated trivially by contraction with the trivial tensor (here  $O_4, \dots, O_{n-1}$ ).

When many atom extraction formulas differ only by a constant, we can replace the constant by an auxiliary term variable. The atoms are then the atomizations of this variable (see Sect. 8.3.2), treated as a categorical variable, with respect to the constant in the extraction query. The advantages are that we can avoid the  $\beta$ -formalism and directly model the categorical distributions.

This also enables a batchwise computation of multiple SPARQL queries, which differ only in one constant.

#### 14.6 Generation of first-order logic worlds

So far we have discussed, how MLNs for FOL Knowledge Bases such as Knowledge Graphs can be built by extracting data. Conversely, any binary tensor can be interpreted as a Knowledge Graph. To be more precise, we follow the intuition that the ones coordinates mark possible worlds compatible with the knowledge about a factored system. Each possible world can then be encoded in a subgraph of the Knowledge Graph representing the world. This amounts to an "inversion" of the data generation process described in the subsection above.

In the previous section we have described a way to extract an effective empirical distribution for the likelihood of a first-order logic world given a HFLN. We now want to investigate methods to reproduce an empirical distribution based on a constructed first-order logic world.

**Definition 76** (Reproduction of Empirical Distributions). *Given an empirical distribution  $\mathbb{P}^D \in \bigotimes_{k \in [d]} \mathbb{R}^2$ , we say that a triple  $(x_W, p, q_{[d]})$  of a FOL world  $x_W$  an importance formula  $p$  and extraction formulas  $q_{[d]} = \{q_k : k \in [d]\}$  reproduces  $\mathbb{P}^D$ , when*

$$\mathbb{P}^D = \left\langle \{ p|_{x_W} \} \cup \{ \beta^{q_k|_{x_W}} : k \in [d] \} \right\rangle [X_{[d]}|\emptyset] .$$

Note that for distribution  $\mathbb{P}$  to be reproducible, it needs to have rational coordinates. If any only if all coordinates are rational, we find a  $m \in \mathbb{N}$  such that  $\text{im}(m \cdot \mathbb{P}) \subset \mathbb{N}$ . We can then interpret  $m$  as the number of samples, and construct

a sample selector map by understanding each coordinate of  $m \cdot \mathbb{P}$  as the number of appearances of the respective world in the samples.

We show different schemes and give examples on Knowledge Graphs, where we provide examples for importance and extraction formulas by SPARQL queries.

#### 14.6.1 Samples by single objects

In the first reproduction scheme we construct datapoints by dedicated objects, which represent a sample, that is we choose a domain  $\mathcal{U} = [m]$ .

**Theorem 96.** *Let there be an empirical distribution  $\mathbb{P}^D$  to a sample selector map  $D$  (see Def. 37), we construct a world  $x_W[L, O]$  with  $d$  unary predicates by*

$$x_W[L, O] = \sum_{k \in [d]} \sum_{j \in [m] : D_k(j)=1} \epsilon_k[L] \otimes \epsilon_j[O] .$$

We further choose a trivial importance query, that is

$$p|_{x_W}[O] = \mathbb{I}[O] ,$$

and extraction queries coinciding with the unary predicates, that is for  $k \in [d]$

$$q_k = g_k .$$

Then, the triple  $(x_W, p, q_{[d]})$  reproduces  $\mathbb{P}^D$ .

*Proof.* By The. 95 it is enough to show, that the data cores constructed from the data extraction process coincide with those of  $\mathbb{P}^D$ . We enumerate to this end the non-vanishing coordinates of  $p|_{x_W}$  by the data variable  $J$  taking values  $j \in [m]$ , as

$$p|_{x_W}[O = j] = 1$$

and choose

$$\tilde{D} = \delta .$$

For arbitrary  $k \in [d]$  and  $j \in [m]$  we now have

$$\begin{aligned} \beta^{D_k}[X_k, J = j] &= \left\langle \beta^{q_k|_{x_W}}[X_k, O], \rho^0[O, J = j] \right\rangle [X_k, J] \\ &= \left\langle \beta^{q_k|_{x_W}}[X_k, O], \epsilon_{\tilde{D}(j)}[O] \right\rangle [X_k, J = j] \\ &= \epsilon_{D_k(j)}[X_k] . \end{aligned}$$

This coincides with the slice of the data core of the CP representation of empirical distributions used in The. 23. Since the slice and the core was arbitrary, the tensor network representations in The. 23 and The. 95 are equal and thus the triple  $(x_W, p, q_{[d]})$  reproduces  $\mathbb{P}^D$ .  $\square$

We now give by the next theorem an example of a Knowledge Graph with SPARQL queries reproducing an arbitrary empirical distribution.

**Theorem 97.** *Let  $\mathbb{P}^D$  be an empirical distribution to the sample selector  $D$ . We construct a Knowledge Graph of the resources  $\mathcal{U} = \{s_j : j \in [m]\} \cup \{C\} \cup \{C_k : k \in [d]\}$ , where  $s_j$  represent samples and  $C_k$  unary predicates, by*

$$\text{RDF}|_{x_W} = \sum_{j \in [m]} \epsilon_{I_{s_j}} O_s \otimes \epsilon_{I_{\text{rdf type}}} O_p \otimes \epsilon_{I_C} O_o + \sum_{j \in [m]} \sum_{k \in [d] : D_k(j)=1} \epsilon_{I_{s_j}} O_s \otimes \epsilon_{I_{\text{rdf type}}} O_p \otimes \epsilon_{I_{C_k}} O_o .$$

We further define an importance formula by the SPARQL query

$$p = \text{SELECT } \{ ?x \} \text{ WHERE } \{ ?x \quad \text{rdf} : \text{type} \quad C \quad . \}$$

and for each  $k \in [d]$  an extraction formula by the query

$$q_k = \text{SELECT } \{ ?x \} \text{ WHERE } \{ ?x \quad \text{rdf} : \text{type} \quad C_k \quad . \} .$$

Then the triple  $(\text{KG}|_{x_W}, p, q_{[d]})$  reproduces  $\mathbb{P}^D$ .

*Proof.* We show the theorem analogously to The. 96, with the slide difference in the importance formula. We have for the grounding of  $p$  on  $\text{KG}|_{x_W}$  that

$$p|_{x_W} [O] = \sum_{j \in [m]} \epsilon_{I_{s_j}} O$$

and enumerate the non-vanishing coordinates by  $J$ .

For each extraction formula we have

$$q_k|_{x_W} [O] = \sum_{j \in [m] : D_k(j)=1} \epsilon_{I_{s_j}} O.$$

It follows that the data cores used in The. 95 are

$$\beta^{D_k} [X_k, j] = \epsilon_0 [X_k] \otimes \left( \sum_{j \in [m] : D_k(j)=0} \epsilon_j [J] \right) + \epsilon_1 [X_k] \otimes \left( \sum_{j \in [m] : D_k(j)=1} \epsilon_j [J] \right)$$

and they thus coincide with those in the decomposition in The. 23. The claim follows therefore with the same argumentation as in the proof of The. 96.  $\square$

Let us provide some more insights on the construction of the reproducing Knowledge Graph in The. 97. By the insertions to the one-hot encodings  $\epsilon_{I_{s_j}} O_s \otimes \epsilon_{I_{\text{rdftype}}} O_p \otimes \epsilon_{I_C} O_o$  we mark each sample representing resource by a class and ensure its appearance as a owl : NamedIndividual in the graph. The insertions  $\epsilon_{I_{s_j}} O_s \otimes \epsilon_{I_{\text{rdftype}}} O_p \otimes \epsilon_{I_{C_k}} O_o$  on the other side encode the sample selecting map, by inserting exactly the assertions corresponding with the respective sample. In this simple Knowledge Graph, Description Logic is expressive enough to represent any formula  $q$  composed of the formulas  $q_0, \dots, q_{d-1}$ .

#### 14.6.2 Samples by pairs of objects

We now instantiate multiple objects for each datapoint, one for each variable of the importance formula, i.e.  $\mathcal{U} = [m] \times [n]$  Label individuals  $s_{j,l}$  by data index and variable index.

**Lemma 24.** *Let there a data map  $D$ , queries  $p, q_{[d]}$  and a first-order logic world containing objects  $s_{j,l}$  for  $j \in [m]$  and  $l \in [n]$  If*

$$p|_{x_W} = \sum_{j \in [m]} \bigotimes_{l \in [n]} \epsilon_{I_{s_{j,l}}} [O_l]$$

and for any  $k \in [d]$

$$q_k|_{x_W} = \sum_{j : D_k(j)=1} \bigotimes_{O_l \in O_{q_k}} \epsilon_{I_{s_{j,l}}} [O_l].$$

Then the tuple  $(\text{KG}|_{x_W}, p, \{q_0, \dots, q_{d-1}\})$  reproduces  $\mathbb{P}^D$ .

*Proof.* We notice, that the grounding of the importance formula is in a basis CP format, since by assumption

$$p|_{x_W} = \sum_{j \in [m]} \bigotimes_{l \in [n]} \epsilon_{I_{s_{j,l}}} [O_l].$$

We choose  $J$  to enumerate the non-vanishing entries and get a term selecting map

$$\tilde{D}_l(j) = I_{s_{j,l}}.$$

From this we have

$$\left\langle \{ \beta^{q_k|_{x_W}} [X_k, O_{q_k}] \} \cup \{ \beta^{\tilde{D}_l} [O_l, J] : l \in [n] \} \right\rangle [X_k, J] = \beta^{D_k} [X_k, J]$$

and the claim follows with the same argumentation as in the proof of The. 96.  $\square$



## 14.7 Discussion

Statistical Models are called Probabilistic Relational Models. Extensions are models that also handle structural uncertainty, i.e. distributions of worlds with varying  $\mathcal{U}$ .

In the emerging area of network science Barabási (2016); Giovanni Russo Vito Latora (2017), statistical models for random graphs are investigated. Statistical Models of first-order logic go beyond the typical single edge type perspective of network science.

**Remark 20** (Alternative Representation of empirical distributions). *So far, we have motivated the representation of empirical distributions based on basis CP decompositions based on data maps. In this section, based on the extraction queries, we have observed that empirical distributions might have more efficient representation formats. In many applications such as the computation of log-likelihoods we can use any representation of the empirical distribution by tensor networks. It is thus not necessary to compute the data cores as above, unless one requires a list of the extracted samples.*

## Part III

# Contraction Calculus

## 15 Introduction into Part III

We frequently worked in Part I and Part II with tensors, which have non-negative coordinates and occasionally are boolean (see Def. 4) or directed (see Def. 12). While boolean tensors have appeared as semantical representation of formulas, directed tensors have appeared mostly as conditional distributions. In this chapter we provide further insights into the situation, where tensors satisfy both. For a schematic depiction of this see Figure 37.

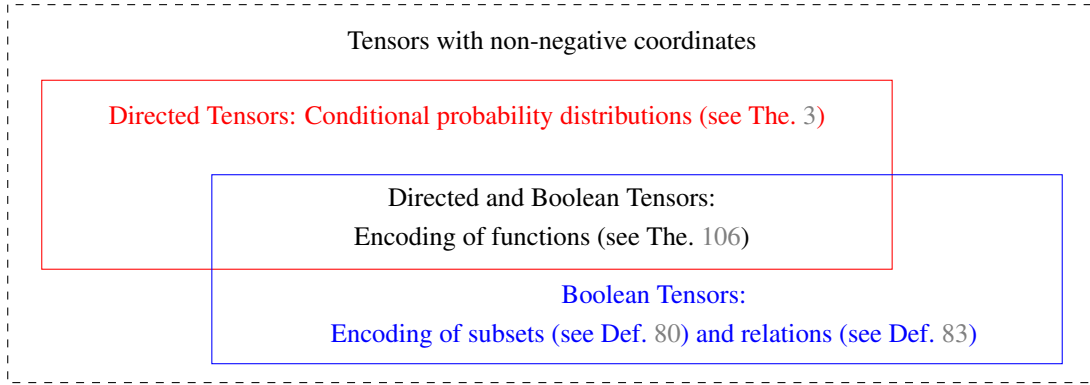


Figure 37: Sketch of the tensors with non-negative coordinates. We investigate in this chapter tensors, which are directed and boolean.

## 16 Coordinate Calculus

In the previous chapters, information to states has been stored in coordinates of a tensor. To distinguish from other schemes of calculus such as the basis calculus (see Chapter 17), we call this scheme of storing and retrieving information the coordinate calculus.

### 16.1 One-hot encodings as basis

Let us first state, that the one-hot encodings, which we have used to motivate tensor representations, build an orthonormal basis of the respective tensor spaces.

**Lemma 25.** *The image of the one-hot encoding map is an orthonormal basis of the tensor space  $\bigotimes_{k \in [d]} \mathbb{R}^{m_k}$ , that is for any  $x_{[d]}, \tilde{x}_{[d]} \in \times_{k \in [d]} [m_k]$  we have*

$$\langle \epsilon_{x_{[d]}} [X_{[d]}], \epsilon_{\tilde{x}_{[d]}} [X_{[d]}] \rangle [\emptyset] = \delta_{x_{[d]}, \tilde{x}_{[d]}} := \begin{cases} 1 & \text{if } x_{[d]} = \tilde{x}_{[d]} \\ 0 & \text{else} \end{cases}.$$

Any element  $\tau \in \bigotimes_{k \in [d]} \mathbb{R}^{m_k}$  has a decomposition

$$\tau [X_{[d]}] = \sum_{x_{[d]} \in \times_{k \in [d]} [m_k]} \tau [X_{[d]} = x_{[d]}] \cdot \epsilon_{x_{[d]}} [X_{[d]}].$$

We notice that the coordinates are the weights to the basis elements in the one-hot decomposition.

*Proof.* The first claim follows from an elementary decomposition of one-hot encodings and the orthogonality of basis vectors as

$$\langle \epsilon_{x_{[d]}} [X_{[d]}], \epsilon_{\tilde{x}_{[d]}} [X_{[d]}] \rangle [\emptyset] = \prod_{k \in [d]} \langle \epsilon_{x_k} [X_k], \epsilon_{\tilde{x}_k} [X_k] \rangle [\emptyset] = \prod_{k \in [d]} \delta_{x_k, \tilde{x}_k} = \delta_{x_{[d]}, \tilde{x}_{[d]}}.$$

To show the second claim, it is enough to notice that for any  $\tilde{x}_{[d]} \in \times_{k \in [d]} [m_k]$  we have

$$\begin{aligned} \sum_{x_{[d]} \in \times_{k \in [d]} [m_k]} \tau [X_{[d]} = x_{[d]}] \cdot \epsilon_{x_{[d]}} [X_{[d]} = \tilde{x}_{[d]}] &= \sum_{x_{[d]} \in \times_{k \in [d]} [m_k]} \tau [X_{[d]} = x_{[d]}] \cdot \delta_{x_{[d]}, \tilde{x}_{[d]}} \\ &= \tau [X_{[d]} = \tilde{x}_{[d]}]. \end{aligned}$$

□

Any tensor can be understood as a coordinate encoding of a real-valued function, as we define next.

**Definition 77.** *Given any real-valued function*

$$q : \times_{k \in [d]} [m_k] \rightarrow \mathbb{R}$$

*we define the coordinate encoding by*

$$\chi^q [X_{[d]}] = \sum_{x_{[d]} \in \times_{k \in [d]} [m_k]} q(x_{[d]}) \cdot \epsilon_{x_{[d]}} [X_{[d]}].$$

In Part I and Part II we did not distinguish between a real-valued function  $q$  and its coordinate encoding  $\tau^q$ , in order to abbreviate notation. Based on coordinate encodings, we now show, that function evaluation can be performed by contractions.

**Theorem 98** (Function evaluation in Coordinate Calculus). *Given any real-valued function*

$$q : \times_{k \in [d]} [m_k] \rightarrow \mathbb{R}$$

*and any input state  $x_{[d]} \in \times_{k \in [d]} [m_k]$ , we have*

$$q(x_{[d]}) = \langle \chi^q [X_{[d]}], \epsilon_{x_{[d]}} [X_{[d]}] \rangle [\emptyset].$$

*Proof.* We use the decomposition in Lem. 25 and have by linearity of contractions for any index tuple  $x_{[d]} \in \times_{k \in [d]} [m_k]$

$$\begin{aligned} \langle \chi^q [X_{[d]}], \epsilon_{x_{[d]}} [X_{[d]}] \rangle [\emptyset] &= \sum_{\tilde{x}_{[d]} \in \times_{k \in [d]} [m_k]} \chi^q [X_{[d]} = \tilde{x}_{[d]}] \cdot \langle \epsilon_{\tilde{x}_{[d]}} [X_{[d]}], \epsilon_{x_{[d]}} [X_{[d]}] \rangle [\emptyset] \\ &= \sum_{\tilde{x}_{[d]} \in \times_{k \in [d]} [m_k]} q(\tilde{x}_{[d]}) \cdot \delta_{\tilde{x}_{[d]}, x_{[d]}} \\ &= q(x_{[d]}) \end{aligned}$$

where we used that one-hot encodings are orthonormal.

□

Coordinate calculus is the representation of real-valued functions as tensors, from which its evaluations can be retrieved by the scheme of The. 98. This is in contrast to the basis calculus scheme to be discussed (see The. 107), where the contraction-based evaluations of functions outputs one-hot encodings.

Tensors of large orders often admit a decomposition by tensor networks. We in the next theorem show, how such a decomposition can be exploited for efficient contractions and in particular coordinate retrieval.

**Theorem 99.** *Given a tensor network  $\tau^{\mathcal{G}}$  on a hypergraph  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ , disjoint subsets  $A, B \subset \mathcal{V}$  and  $x_B \in [m_B]$ , we have*

$$\langle \tau^{\mathcal{G}} \rangle [X_A, X_B = x_B] = \langle \{ \langle \tau^e \rangle [X_{e/B}, X_{e \cap B} = x_{e \cap B}] : e \in \mathcal{E} \} \rangle [X_A] .$$

*Proof.* By definition of contractions we have for any  $x_A$

$$\begin{aligned} \langle \tau^{\mathcal{G}} \rangle [X_A = x_A, X_B = x_B] &= \sum_{x_{\mathcal{V}/(A \cup B)} \in \mathcal{X}_{\mathcal{V}/(A \cup B)}[m_{\mathcal{V}}]} \prod_{e \in \mathcal{E}} \tau^e [X_{e/B} = x_{e/B}, X_{e \cap B} = x_{e \cap B}] \\ &= \langle \{ \langle \tau^e \rangle [\emptyset] X_{e/(A \cup B)}, X_{e \cap A} = x_{e \cap A}, X_{e \cap B} = x_{e \cap B} : e \in \mathcal{E} \} \rangle [\emptyset] \\ &= \langle \{ \langle \tau^e \rangle [\emptyset] X_{e/B}, X_{e \cap B} = x_{e \cap B} : e \in \mathcal{E} \} \rangle [X_A = x_A] \end{aligned}$$

and the claim follows.  $\square$

If we retrieve a single coordinate of a tensor, we have the situation  $A = \emptyset, B = \mathcal{V}$ . In that case, Theorem 99 shows, that the coordinate is the product of the coordinates of the cores.

## 16.2 Coordinatewise Transforms

Let us now discuss a scheme to perform transformations of tensors. We call them coordinatewise, when the target tensor has the same variables as the input tensors, and each coordinate of the target tensor depends only on the respective coordinates of the input tensors.

**Definition 78.** *Let  $h : \mathbb{R}^p \rightarrow \mathbb{R}$  be a function. Then the coordinatewise transform of tensors  $\tau^l [X_{[d]}]$ , where  $l \in [p]$ , under  $q$  is the tensor*

$$h(\tau^0, \dots, \tau^{p-1}) [X_{[d]}]$$

with coordinates

$$h(\tau^0, \dots, \tau^{p-1}) [X_{[d]} = x_{[d]}] = h(\tau^0 [X_{[d]} = x_{[d]}], \dots, \tau^{p-1} [X_{[d]} = x_{[d]}]) .$$

Coordinatewise transforms in case of  $p = 1$  have been indicated by ellipses in the diagrammatic depiction of contractions. We will provide a generic tensor network representation in Chapter 17, see The. 109.

In the following lemma, we state that coordinatewise transforms can be restricted to slices of tensors, when Although this is an obvious fact, this property can tremendously reduce the computational demand of contractions with coordinatewise transforms of tensors.

**Lemma 26.** *For any function  $h : \mathbb{R} \rightarrow \mathbb{R}$ , any tensor  $\tau [X_{[d]}]$  and index  $x_A$ , where  $A \subset [d]$ , we have*

$$h(\tau [X_{[d]}]) [X_{[d]/A}, X_A = x_A] = h(X_{[d]/A}, X_A = x_A) [X_{[d]/A}] .$$

*Proof.* For any state  $x_{[d]/A}$  we have that

$$\begin{aligned} h(\tau [X_{[d]}]) [X_{[d]/A} = x_{[d]/A}, X_A = x_A] &= h(\tau [X_{[d]} = x_{[d]}]) \\ &= h(X_{[d]/A}, X_A = x_A) [X_{[d]/A} = x_{[d]/A}] . \end{aligned}$$

$\square$

**Example 17** (Hadamard products as coordinatewise transforms). *Hadamard products of tensors (see Example 3) are a special way of coordinate calculus, where the transform is the product and thus*

$$\cdot (\tau^0, \dots, \tau^{p-1}) [X_{[d]}] = \langle \{ \tau^l [X_{[d]}] : l \in [p] \} \rangle [X_{[d]}] .$$

These hadamard products are applied in the effective computation of conjunctions, as we will discuss in more detail in Sect. 17.6.

**Example 18** (Exponentiation of energies). In Def. 24 we introduced exponential families, based on the exponentiation of energies. For a statistic  $\mathcal{S}$ , a base measure  $\nu$  and a canonical parameter  $\theta$  we defined

$$\mathbb{P}^{(\mathcal{S}, \theta, \nu)} = \frac{\langle \exp [\langle \sigma^{\mathcal{S}}, \theta \rangle [X_{[d]}], \nu [X_{[d]}] \rangle [X_{[d]}]}{\langle \exp [\langle \sigma^{\mathcal{S}}, \theta \rangle [X_{[d]}], \nu [X_{[d]}] \rangle [\emptyset]}.$$

Both the nominator and the denominator involve a coordinatewise transform of the energy tensor  $\phi^{(\mathcal{S}, \theta, \nu)}[X_{[d]}]$  by the exponentiation. The. 11 provided a transform-free contraction expression by basis encodings, which is the central tool of basis calculus (see Chapter 17).

Let us note, that Lem. 26 enables the energy-based answering of conditional queries, as has been shown in The. 22.

### 16.3 Directed Tensors

Directionality as defined in Def. 12 is a constraint on the structure of a tensor, namely that the contraction leaving only incoming variables open trivializes the tensor. We have motivated such constraints by conditional distributions, see Def. 22, and referred to Markov Networks (see Def. 27) satisfying these by Bayesian Networks (see Def. 31). To support our findings therein, we now discuss in more detail the connection between directed hypergraphs and directed tensors.

**Definition 79** (Directed Hypergraph). A directed hyperedge is a hyperedge, which node set is split into disjoint sets of incoming and outgoing nodes. We say a hypercore  $\tau^e$  decorating a directed hyperedge respects the direction, when it is a conditional probability tensor with respect to the direction of the hyperedge. The hypergraph is acyclic, when there is no nonempty cycle of node tuples  $(v_1, v_2)$ , such that  $v_1$  is an incoming node and  $v_2$  an outgoing node of the same hyperedge.

There can be multiple ways to direct a tensor, with an extreme example being Diracs Delta Tensors to be introduced in the next example. More general examples are basis encodings of invertible functions.

**Example 19** (Dirac Delta Tensors). Given a set of variables  $X_{[d]} = X_0, \dots, X_{d-1}$  with identical dimension  $m$ , Diracs Delta Tensor is the element

$$\delta^{[d], m} [X_{[d]}] \in \bigotimes_{k \in [d]} \mathbb{R}^m$$

with coordinates

$$\delta^{[d], m} [X_{[d]} = x_{[d]}] = \begin{cases} 1 & \text{if } x_0 = \dots = x_{d-1} \\ 0 & \text{else} \end{cases}. \quad (42)$$

The contractions with respect to subsets  $\tilde{\mathcal{V}} \subset [d]$  are

$$\langle \delta^{[d], m} \rangle [X_{\tilde{\mathcal{V}}}] = \begin{cases} m & \text{if } \tilde{\mathcal{V}} = \emptyset \\ \mathbb{I} [X_{\tilde{\mathcal{V}}}] & \text{if } |\tilde{\mathcal{V}}| = 1 \\ \delta^{\tilde{\mathcal{V}}, m} [X_{\tilde{\mathcal{V}}}] & \text{else} \end{cases}. \quad (43)$$

Thus are directed for any orientation of the respective edge with exactly one incoming variable.

We can use Diracs Delta Tensors to represent any contraction of a tensor network on a hypergraph by a tensor network on a graph, as we show next.

**Lemma 27.** Let  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$  be a hypergraph and  $\tau^{\mathcal{G}}$  a tensor network on  $\mathcal{G}$ . We build a graph  $\tilde{\mathcal{G}} = (\tilde{\mathcal{V}}, \tilde{\mathcal{E}} \cup \Delta^{\mathcal{G}})$  and a tensor network  $\tau^{\tilde{\mathcal{G}}}$  by

- Recolored Edges  $\tilde{\mathcal{E}} = \{\tilde{e} : e \in \mathcal{E}\}$  where  $\tilde{e} = \{v^e : v \in e\}$ , which decoration tensor  $\tau^{\tilde{e}}$  has same coordinates as  $\tau^e$
- Nodes  $\tilde{\mathcal{V}} = \bigcup_{e \in \mathcal{E}} \tilde{e}$
- Delta Edges  $\Delta^{\mathcal{G}} = \{\{v\} \cup \{v^e : e \ni v\} : v \in \mathcal{V}\}$  each of which decorated by a delta tensor  $\delta^{\{v^e : e \ni v\}}$

Then we have

$$\langle \tau^{\mathcal{G}} \rangle [X_{\mathcal{V}}] = \langle \tau^{\tilde{\mathcal{G}}} \rangle [X_{\mathcal{V}}].$$

*Proof.* For any  $x_{\mathcal{V}}$  we have

$$\begin{aligned} \langle \tau^{\tilde{G}} \rangle [X_{\mathcal{V}} = x_{\mathcal{V}}] &= \left\langle \{ \tau^{\tilde{e}} [X_{\{v^e: v \in e\}}] : e \in \mathcal{E} \} \cup \{ \delta^{\{v\} \cup \{v^e: e \ni v\}} [X_{\{v^e: e \ni v\}}, X_v = x_v] : v \in \mathcal{V} \} \right\rangle [\emptyset] \\ &= \langle \{ \tau^{\tilde{e}} [X_{\{v^e: v \in e\}} = x_{\{v: v \in e\}}] : e \in \mathcal{E} \} \rangle [\emptyset] \\ &= \langle \tau^{\mathcal{G}} \rangle [X_{\mathcal{V}} = x_{\mathcal{V}}] , \end{aligned}$$

which establishes the claim.  $\square$

### 16.3.1 normalization

Normed tensors (see Def. 13) are directed and directed tensors invariant under normalization wrt their incoming and outgoing variable, as we show next.

**Theorem 100.** *For any tensor network  $\tau^{\mathcal{G}}$  on variables  $\mathcal{V}$  that can be normed with respect to  $\mathcal{V}^{\text{in}}$  and  $\mathcal{V}^{\text{out}}$ , the normalization is directed with  $\mathcal{V}^{\text{in}}$  incoming and  $\mathcal{V}^{\text{out}}$  outgoing.*

*Proof.* We have for any incoming state  $x_{\mathcal{V}^{\text{in}}} \in \times_{v \in \mathcal{V}^{\text{in}}} m_v$  that

$$\langle \langle \tau^{\mathcal{G}} \rangle [\mathcal{V}^{\text{in}} | \mathcal{V}^{\text{out}}], \epsilon_{x_{\mathcal{V}^{\text{in}}}} \rangle [\emptyset] = \frac{\langle \tau^{\mathcal{G}} \cup \{ \epsilon_{x_{\mathcal{V}^{\text{in}}}} \} \rangle [\emptyset]}{\langle \tau^{\mathcal{G}} \cup \{ \epsilon_{x_{\mathcal{V}^{\text{in}}}} \} \rangle [\emptyset]} .$$

By Def. 12,  $\langle \tau^{\mathcal{G}} \rangle [\mathcal{V}^{\text{out}} | \mathcal{V}^{\text{in}}]$  is thus directed.  $\square$

The normalization operation coincides in cases of non-negative tensors with the conditioning of a Markov Network representing a probability distribution.

### 16.3.2 normalization Equations

normalization equations capture certain properties of normalizations of tensors. We first show that any normable tensor is the contraction of its normalization and an accompanying contraction, which generalizes the Bayes The. 4 towards more generic normable tensors.

**Theorem 101** (normalization as a Contraction Equation). *For any on  $\mathcal{V}^{\text{in}}$  normable tensor  $\tau [X_{\mathcal{V}}]$ , where  $\mathcal{V}^{\text{in}} \cup \mathcal{V}^{\text{out}} = \mathcal{V}$ , we have*

$$\langle \tau \rangle [X_{\mathcal{V}}] = \langle \langle \tau \rangle [X_{\mathcal{V}^{\text{out}}} | X_{\mathcal{V}^{\text{in}}}], \langle \tau \rangle [X_{\mathcal{V}^{\text{in}}}] \rangle [X_{\mathcal{V}}] .$$

*Proof.* Let us choose indices  $x_{\mathcal{V}^{\text{in}}}$  and  $x_{\mathcal{V}^{\text{out}}}$ . We have that

$$\langle \tau \rangle [X_{\mathcal{V}^{\text{in}}} = x_{\mathcal{V}^{\text{in}}} | X_{\mathcal{V}^{\text{out}}} = x_{\mathcal{V}^{\text{out}}}] = \frac{\langle \tau \rangle [X_{\mathcal{V}^{\text{in}}} = x_{\mathcal{V}^{\text{in}}}, X_{\mathcal{V}^{\text{out}}} = x_{\mathcal{V}^{\text{out}}}]}{\langle \tau \rangle [X_{\mathcal{V}^{\text{in}}} = x_{\mathcal{V}^{\text{in}}}]}$$

and therefor

$$\langle \tau \rangle [X_{\mathcal{V}^{\text{in}}} = x_{\mathcal{V}^{\text{in}}}, X_{\mathcal{V}^{\text{out}}} = x_{\mathcal{V}^{\text{out}}}] = \langle \tau \rangle [X_{\mathcal{V}^{\text{in}}} = x_{\mathcal{V}^{\text{in}}} | X_{\mathcal{V}^{\text{out}}} = x_{\mathcal{V}^{\text{out}}}] \cdot \langle \tau \rangle [X_{\mathcal{V}^{\text{in}}} = x_{\mathcal{V}^{\text{in}}}]$$

Since the equation holds for arbitrary indices, the claim is established.  $\square$

Based on this property, we now show a generic decomposition scheme of tensors, which generalizes the chain rule of The. 5.

**Theorem 102** (Generic Chain Rule). *For any Tensor  $\tau [X_{\mathcal{V}}]$  and any total order  $\prec$  on the nodes  $\mathcal{V}$  we have*

$$\tau [X_{\mathcal{V}}] = \langle \{ \langle \tau \rangle [X_v | X_{\{\tilde{v}: \tilde{v} \prec v, \tilde{v} \neq v\}}] : v \in \mathcal{V} \} \rangle [X_{\mathcal{V}}] ,$$

*provided that the normalizations exist.*

*Proof.* We apply The. 101 on the tensor

$$\langle \tau \rangle [X_v, X_{\{\tilde{v}: v \prec \tilde{v}, \tilde{v} \neq v\}} | X_{\{\tilde{v}: \tilde{v} \prec v, \tilde{v} \neq v\}}] = x_{\{\tilde{v}: \tilde{v} \prec v, \tilde{v} \neq v\}} ,$$

where  $v \in \mathcal{V}$  and  $x_{\mathcal{V}}$  are chosen arbitrarily. For any  $v \in \mathcal{V}$  we get

$$\langle \tau \rangle [X_v, X_{\{\tilde{v}: v \prec \tilde{v}, \tilde{v} \neq v\}} | X_{\{\tilde{v}: \tilde{v} \prec v, \tilde{v} \neq v\}}] = \langle \langle \tau \rangle [X_{\{\tilde{v}: v \prec \tilde{v}, \tilde{v} \neq v\}} | X_v, X_{\{\tilde{v}: \tilde{v} \prec v, \tilde{v} \neq v\}}], \langle \tau \rangle [X_v | X_{\{\tilde{v}: \tilde{v} \prec v, \tilde{v} \neq v\}}] \rangle [X_{\mathcal{V}}] .$$

Applying this equation iteratively and making use of the commutation of contractions we get for any  $v \in \mathcal{V}$

$$\langle \tau \rangle [X_v, X_{\{\tilde{v}: v \prec \tilde{v}, \tilde{v} \neq v\}} | X_{\{\tilde{v}: \tilde{v} \prec v, \tilde{v} \neq v\}}] = \langle \langle \tau \rangle [X_{\tilde{v}} | X_{\{\tilde{v}: \tilde{v} \prec \tilde{v}, \tilde{v} \neq \tilde{v}\}}] : v \prec \tilde{v} \rangle [X_{\mathcal{V}}] .$$

With the maximal node  $v$ , that is the  $v$ , such that no  $\tilde{v} \in \mathcal{V}$  with  $v \prec \tilde{v}$  and  $v \neq \tilde{v}$  exists, this is the claim.  $\square$

### 16.3.3 Contraction of Directed Tensors

Let us now investigate, which contractions inherit the directionality of the tensors.

**Theorem 103.** *Let  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$  be a directed acyclic hypergraph, such that each node  $v \in e$  appears at most in one hyperedge as an outgoing variable and denote  $\mathcal{V}^{\text{in}}$  as those nodes, which do not appear as outgoing variables. For any tensor network  $\tau^{\mathcal{G}}$  respecting the direction of  $\mathcal{G}$  we have that*

$$\langle \tau^{\mathcal{G}} \rangle [X_{\mathcal{V}^{\text{in}}}] = \mathbb{I} [X_{\mathcal{V}^{\text{in}}}] ,$$

that is  $\langle \tau^{\mathcal{G}} \rangle [X_{\mathcal{V}}]$  is a directed tensor with  $\mathcal{V}^{\text{in}}$  incoming and  $\mathcal{V}/\mathcal{V}^{\text{in}}$  outgoing.

*Proof.* We show the theorem only for the case of hypergraphs, where variables are appearing at most in two hyperedges. If a hypergraph fails to satisfy this assumption, we apply Lem. 27 and add delta tensors copying the variables, which are appearing in multiple tensors, and arrive at a tensor network with nodes appearing in at most two hyperedges.

We show the theorem over induction on the number  $n$  of cores.

$n = 1$ : The claim holds trivially, when  $\tau^{\mathcal{G}}$  consists of a single core.

$n \rightarrow n - 1$ : Let us assume, that the claim holds for graphs with  $n - 1$  hyperedges and let  $\tau^{\mathcal{G}}$  be a tensor network with  $n$  hyperedges. Since the hypergraph is acyclic, we find an edge  $e \in \mathcal{E}$  such that all outgoing nodes of  $e$  are not appearing as an incoming node in any edge. We then apply Theorem 131 and get

$$\begin{aligned} \langle \tau^{\mathcal{G}} \rangle [X_{\mathcal{V}^{\text{in}}}] &= \left\langle \tau^{(\mathcal{V}, \mathcal{E}/\{e\})} \cup \{\tau^e [X_{e^{\text{in}}}, X_{e^{\text{out}}}] \} \right\rangle [X_{\mathcal{V}^{\text{in}}}] \\ &= \left\langle \tau^{(\mathcal{V}, \mathcal{E}/\{e\})} \cup \{\langle \tau^e \rangle [X_{e^{\text{in}}}] \} \right\rangle [X_{\mathcal{V}^{\text{in}}}] \\ &= \left\langle \tau^{(\mathcal{V}, \mathcal{E}/\{e\})} \cup \{\mathbb{I} [X_{e^{\text{in}}}] \} \right\rangle [X_{\mathcal{V}^{\text{in}}}] \\ &= \left\langle \tau^{(\mathcal{V}, \mathcal{E}/\{e\})} \right\rangle [X_{\mathcal{V}^{\text{in}}}] . \end{aligned}$$

We then notice that the hypergraph  $(\mathcal{V}, \mathcal{E}/\{e\})$  has  $n - 1$  hyperedges and each node appears at most once as an incoming and at most once as an outgoing node. Thus, we apply the assumption of the induction and get

$$\langle \tau^{\mathcal{G}} \rangle [X_{\mathcal{V}^{\text{in}}}] = \left\langle \tau^{(\mathcal{V}, \mathcal{E}/\{e\})} \right\rangle [X_{\mathcal{V}^{\text{in}}}] = \mathbb{I} [X_{\mathcal{V}^{\text{in}}}] .$$

□

### 16.4 Proof of Hammersley-Clifford Theorem

Let us now proof the Hammersley-Clifford theorem formulated in Chapter 5 as The. 12. Different to the original statement (see Clifford and Hammersley (1971)), we here proof the analogous statement for hypergraphs, where we have to demand the property of clique-capturing defined in Def. 29. We start with showing the following Lemmata to be exploited in the proof.

**Lemma 28.** *Let  $\tau [X_{\mathcal{V}}]$  be a positive tensor and  $y_{\mathcal{V}}$  an arbitrary index. Then we have*

$$\tau [X_{\mathcal{V}}] = \left\langle \left( \langle \tau \rangle [X_{\mathcal{V}/\bar{\mathcal{V}}}, X_{\bar{\mathcal{V}}} = y_{\bar{\mathcal{V}}}] \right)^{(-1)^{|\bar{\mathcal{V}}| - |\mathcal{V}|}} : \bar{\mathcal{V}} \subset \tilde{\mathcal{V}} \subset \mathcal{V} \right\rangle [X_{\mathcal{V}}] ,$$

where the exponentiation is performed coordinatewise and positivity of  $\tau$  ensures the well-definedness.

*Proof.* It suffices to show, that for an arbitrary index  $x_{\mathcal{V}}$  be an arbitrary index we have

$$\tau [X_{\mathcal{V}} = x_{\mathcal{V}}] = \prod_{\tilde{\mathcal{V}} \subset \mathcal{V}} \prod_{\bar{\mathcal{V}} \subset \tilde{\mathcal{V}}} \left( \langle \tau \rangle [X_{\mathcal{V}/\bar{\mathcal{V}}} = x_{\mathcal{V}/\bar{\mathcal{V}}}, X_{\bar{\mathcal{V}}} = y_{\bar{\mathcal{V}}}] \right)^{(-1)^{|\bar{\mathcal{V}}| - |\mathcal{V}|}} .$$

We do this by applying a logarithm on the right hand side and grouping the terms by  $\bar{\mathcal{V}}$  as

$$\begin{aligned} & \ln \left[ \prod_{\bar{\mathcal{V}} \subset \mathcal{V}} \prod_{\bar{\mathcal{V}} \subset \bar{\mathcal{V}}} \langle \tau \rangle [X_{\mathcal{V}/\bar{\mathcal{V}}} = x_{\mathcal{V}/\bar{\mathcal{V}}}, X_{\bar{\mathcal{V}}} = y_{\bar{\mathcal{V}}}] \right]^{(-1)^{|\bar{\mathcal{V}}| - |\mathcal{V}|}} \\ &= \sum_{\bar{\mathcal{V}} \subset \mathcal{V}} \ln [\langle \tau \rangle [X_{\mathcal{V}/\bar{\mathcal{V}}} = x_{\mathcal{V}/\bar{\mathcal{V}}}, X_{\bar{\mathcal{V}}} = y_{\bar{\mathcal{V}}}] \left( \sum_{\bar{\mathcal{V}} \subset \mathcal{V} : \bar{\mathcal{V}} \subset \bar{\mathcal{V}}} (-1)^{|\bar{\mathcal{V}}| - |\mathcal{V}|} \right) \\ &= \sum_{\bar{\mathcal{V}} \subset \mathcal{V}} \ln [\langle \tau \rangle [X_{\mathcal{V}/\bar{\mathcal{V}}} = x_{\mathcal{V}/\bar{\mathcal{V}}}, X_{\bar{\mathcal{V}}} = y_{\bar{\mathcal{V}}}] \left( \sum_{i \in [|\mathcal{V}| - |\bar{\mathcal{V}}|]} (-1)^i \binom{|\mathcal{V}| - |\bar{\mathcal{V}}|}{i} \right) \end{aligned}$$

Now, by the generic binomial theorem we have that for  $n \in \mathbb{N}, n \neq 0$

$$0 = (1 - 1)^n = \sum_{i \in [n]} (-1)^i \binom{n}{i}.$$

Therefore, the summands for  $\bar{\mathcal{V}} \neq \mathcal{V}$  vanish and we have

$$\begin{aligned} & \ln \left[ \prod_{\bar{\mathcal{V}} \subset \mathcal{V}} \prod_{\bar{\mathcal{V}} \subset \bar{\mathcal{V}}} (\langle \tau \rangle [X_{\mathcal{V}/\bar{\mathcal{V}}} = x_{\mathcal{V}/\bar{\mathcal{V}}}, X_{\bar{\mathcal{V}}} = y_{\bar{\mathcal{V}}}] \right]^{(-1)^{|\bar{\mathcal{V}}| - |\mathcal{V}|}} \\ &= \ln [\tau [X_{\mathcal{V}} = x_{\mathcal{V}}]] \left( \sum_{i \in [0]} (-1)^i \binom{0}{i} \right) \\ &= \ln [\tau [X_{\mathcal{V}} = x_{\mathcal{V}}]]. \end{aligned}$$

Applying the exponential function on both sides establishes the claim.  $\square$

**Lemma 29.** *Let  $\tau$  be a positive tensor,  $\tilde{\mathcal{V}} \subset \mathcal{V}$  and arbitrary subset and  $x_{\tilde{\mathcal{V}}}$  an arbitrary index. When there are  $A, B \in \tilde{\mathcal{V}}$ , such that*

$$\langle \tau \rangle [X_{A,B} | X_{\mathcal{V}/\{A,B\}}] = \langle \langle \tau \rangle [X_A | X_{\mathcal{V}/\{A,B\}}], \langle \tau \rangle [X_B | X_{\mathcal{V}/\{A,B\}}] \rangle [X_{\tilde{\mathcal{V}}}]$$

then

$$\prod_{\bar{\mathcal{V}} \subset \tilde{\mathcal{V}}} (\langle \tau \rangle [X_{\mathcal{V}/\bar{\mathcal{V}}} = x_{\mathcal{V}/\bar{\mathcal{V}}}, X_{\bar{\mathcal{V}}} = y_{\bar{\mathcal{V}}}] \right)^{(-1)^{|\bar{\mathcal{V}}| - |\tilde{\mathcal{V}}|}} = 1.$$

*Proof.* We abbreviate

$$Z_{\bar{\mathcal{V}}} = \langle \tau \rangle [X_{\mathcal{V}/\bar{\mathcal{V}}} = x_{\mathcal{V}/\bar{\mathcal{V}}}, X_{\bar{\mathcal{V}}} = y_{\bar{\mathcal{V}}}] .$$

By reorganizing the sum over  $\bar{\mathcal{V}} \subset \tilde{\mathcal{V}}$  into  $\bar{\mathcal{V}} \subset \tilde{\mathcal{V}}/A \cup B$  we have

$$\prod_{\bar{\mathcal{V}} \subset \tilde{\mathcal{V}}} (Z_{\bar{\mathcal{V}}})^{(-1)^{|\bar{\mathcal{V}}| - |\tilde{\mathcal{V}}|}} = \prod_{\bar{\mathcal{V}} \subset \tilde{\mathcal{V}}/\{A,B\}} \left( \frac{Z_{\bar{\mathcal{V}}} \cdot Z_{\bar{\mathcal{V}} \cup \{A,B\}}}{Z_{\bar{\mathcal{V}} \cup \{A\}} \cdot Z_{\bar{\mathcal{V}} \cup \{B\}}} \right)^{(-1)^{|\bar{\mathcal{V}}| - |\tilde{\mathcal{V}}|}}. \quad (44)$$

From the independence assumption it follows that for any index  $x$

$$\begin{aligned} & \langle \tau \rangle [X_A = x_A | X_{\mathcal{V}/\bar{\mathcal{V}} \cup \{A,B\}} = x_{\mathcal{V}/\bar{\mathcal{V}} \cup \{A,B\}}, X_{\bar{\mathcal{V}}} = y_{\bar{\mathcal{V}}}, X_B = x_B] \\ &= \langle \tau \rangle [X_A = x_A | X_{\mathcal{V}/\bar{\mathcal{V}} \cup \{A,B\}} = x_{\mathcal{V}/\bar{\mathcal{V}} \cup \{A,B\}}, X_{\bar{\mathcal{V}}} = y_{\bar{\mathcal{V}}}] \\ &= \langle \tau \rangle [X_A = x_A | X_{\mathcal{V}/\bar{\mathcal{V}} \cup \{A,B\}} = x_{\mathcal{V}/\bar{\mathcal{V}} \cup \{A,B\}}, X_{\bar{\mathcal{V}}} = y_{\bar{\mathcal{V}}}, X_B = y_B] \end{aligned}$$

Applying this in each squares bracket term of (44) we get

$$\begin{aligned} \frac{Z_{\bar{\mathcal{V}}}}{Z_{\bar{\mathcal{V}} \cup \{A\}}} &= \frac{\langle \tau \rangle [X_A = x_A | X_{\mathcal{V}/\bar{\mathcal{V}} \cup \{A,B\}} = x_{\mathcal{V}/\bar{\mathcal{V}} \cup \{A,B\}}, X_{\bar{\mathcal{V}}} = y_{\bar{\mathcal{V}}}, X_B = x_B]}{\langle \tau \rangle [X_A = y_A | X_{\mathcal{V}/\bar{\mathcal{V}} \cup \{A,B\}} = x_{\mathcal{V}/\bar{\mathcal{V}} \cup \{A,B\}}, X_{\bar{\mathcal{V}}} = y_{\bar{\mathcal{V}}}, X_B = x_B]} \\ &= \frac{\langle \tau \rangle [X_A = x_A | X_{\mathcal{V}/\bar{\mathcal{V}} \cup \{A,B\}} = x_{\mathcal{V}/\bar{\mathcal{V}} \cup \{A,B\}}, X_{\bar{\mathcal{V}}} = y_{\bar{\mathcal{V}}}, X_B = y_B]}{\langle \tau \rangle [X_A = y_A | X_{\mathcal{V}/\bar{\mathcal{V}} \cup \{A,B\}} = x_{\mathcal{V}/\bar{\mathcal{V}} \cup \{A,B\}}, X_{\bar{\mathcal{V}}} = y_{\bar{\mathcal{V}}}, X_B = y_B]} \\ &= \frac{Z_{\bar{\mathcal{V}} \cup \{B\}}}{Z_{\bar{\mathcal{V}} \cup \{A,B\}}}. \end{aligned}$$

Thus, each factor in (44) is trivial, which establishes the claim.  $\square$

We are finally ready to proof the Hammersley-Clifford The. 12 based on the Lemmata above.

*Proof of The. 12.* By Lem. 28 we have for any index  $x_{\mathcal{V}}$

$$\mathbb{P}[X_{\mathcal{V}} = x_{\mathcal{V}}] = \prod_{\tilde{\mathcal{V}} \subset \mathcal{V}} \prod_{\tilde{\mathcal{V}} \subset \tilde{\mathcal{V}}} (\mathbb{P}[X_{\tilde{\mathcal{V}}} = x_{\tilde{\mathcal{V}}}, X_{\mathcal{V}/\tilde{\mathcal{V}}} = y_{\mathcal{V}/\tilde{\mathcal{V}}}] )^{(-1)^{|\tilde{\mathcal{V}}| - |\mathcal{V}|}}$$

For any subset  $\tilde{\mathcal{V}} \subset \mathcal{V}$ , which is not contained in a hyperedge, we find  $A, B \in \tilde{\mathcal{V}}$  such that  $X_A$  is independent on  $X_B$  conditioned on  $X_{\tilde{\mathcal{V}}/\{A,B\}}$ . If no such nodes  $A, B \in \tilde{\mathcal{V}}$  exists,  $\tilde{\mathcal{V}}$  would be contained in a hyperedge, since the hypergraph is assumed to be clique-capturing. By Lem. 29 we then have

$$\prod_{\tilde{\mathcal{V}} \subset \tilde{\mathcal{V}}} (\mathbb{P}[X_{\tilde{\mathcal{V}}} = x_{\tilde{\mathcal{V}}}, X_{\mathcal{V}/\tilde{\mathcal{V}}} = y_{\mathcal{V}/\tilde{\mathcal{V}}}] )^{(-1)^{|\tilde{\mathcal{V}}| - |\mathcal{V}|}} = 1.$$

We label by a function

$$\alpha : \{\tilde{\mathcal{V}} : \exists e \in \mathcal{E} : \tilde{\mathcal{V}} \subset e\} \rightarrow \mathcal{E}$$

the remaining node subsets by a hyperedge containing the subset. We build the tensor

$$\tau^e[X_e] = \prod_{\tilde{\mathcal{V}} : \alpha(\tilde{\mathcal{V}}) = e} \prod_{\tilde{\mathcal{V}} \subset \tilde{\mathcal{V}}} (\mathbb{P}[X_{\tilde{\mathcal{V}}} = x_{\tilde{\mathcal{V}}}, X_{\mathcal{V}/\tilde{\mathcal{V}}} = y_{\mathcal{V}/\tilde{\mathcal{V}}}] )^{(-1)^{|\tilde{\mathcal{V}}| - |\mathcal{V}|}}.$$

and get, that

$$\begin{aligned} \mathbb{P}[X_{\mathcal{V}}] &= \langle \{\tau^e[X_e] : e \in \mathcal{E}\} \rangle [X_{\mathcal{V}}] \\ &= \langle \{\tau^e[X_e] : e \in \mathcal{E}\} \rangle [X_{\mathcal{V}}|\emptyset]. \end{aligned}$$

We have thus constructed a Markov Network with trivial partition function, which contraction coincides with the probability distribution.  $\square$

## 16.5 Differentiation of Contraction

The structured mean field approaches discussed in Chapter 6 used differentiations of the parametrized tensor networks. Let us now develop in more detail, how the contraction of tensor networks with variable cores is differentiated. We capture in additional variables  $Y$  selecting the coordinates of a tensor, which are varied in a differentiation.

**Lemma 30.** *For any tensor network  $\tau^{\mathcal{G}}$  with positive  $\tau^e$  we have*

$$\begin{aligned} \frac{\partial}{\partial \tau^e[Y_e]} \langle \tau^{\mathcal{G}} \rangle [X_{\mathcal{V}}|\emptyset] &= \left\langle \delta[Y_e, X_e], \frac{\langle \tau^{\mathcal{G}} \rangle [X_e]}{\tau^e[X_e]}, \langle \tau^{\mathcal{G}} \rangle [X_{\mathcal{V}/e}|X_e] \right\rangle [Y_e, X_{\mathcal{V}}] \\ &\quad - \langle \tau^{\mathcal{G}} \rangle [X_{\mathcal{V}}|\emptyset] \otimes \left\langle \frac{\langle \tau^{\mathcal{G}} \rangle [Y_e]}{\tau^e[Y_e]} \right\rangle [Y_e]. \end{aligned}$$

*Proof.* By multilinearity of tensor network contractions we have

$$\frac{\partial}{\partial \tau^e[Y_e]} \langle \tau^{\mathcal{G}} \rangle [X_{\mathcal{V}}] = \langle \{\delta[Y_e, X_e]\} \cup \{\tau^{\tilde{e}}[X_{\tilde{e}}] : \tilde{e} \neq e\} \rangle [Y_e, X_{\mathcal{V}}]$$

and thus

$$\frac{\partial}{\partial \tau^e[Y_e]} \langle \tau^{\mathcal{G}} \rangle [\emptyset] = \langle \{\delta[Y_e, X_e]\} \cup \{\tau^{\tilde{e}}[X_{\tilde{e}}] : \tilde{e} \neq e\} \rangle [Y_e].$$



Using both we get

$$\begin{aligned}
 \frac{\partial}{\partial \tau^e [Y_e]} \langle \tau^G \rangle [X_{\mathcal{V}} | \emptyset] &= \frac{\partial}{\partial \tau^e [Y_e]} \frac{\langle \tau^G \rangle [X_{\mathcal{V}}]}{\langle \tau^G \rangle [\emptyset]} \\
 &= \frac{\frac{\partial}{\partial \tau^e [Y_e]} \langle \tau^G \rangle [X_{\mathcal{V}}]}{\langle \tau^G \rangle [\emptyset]} - \frac{\langle \tau^G \rangle [X_{\mathcal{V}}] \frac{\partial}{\partial \tau^e [Y_e]} \langle \tau^G \rangle [\emptyset]}{(\langle \tau^G \rangle [\emptyset])^2} \\
 &= \frac{\langle \{\delta [Y_e, X_e]\} \cup \{\tau^{\tilde{e}} [X_{\tilde{e}}] : \tilde{e} \neq e\} \rangle [Y_e, X_{\mathcal{V}}]}{\langle \tau^G \rangle [\emptyset]} \\
 &\quad - \langle \tau^G \rangle [X_{\mathcal{V}} | \emptyset] \cdot \frac{\langle \{\delta [Y_e, X_e]\} \cup \{\tau^{\tilde{e}} [X_{\tilde{e}}] : \tilde{e} \neq e\} \rangle [Y_e]}{\langle \tau^G \rangle [\emptyset]}
 \end{aligned} \tag{45}$$

For the first term we get with a normalization equation (see Theorem 101) that

$$\begin{aligned}
 \frac{\langle \{\delta [Y_e, X_e]\} \cup \{\tau^{\tilde{e}} [X_{\tilde{e}}] : \tilde{e} \neq e\} \rangle [Y_e, X_{\mathcal{V}}]}{\langle \tau^G \rangle [\emptyset]} &= \frac{\langle \{\delta [Y_e, X_e]\} \cup \{\tau^{\tilde{e}} [X_{\tilde{e}}] : \tilde{e} \in \mathcal{E}\} \rangle [Y_e, X_{\mathcal{V}}]}{\tau^e [X_e] \cdot \langle \tau^G \rangle [\emptyset]} \\
 &= \frac{\langle \delta [Y_e, X_e], \langle \tau^G \rangle [X_{\mathcal{V}} | \emptyset] \rangle [Y_e, X_{\mathcal{V}}]}{\tau^e [X_e]} \\
 &= \frac{\langle \delta [Y_e, X_e], \langle \tau^G \rangle [X_e | \emptyset], \langle \tau^G \rangle [X_{\mathcal{V}/e} | X_e] \rangle [Y_e, X_{\mathcal{V}}]}{\tau^e [X_e]}.
 \end{aligned}$$

Analogously, we have

$$\frac{\langle \{\delta [Y_e, X_e]\} \cup \{\tau^{\tilde{e}} [X_{\tilde{e}}] : \tilde{e} \neq e\} \rangle [Y_e]}{\langle \tau^G \rangle [\emptyset]} = \frac{\langle \delta [Y_e, X_e], \langle \tau^G \rangle [X_e | \emptyset] \rangle [Y_e]}{\tau^e [X_e]}.$$

With (45), we arrive at the claim

$$\begin{aligned}
 \frac{\partial}{\partial \tau^e [Y_e]} \langle \tau^G \rangle [X_{\mathcal{V}} | \emptyset] &= \left\langle \delta [Y_e, X_e], \frac{\langle \tau^G \rangle [X_e]}{\tau^e [X_e]}, \langle \tau^G \rangle [X_{\mathcal{V}/e} | X_e] \right\rangle [Y_e, X_{\mathcal{V}}] \\
 &\quad - \langle \tau^G \rangle [X_{\mathcal{V}} | \emptyset] \otimes \left\langle \frac{\langle \tau^G \rangle [Y_e]}{\tau^e [Y_e]} \right\rangle [Y_e].
 \end{aligned}$$

□

**Lemma 31.** For any function  $q(\tau^e)[X_{\mathcal{V}}]$  we have

$$\begin{aligned}
 \frac{\partial}{\partial \tau^e [Y_e]} \langle \langle \tau^G \rangle [X_{\mathcal{V}} | \emptyset], q(\tau^e)[X_{\mathcal{V}}] \rangle [\emptyset] &= \frac{\langle \tau^G \rangle [X_e = x_e | \emptyset]}{\tau^e [X_e = x_e]} \left( \langle \langle \tau^G \rangle [X_{\mathcal{V}/e} | X_e = x_e], q(\tau^e)[X_{\mathcal{V}}, Y_e] \rangle [\emptyset] \right. \\
 &\quad \left. - \langle \langle \tau^G \rangle [X_{\mathcal{V}} | \emptyset], q(\tau^e)[X_{\mathcal{V}}] \rangle [\emptyset] \right) \\
 &\quad + \left\langle \langle \tau^G \rangle [X_{\mathcal{V}} | \emptyset], \frac{\partial q(\tau^e)[X_{\mathcal{V}}]}{\partial \tau^e [Y_e]} \right\rangle [\emptyset]
 \end{aligned}$$

*Proof.* By product rule of differentiation we have

$$\begin{aligned}
 \frac{\partial}{\partial \tau^e [X_e = x_e]} \langle \langle \tau^G \rangle [X_{\mathcal{V}} | \emptyset], q(\tau^e)[X_{\mathcal{V}}] \rangle [\emptyset] &= \left\langle \frac{\partial}{\partial \tau^e [X_e = x_e]} \langle \tau^G \rangle [X_{\mathcal{V}} | \emptyset], q(\tau^e)[X_{\mathcal{V}}] \right\rangle [\emptyset] \\
 &\quad + \left\langle \langle \tau^G \rangle [X_{\mathcal{V}} | \emptyset], \frac{\partial}{\partial \tau^e [X_e = x_e]} q(\tau^e)[X_{\mathcal{V}}] \right\rangle [\emptyset].
 \end{aligned}$$

The claim now follows with the application of Lem. 30 on the first term. □

## 16.6 Discussion

Representations of linear maps is the typical application of tensors, reason for refering to tensor networks as multilinear algebra.

## 17 Basis Calculus

Basis Calculus stores informations in the selection of basis elements, while coordinate calculus uses the coordinates to each index for storage. While coordinate calculus is more expressive, basis calculus can be exploited in sparse representations of composed functions.

### 17.1 Basis Encoding of Subsets

Based on the concept of one-hot encodings of states we in this chapter develop the construction of encodings to sets, relations and functions. We start with the definition of subset encodings, which represent set memberships in their boolean coordinates.

**Definition 80** (Basis encoding of subsets). *We say that an arbitrary set  $\mathcal{U}$  is enumerated by an enumeration variable  $O_{\mathcal{U}}$  taking values in  $[r_{\mathcal{U}}]$ , when  $r_{\mathcal{U}} = |\mathcal{U}|$  and there is a bijective index interpretation function*

$$I : [r_{\mathcal{U}}] \rightarrow \mathcal{U}.$$

*Given an set  $\mathcal{U}$  enumerated by the variable  $O_{\mathcal{U}}$ , any subset  $\mathcal{V} \subset \mathcal{U}$  is encoded by the tensor  $\epsilon_{\mathcal{V}}[O]$  defined for  $o \in [|\mathcal{U}|]$  as*

$$\epsilon_{\mathcal{V}}[O = o] = \begin{cases} 1 & \text{if } I(o) \in \mathcal{V} \\ 0 & \text{else} \end{cases}.$$

In a one-hot basis decomposition we have

$$\epsilon_{\mathcal{V}}[O] := \sum_{o \in [|\mathcal{U}|] : I(o) \in \mathcal{V}} \epsilon_o[O].$$

The inclusion of subsets is represented by the partial ordering of tensors. Let us first define this property for arbitrary tensors.

**Definition 81** (Partial ordering of tensors). *We say that two tensors  $f[X_{[d]}]$  and  $h[X_{[d]}]$  attached with the same variables are partially ordered, denoted by*

$$f \prec h,$$

*if for all  $x_{[d]} \in \times_{k \in [d]} [m_k]$*

$$f[X_{[d]} = x_{[d]}] \leq h[x_{[d]}].$$

For boolean tensors, the partially ordering is equal to a subset relation of the coordinates with value 1, as we show next.

**Theorem 104.** *Let  $\mathcal{U}$  be an arbitrary set enumerated by the variable  $O$  and index interpretation function  $I$ . For two subsets  $\mathcal{U}^0, \mathcal{U}^1$  of  $\mathcal{U}$  we have*

$$\mathcal{U}^0 \subset \mathcal{U}^1$$

*if and only if*

$$\epsilon_{\mathcal{U}^0}[O] \prec \epsilon_{\mathcal{U}^1}[O].$$

*Proof.* We have  $\mathcal{U}^0 \subset \mathcal{U}^1$  if and only if

$$\forall o \in [|\mathcal{U}|] (I(o) \in \mathcal{U}^0) \Rightarrow (I(o) \in \mathcal{U}^1),$$

which is equal to

$$\forall o \in [|\mathcal{U}|] (\epsilon_{\mathcal{U}^0}[O = o] = 1) \Rightarrow (\epsilon_{\mathcal{U}^1}[O = o] = 1).$$

Since subset encodings are boolean tensors, this is equivalent to

$$\epsilon_{\mathcal{U}^0}[O] \prec \epsilon_{\mathcal{U}^1}[O].$$

□

### 17.1.1 Binary Relations

Since relations are subsets of cartesian products between two sets, their encoding is a straightforward generalization of Def. 80.

**Definition 82** (Basis encoding of binary relations). *A relation between two finite sets  $\mathcal{U}^{\text{in}}$  and  $\mathcal{U}^{\text{out}}$  is a subset of their cartesian product*

$$\mathcal{R} \subset \mathcal{U}^{\text{in}} \times \mathcal{U}^{\text{out}}.$$

*Given an enumeration of  $\mathcal{U}^{\text{in}}$  and  $\mathcal{U}^{\text{out}}$  by the categorical variables  $O_{\text{in}}$  and  $O_{\text{out}}$  and interpretation maps  $I_{\text{in}}, I_{\text{out}}$ , we define the basis encoding of this subset as the tensor  $\epsilon_{\mathcal{R}}[O_{\text{in}}, O_{\text{out}}]$  with the coordinates*

$$\epsilon_{\mathcal{R}}[O_{\text{in}} = o_{\text{in}}, O_{\text{out}} = o_{\text{out}}] = \begin{cases} 1 & \text{if } (I_{\text{in}}(o_{\text{in}}), I_{\text{out}}(o_{\text{out}})) \in \mathcal{R} \\ 0 & \text{else} \end{cases}.$$

The basis encoding has a decomposition into one-hot encodings as

$$\epsilon_{\mathcal{R}}[O_{\text{in}}, O_{\text{out}}] = \sum_{o_{\text{in}}, o_{\text{out}} : (I_{\text{in}}(o_{\text{in}}), I_{\text{out}}(o_{\text{out}})) \in \mathcal{R}} \epsilon_{o_{\text{in}}}[O_{\text{in}}] \otimes \epsilon_{o_{\text{out}}}[O_{\text{out}}].$$

basis encodings have a matrix structure by the cartesian product, which can be further folded to tensors, when the sets itself are cartesian products. The basis encoding is a bijection between the relations of two sets and the boolean tensors with their enumeration variables.

### 17.1.2 Higher order relations

We can extend this contraction to relations of higher order, and arrive at encoding schemes usable for relational databases.

**Definition 83** (Basis encoding of  $d$ -ary relations). *Given sets  $\mathcal{U}^k$  for  $k \in [d]$ , a  $d$ -ary relation is a subset of a their cartesian product, that is*

$$\mathcal{R} \subset \bigtimes_{k \in [d]} \mathcal{U}^k.$$

*Given an enumeration of each set  $\mathcal{U}^k$  by a variable  $O_k$  and an interpretation map  $I_k$ , we define the basis encoding of the relation as the tensor  $\epsilon_{\mathcal{R}}[O_{[d]}]$  with coordinates*

$$\epsilon_{\mathcal{R}}[O_{[d]} = o_{[d]}] = \begin{cases} 1 & \text{if } (I_0(o_0), \dots, I_{d-1}(o_{d-1})) \in \mathcal{R} \\ 0 & \text{else} \end{cases}.$$

Let there be for  $k \in [d]$  sets  $\mathcal{U}^k$  of truth assignments to the  $k$ -th atom, which are all enumerated by [2]. A propositional formula then corresponds with a  $d$ -ary relation and we directly defined them in Def. 43 by their basis encoding.

**Theorem 105.** *The encoding of any  $d$ -ary relation*

$$\mathcal{R} = \{x_{[d]}^i : i \in [n]\} \subset \bigtimes_{k \in [d]} \{\text{False}, \text{True}\}$$

*where the objects in  $\{\text{False}, \text{True}\}$  are enumerated by  $X_k$  with the standard index interpretation function (see Sect. 7.1)*

$$I(\text{True}) = 1 \quad \text{and} \quad I(\text{False}) = 0,$$

*coincides with the propositional formula*

$$f[X_{[d]}] = \bigvee_{i \in [n]} Z_{x_{[d]}^i}^{\wedge}.$$

*Proof.* By definition, the encoding  $\epsilon_{\mathcal{R}}$  is decomposed as

$$\epsilon_{\mathcal{R}}[X_{[d]}] = \sum_{i \in [n]} \epsilon_{x_{[d]}^i}[X_{[d]}].$$

By The. 42 this is equal to

$$\tau[X_{[d]}] = \left( \bigvee_{x_{[d]} : \tau[X_{[d]} = x_{[d]}] = 1} Z_{\{k: x_k=0\}, \{k: x_k=0\}}^{\wedge} \right) [X_{[d]}] = \bigvee_{i \in [n]} Z_{x_{[d]}^i}^{\wedge}.$$

□

**Example 20** (Relational Databases). *Relational Databases can be encoded as tensors using the relation encoding scheme. Each column is thereby understood as an enumeration variable, whose values form the sets  $\mathcal{U}^k$ .*

Let us notice, that the dimensionality of the tensor space used for representing a relation is

$$\prod_{k \in [d]} |\mathcal{U}^k|$$

and therefore growing exponentially with the number of variables. Relations are however often sparse, in the sense that

$$|\mathcal{R}| \ll \prod_{k \in [d]} |\mathcal{U}^k|.$$

It is therefore often beneficially to choose sparse encoding schemes, for example by restricted CP formats (see Chapter 18) to represent  $\epsilon_{\mathcal{R}}$ .

## 17.2 Basis Encoding of Functions

Let us now restrict to relations, which have an expression by functions. We in this section then show, how contractions of their encodings can be exploited in function evaluation.

### 17.2.1 Basis encoding of Functions

**Definition 84** (Basis encoding of maps). *Any map*

$$q : \mathcal{U}^{\text{in}} \rightarrow \mathcal{U}^{\text{out}}$$

*can be represented by a relation*

$$\mathcal{R}^q := \{(x, q(x)) : x \in \mathcal{U}^{\text{in}}\} \subset \mathcal{U}^{\text{in}} \times \mathcal{U}^{\text{out}}.$$

*Given an enumeration of the sets by  $O_{\text{in}}$  and  $O_{\text{out}}$  we define the basis encoding of  $q$  as the tensor*

$$\beta^q [O_{\text{out}}, O_{\text{in}}] = \epsilon_{\mathcal{R}^q} [O_{\text{in}}, O_{\text{out}}].$$

**Remark 21** (Reduction to images). *When  $q$  maps into a set of infinite cardinality, we restrict  $\mathcal{U}^{\text{out}}$  to the image of  $q$  and enumerate the image by a variable  $O_q$ . This scheme is applied, when  $q$  is itself a tensor, i.e.  $\mathcal{U}^{\text{out}} = \mathbb{R}$ . While the variable  $O_q$  can in general be of the same cardinality as the domain set  $\mathcal{U}^{\text{in}}$ , it will be valued in  $[2]$  when considering boolean tensors.*

We notice, that any basis representation of a function is also a directed tensor with incoming variables to the domain and outgoing variables to the image. It furthermore holds, that the set of directed and boolean tensors is characterized by the basis encoding of functions. This is shown in the next theorem, by the claim that any boolean tensor which is directed is the basis representation of a function.

**Theorem 106.** *Let  $\mathcal{U}^{\text{in}}, \mathcal{U}^{\text{out}}$  be sets and  $\mathcal{R} \subset \mathcal{U}^{\text{in}} \times \mathcal{U}^{\text{out}}$  a relation. If and only if there exists a map  $q : \mathcal{U}^{\text{in}} \rightarrow \mathcal{U}^{\text{out}}$  such that  $\mathcal{R} = \mathcal{R}^q$ , the basis encoding  $\beta^q$  is a directed tensor with  $O_{\text{in}}$  incoming and  $O_{\text{out}}$  outgoing.*

*Proof.* " $\Rightarrow$ ": When  $q$  is a function, we have for any  $o_{\text{in}} \in [r_{\text{in}}]$

$$\sum_{o_{\text{out}} \in [r_{\text{out}}]} \beta^q [O_{\text{out}} = o_{\text{out}}, O_{\text{in}} = o_{\text{in}}] = \beta^q [O_{\text{out}} = I_{\text{out}}^{-1}(q(I_{\text{in}}(o_{\text{in}}))), O_{\text{in}} = o_{\text{in}}] = 1.$$

Thus,  $\beta^q [O_{\text{out}}, O_{\text{in}}]$  is a directed tensor with variables  $O_{\text{in}}$  incoming and  $O_{\text{out}}$  outgoing.

" $\Leftarrow$ ": Conversely let there be a relation  $\mathcal{R}$ , such that  $\beta^{\mathcal{R}}$  is directed. To this end, we observe that for any  $o_{\text{in}} \in [r_{\text{in}}]$  the tensor

$$\epsilon_{\mathcal{R}} [O_{\text{in}} = o_{\text{in}}, O_{\text{out}}]$$

is a boolean tensor with coordinate sum one and therefore a basis vector. It follows that the function  $q : \mathcal{U}^{\text{in}} \rightarrow \mathcal{U}^{\text{out}}$  defined for  $x \in \mathcal{U}^{\text{in}}$  as

$$q(x) = I_{\text{out}}(\epsilon^{-1}(\epsilon_{\mathcal{R}} [O_{\text{in}} = I_{\text{in}}(x), O_{\text{out}}]))$$

is well-defined. We then have by construction

$$\begin{aligned}\beta^q [O_{\text{out}}, O_{\text{in}}] &= \sum_{o_{\text{in}} \in [r_{\text{in}}]} \epsilon_{q(o_{\text{in}})} [O_{\text{out}}] \otimes \epsilon_{o_{\text{in}}} [O_{\text{in}}] \\ &= \sum_{o_{\text{in}} \in [r_{\text{in}}]} \epsilon_{\mathcal{R}} [O_{\text{in}} = o_{\text{in}}, O_{\text{out}}] \otimes \epsilon_{o_{\text{in}}} [O_{\text{in}}] \\ &= \epsilon_{\mathcal{R}} [O_{\text{out}}, O_{\text{in}}]\end{aligned}$$

and therefore by Def. 84  $\mathcal{R} = \mathcal{R}^q$ . □

We are specially interested in sets of states of a factored system, which amounts to the case in Def. 14. Those state sets have a decomposition into a cartesian product of  $d$  sets

$$\mathcal{U} = \bigtimes_{k \in [d]} [m_k].$$

The most obvious enumeration of the set  $\mathcal{U}$  is therefore by the collection of state variables  $\{X_k : k \in [d]\}$ . Functions between states of factored systems with  $d_{\text{in}}$  and  $d_{\text{out}}$  state variables can be represented by  $d_{\text{in}} + d_{\text{out}}$ -ary relations and Def. 84 has an obvious generalization to this case with multiple enumeration variables.

### 17.2.2 Function Evaluation

We now justify the nomenclature of basis calculus, by showing that contraction with basis elements produce the one-hot encoded function evaluation.

**Theorem 107** (Function evaluation in Basis Calculus). *Retrieving the value of the function  $q$  at a specific state is then the contraction of the tensor representation with the one-hot encoded state. For any  $u \in \mathcal{U}^{\text{in}}$  we have*

$$\epsilon_{I_{\text{out}}^{-1}(q(u))} [O_{\text{out}}] = \langle \beta^f [O_{\text{out}}, O_{\text{in}}], \epsilon_{I_{\text{in}}(u)} [O_{\text{in}}] \rangle [O_{\text{out}}].$$

Thus, we can retrieve the function evaluation by the inverse one-hot mapping as

$$q(u) = \epsilon^{-1}(\langle \beta^f [O_{\text{out}}, O_{\text{in}}], \epsilon_{I_{\text{in}}(u)} [O_{\text{in}}] \rangle [O_{\text{out}}]).$$

*Proof.* From the representation

$$\beta^q [O_{\text{out}}, O_{\text{in}}] = \sum_{o_{\text{in}} \in [r_{\text{in}}]} \epsilon_{(I_{\text{out}}^{-1} \circ q \circ I_{\text{in}}) o_{\text{in}}} [O_{\text{in}}] \otimes \epsilon_{o_{\text{in}}} [O_{\text{in}}]$$

and the orthonormality of the one-hot encodings of the input enumeration we get

$$\langle \beta^f [O_{\text{out}}, O_{\text{in}}], \epsilon_{I_{\text{in}}(u)} [O_{\text{in}}] \rangle [O_{\text{out}}] = \epsilon_{I_{\text{out}}^{-1}(q(u))} [O_{\text{out}}].$$

□

In comparison with the Coordinate Calculus scheme (see The. 98), the Basis Calculus produces basis vectors of a functions evaluation instead of scalars. While this seems to produce unnecessary redundancy in representing a function, we will see in the following section, that this scheme is efficient in representing compositions of functions.

## 17.3 Calculus of basis encodings

We now show the utility of basis encodings for functions, by developing tensor network representation to composed functions. We in this section use the notation of factored system representation, as developed in Part I and enumerate states of factored systems by variables  $X$  with states in  $[m]$ , instead of combinations of variables  $O$  with index interpretation functions  $I$  enumerating arbitrary sets.

### 17.3.1 Composition of function

We have already used (see The. 41), that combination of propositional formulas by connectives can be represented by contractions. We now show in a more general perspective, that in basis calculus, any composition of functions in its basis encoding the contraction of the encoded functions.

**Theorem 108** (Composition of Functions). *Let there be two maps between factored systems*

$$q : \bigtimes_{v \in \mathcal{V}^1} [m_v] \rightarrow \bigtimes_{v \in \mathcal{V}^2} [m_v]$$

and

$$g : \bigtimes_{v \in \mathcal{V}^2} [m_v] \rightarrow \bigtimes_{v \in \mathcal{V}^3} [m_v]$$

with the image system of  $q$  is the domain system of  $g$ . Then the basis encoding of the composition

$$g \circ q : \bigtimes_{v \in \mathcal{V}^1} [m_v] \rightarrow \bigtimes_{v \in \mathcal{V}^3} [m_v]$$

is the contraction

$$\beta^{g \circ q} [X_{\mathcal{V}^3}, X_{\mathcal{V}^1}] = \langle \beta^g [X_{\mathcal{V}^3}, X_{\mathcal{V}^2}], \beta^q [X_{\mathcal{V}^2}, X_{\mathcal{V}^1}] \rangle [X_{\mathcal{V}^3}, X_{\mathcal{V}^1}].$$

*Proof.* By definition we have the basis encoding of the composition as

$$\beta^{g \circ q} [X_{\mathcal{V}^3}, X_{\mathcal{V}^1}] = \sum_{x_{\mathcal{V}^1} \in \times_{v \in \mathcal{V}^1} [m_v]} \epsilon_{(g \circ q)(x_{\mathcal{V}^1})} [X_{\mathcal{V}^3}] \otimes \epsilon_{x_{\mathcal{V}^1}} [X_{\mathcal{V}^1}].$$

By using a similar representation for  $\beta^g$  and  $\beta^q$  we now show, that this coincides with the contraction of these basis encodings with closed variables  $X_{\mathcal{V}^2}$ . By the linearity of the contraction operation we get

$$\begin{aligned} \langle \beta^q, \beta^g \rangle [X_{\mathcal{V}^3}, X_{\mathcal{V}^1}] &= \sum_{x_{\mathcal{V}^1} \in \times_{v \in \mathcal{V}^1} [m_v]} \sum_{x_{\mathcal{V}^2} \in \times_{v \in \mathcal{V}^2} [m_v]} \langle (\epsilon_{g(x_{\mathcal{V}^2})} [X_{\mathcal{V}^3}] \otimes \epsilon_{x_{\mathcal{V}^2}} [X_{\mathcal{V}^2}]), \\ &\quad (\epsilon_{q(x_{\mathcal{V}^1})} [X_{\mathcal{V}^2}] \otimes \epsilon_{x_{\mathcal{V}^1}} [X_{\mathcal{V}^1}]) \rangle [X_{\mathcal{V}^3}, X_{\mathcal{V}^1}] \\ &= \sum_{x_{\mathcal{V}^1} \in \times_{v \in \mathcal{V}^1} [m_v]} \delta_{x_{\mathcal{V}^2}, x_{\mathcal{V}^1}} \cdot \epsilon_{g(x_{\mathcal{V}^2})} [X_{\mathcal{V}^3}] \otimes \epsilon_{x_{\mathcal{V}^1}} [X_{\mathcal{V}^1}] \\ &= \sum_{x_{\mathcal{V}^1} \in \times_{v \in \mathcal{V}^1} [m_v]} \epsilon_{(g \circ q)(x_{\mathcal{V}^1})} [X_{\mathcal{V}^3}] \otimes \epsilon_{x_{\mathcal{V}^1}} [X_{\mathcal{V}^1}] \\ &= \beta^{g \circ q} [X_{\mathcal{V}^3}, X_{\mathcal{V}^1}], \end{aligned}$$

where we exploited the orthonormality of the one-hot encodings to the states of  $X_{\mathcal{V}^2}$ , which contraction thus results in the delta symbol  $\delta$  applied on the respective states.  $\square$

We can use The. 108 iteratively to further decompose the function  $g$ . In this way, the basis encoding of a function consistent of multiple compositions can be represented as the contractions of all the functions. This has been applied in The. 41 to efficiently represent propositional formulas, for which syntactical expressions are given.

### 17.3.2 Compositions with real functions

We here investigate how the composition of a tensor

$$\tau : \bigtimes_{k \in [d]} [m_k] \rightarrow \mathbb{R}$$

with arbitrary functions

$$h : \mathbb{R} \rightarrow \mathbb{R}$$

can be represented. This is for example relevant, when representing coordinatewise tensor transforms (see Sect. 16.2) based on tensor network contractions. To this end we understand the tensor  $\tau [X_{[d]}]$  as a map of the states  $\bigtimes_{k \in [d]} [m_k]$  onto its by a variable  $O_\tau$  and index interpretation  $I$  enumerated image  $\text{im}(\tau)$ . We then define the restriction of  $h$  onto  $\text{im}(\tau)$  as the tensor  $h|_{\text{im}(\tau)} [O_\tau]$  with coordinates  $o_\tau$

$$h|_{\text{im}(\tau)} [O_\tau = o_\tau] = (h \circ I)(o_\tau).$$

Let us now show, how contractions with these vectors represents compositions with tensors.

**Theorem 109.** *The coordinatewise transform of any tensor  $\tau$  (see Def. 78) by a real function  $h$  is the contraction (see Figure 38)*

$$h(\tau)[X_{[d]}] = \langle \beta^\tau [O_\tau, X_{[d]}], h|_{\text{im}(\tau)} [O_\tau] \rangle [X_{[d]}] .$$

*Proof.* By the basis calculus The. 107 we have for any state  $x_{[d]} \in \times_{k \in [d]} [m_k]$ , that

$$\begin{aligned} \langle \beta^\tau [O_\tau, X_{[d]}], h|_{\text{im}(\tau)} [O_\tau] \rangle [X_{[d]} = x_{[d]}] &= \langle \beta^\tau [O_\tau, X_{[d]} = x_{[d]}], h|_{\text{im}(\tau)} [O_\tau] \rangle [\emptyset] \\ &= \left\langle \epsilon_{I_{\tau[X_{[d]} = x_{[d]}]}} [O_\tau], h|_{\text{im}(\tau)} [O_\tau] \right\rangle [\emptyset] \\ &= h(\tau)[X_{[d]} = x_{[d]}] . \end{aligned}$$

Since both tensors coincide on all coordinates, they are equal.  $\square$

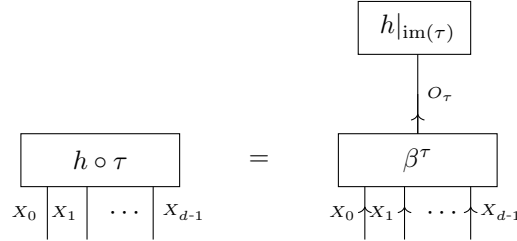


Figure 38: Representation of the composition of a tensor  $\tau$  with a real function  $h$ .

**Corollary 13.** *For any tensor  $\tau [X_{[d]}]$  we have*

$$\tau [X_{[d]}] = \langle \beta^\tau [O_\tau, X_{[d]}], \text{Id}|_{\text{im}(\tau)} [O_\tau] \rangle [X_{[d]}] .$$

*Proof.* This follows from The. 109 using  $h = \text{Id}$  and by noticing that

$$\tau [X_{[d]}] = \text{Id}(\tau)[X_{[d]}] .$$

$\square$

**Remark 22** (Tranform of basis into coordinate encodings). *Cor. 13 states in particular the transformation of the basis encoding of a function into its coordinate encoding. Given a real function  $q : \times_{k \in [d]} [m_k] \rightarrow \mathbb{R}$  we have*

$$\chi^q [X_{[d]}] = \langle \beta^q [Y_q, X_{[d]}], I_q(Y_q) \rangle [X_{[d]}] ,$$

where  $Y_q$  is a variable, which enumerates the image of  $q$  with the interpretation  $I_q(Y_q)$ .

**Corollary 14.** *For any tensor  $\tau$ , which is directed with  $X_{[d]}$  incoming, we have*

$$\mathbb{I} [X_{[d]}] = \langle \beta^\tau \rangle [X_{[d]}] .$$

*Proof.* This follows from The. 109 using  $h = \mathbb{I}$  and by noticing that

$$\mathbb{I} [X_{[d]}] = \mathbb{I}(\tau)[X_{[d]}] .$$

$\square$

### 17.3.3 Decomposition in case of structured images

When a set is structured as the cartesian product of other sets, that is

$$\mathcal{U}^{\text{out}} = \times_{k \in [d]} \mathcal{U}^k ,$$

we can enumerate it by a collection  $\{O_k : k \in [d]\}$  of enumeration variables, each with respective index interpretation maps. When the image of a function admits such a cartesian representation, we now show that the basis encoding can be represented by a contraction of basis encodings to each image coordinate.

**Theorem 110.** *Let  $q$  be a function between factored systems*

$$q : [m] \rightarrow \bigtimes_{k \in [d]} [m_k]$$

*and denote by*

$$q_k : [m] \rightarrow [m_k]$$

*the image coordinate restrictions of  $q$ , that is we have  $q = (q_0, \dots, q_{d-1})$ . Let us assign the variable  $X$  to the factored system in the domain system of  $q$  and the variables  $X_k$  for  $k \in [d]$  to the image system of  $q$ . We can then decompose the basis encoding of  $q$  into the basis encodings of its image coordinate restrictions, that is*

$$\beta^q [X_{[d]}, X] = \langle \{ \beta^{q_k} [X_k, X] : k \in [d] \} \rangle [X_{[d]}, X] .$$

*Proof.* For any  $x \in [m]$  we have

$$\begin{aligned} \beta^q [X_{[d]}, X = x] &= \epsilon_{q(x)} [X_{[d]}] \\ &= \bigotimes_{k \in [d]} \beta^{q_k} [X_k, X = x] \\ &= \langle \{ \beta^{q_k} [X_k, X = x] : k \in [d] \} \rangle [X_{[d]}] \\ &= \langle \{ \beta^{q_k} [X_k, X] : k \in [d] \} \rangle [X_{[d]}, X = x] \end{aligned}$$

and therefore equality of both tensors.  $\square$

In Chapter 18 we will apply The. 110 in The. 126 to show sparse basis CP decompositions to  $\beta^q$ . These decompositions are then applied for efficient the representation of empirical distribution, which involve the basis encoding of data maps (see Example 27), and for exponential families, which statistics have images, which are included in cartesian products of the images to each coordinate (see Example 28).

## 17.4 Selection Encodings

Selection encodings as introduced in Def. 15 are best understood in terms of linear mapping interpretations of tensors. We will first provide by basis representations a generic relation between the coordinatewise tensor definitions in this work and linear maps.

We then show the utility of this perspective in the representation of composed linear functions. The results are applicable in the exponential family theory, in the tensor representation of energies and means.

### 17.4.1 Basis representations of linear maps

Basis representations are standard linear algebra tools, where matrices are understood as linear maps between vector spaces. The state sets  $\bigtimes_{k \in [d]} [m_k]$  can be interpreted as an enumeration of basis elements  $\epsilon_x$  of the tensor space  $\bigotimes_{k \in [d]} \mathbb{R}^{m_k}$ . Along this interpretation, tensors have an interpretation as maps between tensor spaces. Any tensor and any partition of its variables into two sets can be interpreted as the basis elements of a linear map between the tensor spaces of the respective variables. Tensor valued functions on state sets  $\bigtimes_{k \in [d]} [m_k]$  are an intermediate representation.

**Definition 85.** *Let there be two tensor spaces  $V_1$  and  $V_2$  with basis by sets  $\mathcal{U}^1 \subset V_1$  and  $\mathcal{U}^2 \subset V_2$  of cardinality  $r_1$  and  $r_2$ , which are enumerated by variables  $O_1, O_2$  and index interpretation functions  $I_1, I_2$ . The basis representation of any linear map  $F \in \mathbb{L}(V_1, V_2)$  is then the tensor*

$$\tau^q [O_1, O_2] \in \mathbb{R}^{r_1} \otimes \mathbb{R}^{r_2}$$

*defined for  $o_1 \in [r_1]$  and  $o_2 \in [r_2]$  by*

$$\tau^F [O_1 = o_1, O_2 = o_2] = \langle F^{I_1(o_1)}, I_2(o_2) \rangle [\emptyset] .$$

Basis representations for compositions of linear functions can be computed via contractions of the respective basis representations, as we show next.



**Theorem 111.** *If  $F^1$  is a linear function between  $V_1$  and  $V_2$  and  $F^2$  between  $V_2$  and  $V_3$ , and let  $O_1$ ,  $O_2$  and  $O_3$  be enumerations of orthonormal bases in the spaces with index interpretation functions  $I_1$ ,  $I_2$  and  $I_3$ . We have*

$$\tau^{F^2 \circ F^1} [O_1, O_3] = \left\langle \tau^{F^2} [O_2, O_3], \tau^{F^1} [O_1, O_2] \right\rangle [O_1, O_3] .$$

*Proof.* For arbitrary  $o_1 \in [r_1]$  and  $o_3 \in [r_3]$  we have to show that

$$\tau^{F^2 \circ F^1} [O_1 = o_1, O_3 = o_3] = \left\langle \tau^{F^2} [O_2, O_3 = o_3], \tau^{F^1} [O_1 = o_1, O_2] \right\rangle [\emptyset] .$$

By definition we have

$$\tau^{F^2 \circ F^1} [O_1 = o_1, O_3 = o_3] = \langle F^2 \circ F^1(I_1(o_1)), I_3(o_3) \rangle [\emptyset] .$$

Decomposing the linear maps using their basis representation we get

$$\begin{aligned} \langle F^2 \circ F^1(I_1(o_1)), I_3(o_3) \rangle [\emptyset] &= \left\langle F^2 \left( \sum_{o_2 \in [r_2]} \tau^{F^1} [O_1 = o_1, O_2 = o_2] \cdot I_2(o_2) \right), I_3(o_3) \right\rangle [\emptyset] \\ &= \sum_{o_2 \in [r_2]} \left\langle F^2 \left( \tau^{F^1} [O_1 = o_1, O_2 = o_2] \cdot I_2(o_2) \right), I_3(o_3) \right\rangle [\emptyset] \\ &= \sum_{o_2 \in [r_2]} \left\langle \tau^{F^1} [O_1 = o_1, O_2 = o_2] \cdot \tau^{F^2} [O_2 = o_2, O_1 = o_1] \right\rangle [\emptyset] \\ &= \left\langle \tau^{F^2} [O_2, O_3 = o_3], \tau^{F^1} [O_1 = o_1, O_2] \right\rangle [\emptyset] . \end{aligned}$$

Therefore, both tensors are equivalent.  $\square$

For basis representations we thus have a similar composition theorem as for basis encodings of arbitrary functions (see The. 108). What is more, one can understand each basis encodings as a basis representation of a linear function. Along this line, the composition theorem The. 111 as the principle of linear algebra, which underlies The. 108. A typical interpretation of The. 111 is matrix multiplication, where matrices understood since matrices are basis representations of linear maps.

**Example 21** (Basis encodings as basis representations). *Let us justify that we referred to the contractions of basis encodings by basis calculus, by describing that basis encodings are a special case of basis representations. To that end, we understand the sets  $\mathcal{U}^{\text{in}} = [r_{\text{in}}]$  and  $\mathcal{U}^{\text{out}} = [r_{\text{out}}]$  as labels of a basis in  $\mathbb{R}^{r_{\text{in}}}$  and  $\mathbb{R}^{r_{\text{out}}}$ . Then, given a relation  $\mathcal{R} \subset \mathcal{U}^{\text{in}} \times \mathcal{U}^{\text{out}}$ , we define a linear map  $F : \mathbb{R}^{r_{\text{in}}} \rightarrow \mathbb{R}^{r_{\text{out}}}$  through the action on the  $i \in [r_{\text{in}}]$ -th basis element as*

$$F(\epsilon_i) = \sum_{j \in [r_{\text{out}}] : (i,j) \in \mathcal{R}} \epsilon_j .$$

Comparing the coefficients of the basis representation of  $F$  and the basis encoding  $\mathcal{R}$  we get

$$\tau^F [O_{\text{out}} = o_{\text{out}}, O_{\text{in}} = o_{\text{in}}] = \epsilon_{\mathcal{R}} [O_{\text{out}} = o_{\text{out}}, O_{\text{in}} = o_{\text{in}}] .$$

#### 17.4.2 Selection encodings as basis representations

Selection encodings (see Def. 15) are related to basis representations of linear maps as we show in the next theorem.

**Theorem 112.** *Let there be tensor spaces  $\times_{k \in [d]} [m_k]$  and  $\otimes_{s \in [n]} \mathbb{R}^{p_s}$  with basis by the one-hot encodings, enumerated by the categorical variables  $X_{[d]}$  and  $L_{[n]}$  with index interpretation functions by the one-hot map  $\epsilon$ . Given a function*

$$q : \times_{k \in [d]} [m_k] \rightarrow \otimes_{s \in [n]} \mathbb{R}^{p_s}$$

*we define a linear map  $F^q \in \mathbb{L}(\otimes_{k \in [d]} \mathbb{R}^{m_k}, \otimes_{s \in [n]} \mathbb{R}^{p_s})$  by the action on the basis elements to  $x_{[d]} \in \times_{k \in [d]} [m_k]$  as the tensors*

$$F^q(\epsilon_{x_{[d]}}) := q(x_{[d]})$$

*carrying the variables  $L_{[n]}$ . We then have*

$$\sigma^q [X_{[d]}, L_{[n]}] = \tau^{F^q} [X_{[d]}, L_{[n]}] .$$

*Proof.* We show equality on each slice with respect to the variables  $X_{[d]}$  and therefore choose arbitrary  $x_{[d]}$ . It holds by definition of selection encodings and the map  $F^q$  that

$$\sigma^q [X_{[d]} = x_{[d]}, L_0, \dots, L_{n-1}] = q(x_{[d]})[L_{[n]}] = F^q(\epsilon_{x_{[d]}})[L_{[n]}].$$

We further have

$$\begin{aligned} F^q(\epsilon_{x_{[d]}})[L_{[n]}] &= \sum_{l_{[n]}} \langle F^q(\epsilon_{x_{[d]}})[L_{[n]} = l_{[n]}], \epsilon_{l_{[n]}} [L_{[n]}] \rangle [\emptyset] \cdot \epsilon_{l_{[n]}} [L_{[n]}] \\ &= \sum_{l_{[n]}} \tau^{F^q} [X_{[d]} = x_{[d]}, L_{[n]} = l_{[n]}] \cdot \epsilon_{l_{[n]}} [L_{[n]}] \\ &= \tau^{F^q} [X_{[d]} = x_{[d]}, L_{[n]}] . \end{aligned}$$

For arbitrary  $x_{[d]}$  the slices of  $\sigma^q$  and  $\tau^{F^q}$  thus coincide, which proofs the equivalence of both tensors.  $\square$

While basis encoding works for maps from  $\times_{k \in [d]} [m_k]$  to arbitrary sets (which are enumerated), selection encodings as introduced in Def. 15 require and exploit that their image is embedded in a tensor space.

Given a selection encoding of a function, the function is retrieved by slicing with respect to the

$$q(x) = \sigma^q [X = x, L] .$$

More generally, we show in the next Lemma how to construct to any tensor and any partition of its variables functions by slicing operations, such that the tensor is the selection encoding of the function.

**Theorem 113.** *Let  $\tau [X_{\mathcal{V}}]$  be a tensor in  $\bigotimes_{v \in \mathcal{V}} \mathbb{R}^{m_v}$  and let  $A, B$  be a disjoint partition of  $\mathcal{V}$ , that is  $A \dot{\cup} B = \mathcal{V}$ . Then the function*

$$q : \times_{v \in A} [m_v] \rightarrow \bigotimes_{v \in B} \mathbb{R}^{m_v}$$

*defined for  $x_A \in \times_{v \in A} [m_v]$  as*

$$q(x_A) := \tau [X_A = x_A, X_B]$$

*obeys*

$$\sigma^q [X_A, X_B] = \tau [X_{\mathcal{V}}] ,$$

*where we understand the variables  $X_B$  as selection variables.*

*Proof.* We have for any  $x_B$  that

$$\begin{aligned} \sigma^q [X_A, X_B = x_B] &= \sum_{x_A \in \times_{v \in A} [m_v]} \epsilon_{x_A} [X_A] \otimes q(x_A)[X_B = x_B] \\ &= \sum_{x_A \in \times_{v \in A} [m_v]} \epsilon_{x_A} [X_A] \otimes \tau [X_A = x_A, X_B = x_B] \\ &= \tau [X_A, X_B = x_B] \end{aligned}$$

and the equivalence follows.  $\square$

**Example 22** (Markov Logic Networks and Proposal Distributions). *While the statistic of MLN (namely  $\mathcal{H}$ ) and the proposal distribution (namely  $\mathcal{H}^T$ ) have a common selection encoding, both result from the inverse selection encoding described in The. 113. We can construct  $\mathcal{H}^T$  by first building the selection encoding to  $\mathcal{H}$  and then applying the construction of The. 113 with  $A = L$  and  $B = X_{[d]}$ .*

We use selection encodings to represent weighted sums of functions, based on the next theorem.

**Theorem 114** (Weighted formula sums by selection encodings). *Let  $\mathcal{S}$  be a tensor valued function from  $\times_{k \in [d]} [m_k]$  to  $\mathbb{R}^p$  with image coordinates  $\mathcal{S}_l$  and let  $\theta [L]$  be a tensor. Then*

$$\left( \sum_{l \in [p]} \theta [L = l] \cdot \mathcal{S}_l \right) [X_{[d]}] : \times_{k \in [d]} [m_k] \rightarrow \mathbb{R}$$

is represented as

$$\left( \sum_{l \in [p]} \theta [L = l] \cdot S_l \right) [X_{[d]}] = \langle \sigma^S [X_{[d]}, L], \theta [L] \rangle [X_{[d]}] .$$

*Proof.* The representation holds, since for any  $x_{[d]} \in \times_{k \in [d]} [m_k]$  we have

$$\langle \sigma^S [X_{[d]}, L], \theta [L] \rangle [X_{[d]}] = x_{[d]} = \sum_{l \in [p]} q(L = l) \cdot S_l [X_{[d]}] .$$

□

This theorem shows, that while relation encodings can represent any composition with another function by a contractions, selection encodings can be used to represent linear transforms. To see this, we interpret  $S$  and  $q$  in Theorem 114 as basis decompositions of linear maps.

### 17.5 Indicator features to functions

We here provide a subspace perspective for the sparse representation of decomposable functions as tensor networks in the  $\beta$ -basis encoding scheme of basis calculus.

**Definition 86.** Given any function  $q : \mathcal{U}^{\text{in}} \rightarrow \mathcal{U}^{\text{out}}$  and index interpretation functions  $I_{\text{in}}, I_{\text{out}}$  enumerating  $\mathcal{U}^{\text{in}}$  and  $\mathcal{U}^{\text{out}}$  we define the corresponding indicator subspace of  $\mathbb{R}^{|\mathcal{U}^{\text{in}}|}$  as

$$V^q = \{ \langle \beta^q [O_{\text{out}}, O_{\text{in}}], \alpha [O_{\text{out}}] \rangle [O_{\text{out}}, O_{\text{in}}] : \alpha [O_{\text{out}}] \in \mathbb{R}^{r_{\text{out}}} \} .$$

For any function  $q$  we have  $\mathbb{I} [X_{[d]}] \in V^q$ , when choosing  $\alpha [O_{\text{out}}] = \mathbb{I} [O_{\text{out}}]$ .

We now characterize the indicator subspace as a span of boolean indicator features, indicating whether the function value coincides with an element  $y \in \text{im}(q)$ . Each indicator feature is a tensor  $\mathbb{I}_{q=y} [O_{\text{in}}]$  with coordinates

$$\mathbb{I}_{q=y} [O_{\text{in}} = o_{\text{in}}] = \begin{cases} 1 & \text{if } q(I_{\text{in}}(o_{\text{in}})) = y \\ 0 & \text{else} \end{cases} .$$

**Lemma 32.** A value subspace is spanned by the boolean indicator features of the underlying function, that is

$$V^q = \text{span} (\mathbb{I}_{q=y} [O_{\text{in}}] : y \in \text{im}(q)) .$$

*Proof.* By linearity of contractions we have

$$V^q = \text{span} (\langle \beta^q [O_{\text{out}}, O_{\text{in}}] \rangle [O_{\text{out}}], \epsilon_{o_{\text{out}}} [O_{\text{out}}] : o_{\text{out}} \in [r_{\text{out}}]) .$$

It further holds that

$$\mathbb{I}_{q=I_{o_{\text{out}}}} [O_{\text{in}}] = \langle \beta^q [O_{\text{out}}, O_{\text{in}}] \rangle [O_{\text{out}}], \epsilon_{o_{\text{out}}} [O_{\text{out}}]$$

. The claim follows as a combination of both equations. □

#### 17.5.1 Connections with computable families

We now apply indicator spaces to provide further intuition into computable families (see Def. 26).

**Lemma 33.** For any statistic  $S : \times_{k \in [d]} [m_k] \rightarrow \mathbb{R}^p$  we have

$$\Lambda^{S, \mathcal{G}^{\text{max}}} = V^S \cup \mathbb{S}$$

*Proof.* For any non-negative  $\mathbb{P} [X_{[d]}]$  we have  $\mathbb{P} [X_{[d]}] \in \Lambda^{S, \mathcal{G}^{\text{max}}}$  if and only if  $\langle \mathbb{P} [X_{[d]}] \rangle [\emptyset] = 1$  and there is an activation core  $\alpha [Y_{[p]}]$  with

$$\mathbb{P} [X_{[d]}] = \langle \beta^S [Y_{[p]}] \rangle [X_{[d]}] .$$

This is equivalent to  $\mathbb{P} [X_{[d]}] \in V^S \cup \mathbb{S}$ . □

In the other extreme of tensor network formats for the activation tensor, we have the elementary graph  $\mathcal{G}^{\text{EL}}$ . To provide a characterization of  $\Lambda^{\mathcal{S}, \mathcal{G}^{\text{EL}}}$ , we understand any statistic as a cartesian product of its features  $S_l$ . For cartesian products we can show that the distributions realizable with elementary activation tensors coincide with the subspace contraction of the spaces  $V^{S_l}$  to be introduced next.

**Definition 87.** Given two subspaces  $V^1, V^2$  of tensors with variables  $X_{V^1}, X_{V^2}$ , their contraction is

$$\langle V^1, V^2 \rangle [X_{\bar{V}}] = \{ \langle \tau^1 [X_{V^1}], \tau^2 [X_{V^2}] \rangle [X_{\bar{V}}] : \tau^1 [X_{V^1}] \in V^1, \tau^2 [X_{V^2}] \in V^2 \} .$$

We notice, that the contraction of two subspaces is in general not a subspace. This fact will in the following become clearer, where we characterize contractions of subspaces with elementarily computable distributions, which are known to not be linear subspaces.

The cartesian product of functions on the same input set  $\mathcal{U}^{\text{in}}$  and output sets  $\mathcal{U}^{1, \text{out}}, \mathcal{U}^{2, \text{out}}$  is the function

$$(q, g) : \mathcal{U}^{\text{in}} \rightarrow \mathcal{U}^{1, \text{out}} \times \mathcal{U}^{2, \text{out}}$$

with

$$(q, g)(z) = (q(z), g(z)) .$$

Its indicator subspace is the contraction of indicator subspaces

$$\begin{aligned} \langle V^q, V^g \rangle [O_{\text{in}}] &= \langle \text{span}(\{\mathbb{I}_{q=y} [O_{\text{in}}] : y \in \text{im}(q)\}), \text{span}(\{\mathbb{I}_{g=y} [O_{\text{in}}] : y \in \text{im}(g)\}) \rangle [O_{\text{in}}] \\ &= \Lambda^{\{q, g\}, \mathcal{G}^{\text{EL}}} . \end{aligned}$$

We generalize this in the next lemma to cartesian products of multiple features using that any statistic  $\mathcal{S}$  is the cartesian product of its features  $S_l$ .

**Lemma 34.** For any statistic  $\mathcal{S} : \times_{k \in [d]} [m_k] \rightarrow \mathbb{R}^p$  we have

$$\Lambda^{\mathcal{S}, \mathcal{G}^{\text{EL}}} = \langle V^{S_l} : l \in [p] \rangle [X_{[d]}] \cup \mathbb{S} .$$

where  $\mathbb{S}$  is the sphere of normed tensors in  $\bigotimes_{k \in [d]} \mathbb{R}^{m_k}$ .

*Proof.* For any non-negative tensor  $\mathbb{P} [X_{[d]}]$  we have  $\mathbb{P} [X_{[d]}] \in \Lambda^{\mathcal{S}, \mathcal{G}^{\text{EL}}}$  if and only if

$$\langle \mathbb{P} [X_{[d]}] \rangle [\emptyset] = 1$$

and there exist activation cores  $\alpha^l [Y_l]$  such that

$$\mathbb{P} [X_{[d]}] = \left\langle \bigcup_{l \in [p]} \{ \beta^{S_l} [Y_l, X_{[d]}], \alpha^l [Y_l] \} \right\rangle [X_{[d]}] .$$

Since for any  $l \in [p]$  we have

$$\{ \langle \beta^{S_l} [Y_l, X_{[d]}], \alpha^l [Y_l] \rangle [X_{[d]}] : \alpha^l [Y_l] \in \mathbb{R}^{n_l} \} = V^{S_l}$$

$\mathbb{P} [X_{[d]}] \in \Lambda^{\mathcal{S}, \mathcal{G}^{\text{EL}}}$  is equal to

$$\mathbb{P} [X_{[d]}] \in \langle V^{S_l} : l \in [p] \rangle [X_{[d]}] .$$

□

Since always  $\mathbb{I} [X_{[d]}] \in V^q$ , we have for any subspace  $V^1$

$$\langle V^1 \rangle [X_{\bar{V}}] \subset \langle V^1, V^q \rangle [X_{\bar{V}}] .$$

We can therefore understand the addition of a feature to a set of feature as a monotoneously increasing set of elementarily computable distributions, that is

$$\Lambda^{\mathcal{S}, \mathcal{G}^{\text{EL}}} \subset \Lambda^{\mathcal{S} \cup \{q\}, \mathcal{G}^{\text{EL}}} .$$

We now relate the by a function  $q$  computable distributions with those, which are computable by its indicator features with elementary activation cores.

**Lemma 35.** For any function  $q : \times_{k \in [d]} [m_k] \rightarrow \mathcal{U}^{\text{out}}$  we have

$$\Lambda^{q, \mathcal{G}^{\text{max}}} = \Lambda^{\{\mathbb{I}_{q=y} : y \in \text{im}(q)\}, \mathcal{G}^{\text{EL}}}$$

*Proof.* " $\subset$ " For any  $\mathbb{P}[X_{[d]}] \in \Lambda^{q, \mathcal{G}^{\text{max}}}$  we find an activation core  $\alpha[Y]$  such that

$$\mathbb{P}[X_{[d]}] = \langle \beta^q[Y, X_{[d]}], \alpha[Y] \rangle [X_{[d]}] .$$

Given this activation tensor  $\alpha[Y]$  we construct an elementary activation tensor

$$\bigotimes_{l \in [p]} \alpha^l[Y_l]$$

which reproduces  $\mathbb{P}[X_{[d]}]$  by contraction with the basis encoding of the indicator features to  $q$ . To this end, we define for  $l \in [p]$  two-dimensional leg vectors

$$\alpha^l[Y_l] = \begin{bmatrix} \alpha[Y = l] \\ 1 \end{bmatrix} [Y_l] .$$

For these head variables we have for any index tuple  $x_{[d]}$

$$\begin{aligned} \left\langle \bigcup_{l \in [p]} \{\beta^{\mathbb{I}_{q=I_l}}[Y_l], \alpha[Y = l]\} \right\rangle [X_{[d]} = x_{[d]}] &= \alpha^{q(x_{[d]})}[Y_l = 1] \cdot \prod_{l \in [p] : q(x_{[d]}) \neq l} \alpha^{q(x_{[d]})}[Y_l = 0] \\ &= \alpha^{q(x_{[d]})}[Y_l = 1] \\ &= \langle \beta^q[Y, X_{[d]}], \alpha[Y] \rangle [X_{[d]} = x_{[d]}] \\ &= \mathbb{P}[X_{[d]} = x_{[d]}] . \end{aligned}$$

" $\supset$ " Conversely, given  $\mathbb{P}[X_{[d]}] \in \Lambda^{\{\mathbb{I}_{q=y} : y \in \text{im}(q)\}, \mathcal{G}^{\text{EL}}}$  we find an elementary activation tensor

$$\bigotimes_{l \in [p]} \alpha^l[Y_l]$$

such that

$$\mathbb{P}[X_{[d]}] = \left\langle \bigcup_{l \in [p]} \{\beta^{\mathbb{I}_{q=I_l}}[Y_l], \alpha[Y = l]\} \right\rangle [X_{[d]}]$$

If  $\alpha^l[Y_l = 0] = 0$  for an  $l$ , then we have  $\mathbb{P}[X_{[d]}] = \langle \mathbb{I}_{q=I_l}[X_{[d]}] \rangle [X_{[d]}|\emptyset]$ , which is an element of  $\Lambda^{q, \mathcal{G}^{\text{max}}}$  since then

$$\mathbb{P}[X_{[d]}] = \langle \beta^q Y, X_{[d]}, \epsilon_l[Y] \rangle [X_{[d]}|\emptyset] .$$

In all other cases we multiply scalars to the leg tensors such that  $\alpha^l[Y_l = 0] = 1$ . Notice, that the product of all such scalars needs to be 1 since  $\langle \mathbb{P}[X_{[d]}] \rangle [\emptyset] = 1$ . We then construct an activation core  $\alpha[Y]$

$$\alpha[Y = l] = \alpha^l[Y_l = 1]$$

and have with a similar argument as in the converse proof direction

$$\mathbb{P}[X_{[d]}] = \langle \beta^q Y, X_{[d]}, \alpha[Y] \rangle [X_{[d]}] .$$

□

### 17.5.2 Composition of functions

Let  $q : \mathcal{U}^1 \rightarrow \mathcal{U}^2$  and  $h : \mathcal{U}^2 \rightarrow \mathcal{U}^3$  be arbitrary functions, then the indicator of their composition obeys

$$V^{h \circ q} \subset V^q .$$

**Example 23** (Propositional Formulas). Each formula defines by its basis encoding the subspace of  $\bigotimes_{k \in [d]} \mathbb{R}^2$

$$V^f = \text{span}(\neg f[X_{[d]}], f[X_{[d]}]) .$$

For composition of formulas  $f, h$  with a connective  $\circ$  acting on their images we have

$$V^{\circ(f, h)} \subset V^{(f, h)} .$$

A connective  $\circ$  can thus be understood as a selection of a two-dimensional subspace in the four-dimensional subspace of the cartesian product of the connected formulas. The contraction of atomic subspaces further span the space

$$\bigotimes_{k \in [d]} \mathbb{R}^2 = \text{span}(\langle V^{X_k} : k \in [d] \rangle [X_{[d]}]) .$$

### 17.5.3 Effective Representation of Partition Statistics

**Definition 88.** We call a statistic  $\mathcal{S} : \times_{k \in [d]} [m_k] \rightarrow \times_{l \in [p]} [2]$  a partition statistic if

$$\langle \sigma^{\mathcal{S}} [X_{[d]}, L] \rangle [X_{[d]}] = \mathbb{I} [X_{[d]}] .$$

We now show, that partition statistics are exactly those statistics, which are collections of indicator features to a function.

**Lemma 36.** A statistic  $\mathcal{S}$  is a partition statistic, if and only if there exists a function  $q : \times_{k \in [d]} [m_k] \rightarrow \mathcal{U}^{\text{out}}$  with an image interpretation  $I : [p] \rightarrow \mathcal{U}^{\text{out}}$  such that for all  $l \in [p]$

$$\sigma^{\mathcal{S}} [X_{[d]}, L = l] = \mathbb{I}_{q=I(l)} [X_{[d]}] .$$

*Proof.* " $\Leftarrow$ " Given any function  $q : \times_{k \in [d]} [m_k] \rightarrow \mathcal{U}^{\text{out}}$  we have

$$\sum_{l \in [p]} \mathbb{I}_{q=I(l)} [X_{[d]}] = 1$$

and the selection tensor of indicator features is thus a partition statistic.

" $\Rightarrow$ " Conversely, given a partition statistic  $\mathcal{S}$ , we can construct a function  $q_{\mathcal{S}}$  such that the partition statistic coincides with the collection of indicator features to the function. To this end we notice that for any index tuple  $x_{[d]}$  there is a unique  $l \in [p]$  such that

$$\sigma^{\mathcal{S}} [X_{[d]} = x_{[d]}, L = l] = 0 .$$

This follows from  $\sigma^{\mathcal{S}} [X_{[d]}, L]$  being boolean by assumption and  $\langle \sigma^{\mathcal{S}} [X_{[d]}, L] \rangle [X_{[d]} = x_{[d]}] = 1$ . We define a function  $q_{\mathcal{S}} : \times_{k \in [d]} [m_k] \rightarrow [p]$  with image interpretation by the identity on  $[p]$  coordinatewise by

$$q_{\mathcal{S}}(X_{[d]} = x_{[d]}) = l$$

and have for each  $l \in [p]$

$$\sigma^{\mathcal{S}} [X_{[d]}, L = l] = \mathbb{I}_{q_{\mathcal{S}}=l} [X_{[d]}]$$

□

**Example 24** (Edge statistics of Markov Networks). Each edge statistics  $\mathcal{S}_e$  of a Markov Network (see The. 15) defines a partition statistic

$$\mathcal{S}_e(x_{\mathcal{V}}) = x_e .$$

This holds since for any  $x_{\mathcal{V}}$  we have

$$\begin{aligned} \langle \sigma^{\mathcal{S}_e} X_{\mathcal{V}}, L_e \rangle [X_{\mathcal{V}} = x_{\mathcal{V}}] &= \sum_{l_e \in [m_e]} \sigma^{\mathcal{S}_e} X_{\mathcal{V}} = x_{\mathcal{V}}, L_e = l_e \\ &= \sigma^{\mathcal{S}_e} X_{\mathcal{V}} = x_{\mathcal{V}}, L_e = X_e \\ &= 1 . \end{aligned}$$

The corresponding indicators stated in Lem. 36 are enumerated by the image elements  $l_e \in m_e$  and given by

$$\mathcal{S}_{e, l_e} [x_{\mathcal{V}}] = \begin{cases} 1 & \text{if } x_e = l_e \\ 0 & \text{else} \end{cases} .$$

We have already observed in Chapter 5, that Markov Networks have a tensor-network representation involving the selection encodings of their edge statistics. This efficiency gain compared with featurewise relational encodings is in the following generalized to arbitrary partition statistics.

**Theorem 115.** For any partition statistic and any elementary activation tensor  $\bigotimes_{l \in [p]} \alpha^l [Y_l]$  we have

$$\left\langle \bigcup_{l \in [p]} \{ \beta^{\mathcal{S}_l} [Y_l, X_{[d]}], \alpha^l [Y_l] \} \right\rangle [X_{[d]}] = \langle \sigma^{\mathcal{S}} [X_{[d]}, L], \alpha [Y] \rangle [X_{[d]}]$$

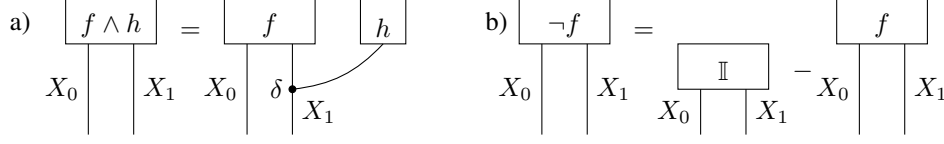


Figure 39: Decomposition schemes by hybrid calculus, using coordinatewise transforms of tensors (see Def. 78). a) Conjunction performed by coordinatewise multiplications. b) Negations performed by coordinatewise subtraction from one.

where

$$\alpha[Y] = \begin{cases} \left( \prod_{l \in [p]} \alpha^l[Y_l = 0] \right) \sum_{l \in [p]} \alpha^l[Y_l = 1] \cdot \epsilon_l[Y] & \text{if } \forall_{l \in [p]} : \alpha^l[Y_l = 0] \neq 0 \\ \alpha^l[Y_l = 1] \cdot \prod_{l \neq \bar{l}} \alpha^l[Y_l = 0] & \text{if } \alpha^{\bar{l}}[Y_{\bar{l}} = 0] = 0 \end{cases}.$$

We notice, that by definition  $\alpha[Y] = 0[Y]$  if there is more than one  $l \in [p]$  with  $\alpha^l[Y_l = 0] = 0$ .

*Proof.* Lem. 36 implies, that there is a function  $q$  such that the features of the partition statistics are the indicator features to  $q$ . Now, as in the proof of Lem. 35 we transform the activation cores to arrive at the statement.  $\square$

The. 115 is especially useful, when instantiating the distribution of a Hybrid Logic Network with a subset  $\tilde{\mathcal{F}} \subset \mathcal{F}$  of features satisfying

$$\sum_{f \in \tilde{\mathcal{F}}} f[X_{[d]}] \prec \mathbb{I}[X_{[d]}].$$

If  $\sum_{f \in \tilde{\mathcal{F}}} f[X_{[d]}] \neq \mathbb{I}[X_{[d]}]$  we can add a dummy feature  $\mathbb{I}[X_{[d]}] - \sum_{f \in \tilde{\mathcal{F}}} f[X_{[d]}]$  to  $\tilde{\mathcal{F}}$  with trivial activation core to get a partition statistics. Now, given a partition statistic  $\tilde{\mathcal{F}}$  we can instantiate the to  $\tilde{\mathcal{F}}$  corresponding tensors in the tensor network representation of The. 11 by the selection encoding  $\sigma^{\tilde{\mathcal{F}}}[X_{[d]}, L]$  as

$$\left\langle \sigma^{\tilde{\mathcal{F}}}[X_{[d]}, L], \exp[\theta[L]] \right\rangle [X_{[d]}] = \left\langle \bigcup_{l \in [p]} \{ \beta^{f_l}[Y_l, X_{[d]}], \exp[\theta[L = l] \cdot I_l(Y_l)] \} \right\rangle [X_{[d]}].$$

## 17.6 Hybrid Basis and Coordinate Calculus

In some situations, we can perform basis calculus more effectively by avoiding image enumeration variables, and instead apply coordinatewise transforms on tensors (see Def. 78). As we show here, these include conjunctions, which correspond with coordinatewise multiplication, and negation, which correspond with coordinatewise subtraction from the trivial tensor. Such schemes are applied for example in Tsilonis et al. (2024) in batchwise logical inference.

**Theorem 116.** For any formulas  $f, h$  we have

$$\langle \beta^\wedge[Y_{f \wedge h}, Y_f, Y_h], \epsilon_1[Y_{f \wedge h}] \rangle [Y_f, Y_h] = \epsilon_1[Y_f] \otimes \epsilon_1[Y_h].$$

In particular, it holds that (see Figure 39a)

$$(f \wedge h)[X_{[d]}] = \langle f, h \rangle [X_{[d]}].$$

*Proof.* We decompose

$$\beta^\wedge[Y_{f \wedge h}, Y_f, Y_h] = \epsilon_1[Y_{f \wedge h}] \otimes \epsilon_1[Y_f] \otimes \epsilon_1[Y_h] + \epsilon_0[Y_{f \wedge h}] (\mathbb{I}[Y_f, Y_h] - \epsilon_1[Y_f] \otimes \epsilon_1[Y_h])$$

and get the first claim as

$$\begin{aligned} \langle \beta^\wedge[Y_{f \wedge h}, Y_f, Y_h], \epsilon_1[Y_{f \wedge h}] \rangle [Y_f, Y_h] &= \langle \epsilon_1[Y_{f \wedge h}] \otimes \epsilon_1[Y_f] \otimes \epsilon_1[Y_h], \epsilon_1[Y_{f \wedge h}] \rangle [Y_f, Y_h] \\ &= \epsilon_1[Y_f] \otimes \epsilon_1[Y_h]. \end{aligned}$$

To show the second claim we use

$$\begin{aligned} (f \wedge h)[X_{[d]}] &= \langle \beta^f[Y_f, X_{[d]}], \beta^h[Y_h, X_{[d]}], \beta^\wedge[Y_{f \wedge h}, Y_f, Y_h], \epsilon_1[Y_{f \wedge h}] \rangle [X_{[d]}] \\ &= \langle \beta^f[Y_f, X_{[d]}], \beta^h[Y_h, X_{[d]}], (\epsilon_1[Y_f] \otimes \epsilon_1[Y_h]) \rangle [X_{[d]}] \\ &= \langle f, h \rangle [X_{[d]}]. \end{aligned}$$

$\square$

A similar decomposition holds for negations, as we show next.

**Theorem 117.** *For any formula  $f$  we have*

$$\langle \beta^\neg [Y_{\neg f}, Y_f], \epsilon_1 [Y_{\neg f}] \rangle [Y_f] = \epsilon_0 [Y_f] = \mathbb{I} [Y_f] - \epsilon_1 [Y_f] .$$

and

$$\langle \beta^\neg [Y_f, Y_{\neg f}], \epsilon_0 [Y_{\neg f}] \rangle [Y_f] = \epsilon_1 [Y_f] .$$

*In particular, it holds that (see Figure 39b)*

$$(\neg f)[X_{[d]}] = \mathbb{I} [X_{[d]}] - f [X_{[d]}] .$$

*Proof.* We have

$$\beta^\neg [Y_{\neg f}, Y_f] = \epsilon_1 [Y_{\neg f}] \otimes \epsilon_0 [Y_f] + \epsilon_0 [Y_{\neg f}] \otimes \epsilon_1 [Y_f]$$

and therefore get the second claim by contraction with  $\epsilon_0 [Y_{\neg f}]$

$$\begin{aligned} \langle \beta^\neg [Y_f, Y_{\neg f}], \epsilon_0 [Y_{\neg f}] \rangle [Y_f] &= \langle \epsilon_1 [Y_{\neg f}] \otimes \epsilon_0 [Y_f], \epsilon_0 [Y_{\neg f}] \rangle [Y_f] + \langle \epsilon_0 [Y_{\neg f}] \otimes \epsilon_1 [Y_f], \epsilon_0 [Y_{\neg f}] \rangle [Y_f] \\ &= \epsilon_1 [Y_f] . \end{aligned}$$

The first equation of the first claim follows similarly by contraction with  $\epsilon_1 [Y_{\neg f}]$  as

$$\langle \beta^\neg [Y_{\neg f}, Y_f], \epsilon_1 [Y_{\neg f}] \rangle [Y_f] = \epsilon_0 [Y_f]$$

and the second equation using that  $\mathbb{I} [X] = \epsilon_0 [X] + \epsilon_1 [X]$  and hence

$$\epsilon_0 [Y_f] = \mathbb{I} [Y_f] - \epsilon_1 [Y_f] .$$

To show the third claim, we contract the computation tensor  $\beta^f [Y_f, X_{[d]}]$  to the formula  $f$  on both sides of the first claim and get

$$\langle \beta^f [Y_f, X_{[d]}], \beta^\neg [Y_{\neg f}, Y_f], \epsilon_1 [Y_{\neg f}] \rangle [Y_f] \langle \beta^f [Y_f, X_{[d]}], \mathbb{I} [Y_f] \rangle [Y_f] - \langle \beta^f [Y_f, X_{[d]}], \epsilon_1 [Y_f] \rangle [Y_f] .$$

We simplify this equation using the trivial head contraction identity of directed tensors of Cor. 14 and Lem. 11 stating that

$$f [X_{[d]}] = \langle \beta^f [Y_f, X_{[d]}], \epsilon_1 [Y_f] \rangle [X_{[d]}] \quad \text{and} \quad \neg f [X_{[d]}] = \langle \beta^f [Y_f, X_{[d]}], \epsilon_0 [Y_f] \rangle [X_{[d]}]$$

and arrive at the third claim

$$(\neg f)[X_{[d]}] = \mathbb{I} [X_{[d]}] - f [X_{[d]}] .$$

□

These theorems provide a mean to represent logical formulas by sums of one-hot encodings. Since any propositional formula can be represented by compositions of negations and conjunctions, they are universal. We further notice, that the resulting decomposition is a basis+ CP format, as further discussed in Chapter 18. In Figure 40 we provide an example of this decomposition.

## 17.7 Applications in Machine Learning

Basis calculus provides a tool suited to represent decompositions of function by efficient tensor networks. The decomposition of function into smaller components, called neurons, is often referred to as the neural paradigm of machine learning. Our model of the neural paradigm are tensor network decompositions, seen as decomposition of functions into smaller functions, which take each other as input. Summations along input axis are avoided, when having directed and boolean tensor networks with basis calculus interpretation. Inference is then performed by contractions with basis tensors representing the input, as shown in The. 107. These contractions can further be executed neuron-wise.

What is more, basis calculus provides an efficient scheme to represent symbols such as logical connectives by their basis encodings. Basis calculus is therefore a tool for neuro-symbolic AI Garcez et al. (2019); Sarker et al. (2022); Marra et al. (2024).



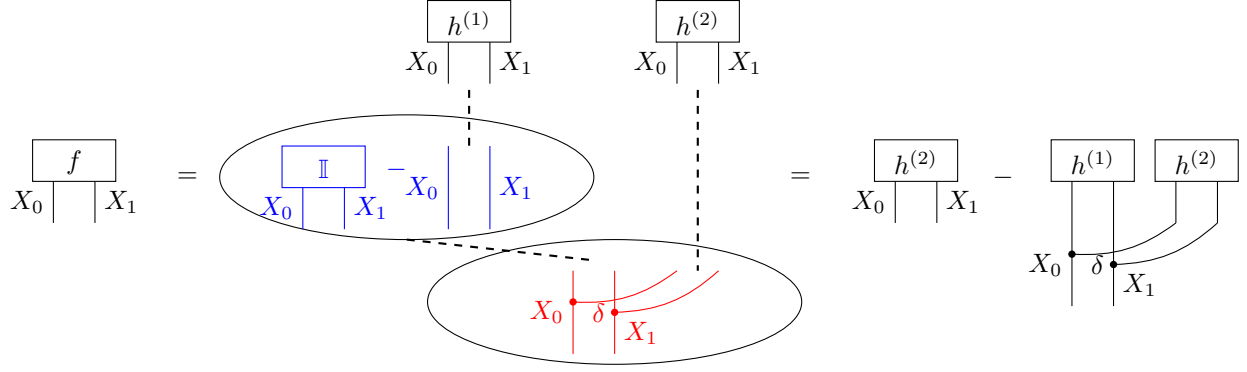


Figure 40: Example of a decomposition by hybrid calculus of a formula  $f[X_1, X_2] = \neg h^{(1)}[X_1, X_2] \wedge h^{(2)}[X_1, X_2]$  into a sum of contractions.

## 18 Sparse Calculus

We in this chapter develop sparse tensor representation formats based on constrained CP formats. Our motivation for these formats result from the connection to encoding mechanisms, which we have applied in Part I and Part II, and to sparse optimization formats.

### 18.1 CP Decomposition

The CP decomposition is one way to generalize the ranks of matrices to tensors. It is oriented on the Singular Value Decomposition of matrices, providing a representation of the matrix as a weighed sum of the tensor product of singular vectors. Given a matrix  $M[X_0, X_1]$ , we enumerate its singular values by  $I$  taking values in  $[n]$  and store them in a vector  $\lambda[I]$ . With the corresponding singular vectors by  $\rho^0[X_0, I]$  and  $\rho^1[X_1, I]$ , the singular value decomposition of  $M$  is

$$M[X_0, X_1] = \sum_{i \in [n]} \lambda[I = i] \cdot \rho^0[X_0, I = i] \otimes \rho^1[X_1, I = i] .$$

Here the smallest  $n$  such that this decomposition exists, is the matrix rank  $\text{rank}(M)$ . In contraction notation we abbreviate this to

$$M[X_0, X_1] = \langle \lambda[I], \rho^0[X_0, I], \rho^1[X_1, I] \rangle [X_0, X_1] .$$

Given a tensor of higher order, a generalization of this decomposition is a tensor product over multiple vectors, as we define next.

**Definition 89.** A CP decomposition of size  $n$  of a tensor  $\tau[X_{[d]}] \in \bigotimes_{k \in [d]} \mathbb{R}^{m_k}$  is a collections of a scalar core  $\lambda[I]$  and leg cores  $\rho^k[I, X_k]$  for  $k \in [d]$ , where  $I$  is an enumeration variable taking values in  $[n]$ , such that

$$\tau[X_{[d]}] = \langle \{\lambda[I]\} \cup \{\rho^k[X_k, I] : k \in [d]\} \rangle [X_{[d]}] .$$

We say that the CP Decomposition is

- *directed*, when for each  $k$  the core  $\rho^k$  is directed with  $I$  incoming and  $X_k$  outgoing.
- *boolean*, when for each  $k$  the core  $\rho^k$  is boolean.
- *basis*, where we demand both properties, that is for each  $k \in [d]$  and  $i \in [n]$

$$\rho^k[X_k, I = i] \in \{\epsilon_{[x_k]}[X_k] x_k \in [m_k]\} .$$

- *basis+*, when for each  $k \in [d]$  and  $i \in [n]$

$$\rho^k[X_k, I = i] \in \{\epsilon_{[x_k]}[X_k] x_k \in [m_k]\} \cup \{\mathbb{I}[X_k]\} .$$

We denote by  $\text{rank}(\tau)$ , respectively  $\text{rank}^{\text{bin}}(\tau)$ ,  $\text{rank}^{\text{bas}}(\tau)$  and  $\text{rank}^{\text{bas+}}(\tau)$  the minimal cardinality such that  $\tau$  has a CP Decomposition, respectively with directed cores, boolean cores, basis cores and basis+ cores.

All ranks have a naive bound by the space dimension, which is obvious from the coordinate decomposition (see Chapter 16)

$$\tau [X_{[d]}] = \sum_{x_{[d]} \in \times_{k \in [d]} [m_k]} \tau [X_{[d]} = x_{[d]}] \cdot \bigotimes_{k \in [d]} \epsilon_k [X_k] .$$

If we construct  $i$  as an enumeration of the coordinates in  $\times_{k \in [d]} [m_k]$ , that is  $n = \prod_{k \in [d]} m_k$ , this is a CP decomposition, which is basis and therefore also directed, boolean and basis+.

CP decomposition as a tensor network format come with some drawbacks. The set of tensors with a fixed rank are not closed [Beylkin and Mohlenkamp \(2005\)](#) and approximation problems are often ill posed [de Silva and Lim \(2008\)](#). Since as a consequence their numerical treatment comes with many problems [Espig et al. \(2012\)](#), alternative formats have gained popularity. Common formats are the TUCKER-format originally introduced in [Hitchcock \(1927\)](#), and often referred to as higher-order singular value decomposition, and the more recently developed TT and HT decomposition formats (see Chapter 2). Given a HT the best approximation of a tensor always exists (Theorem 11.58 in [Hackbusch \(2012\)](#)).

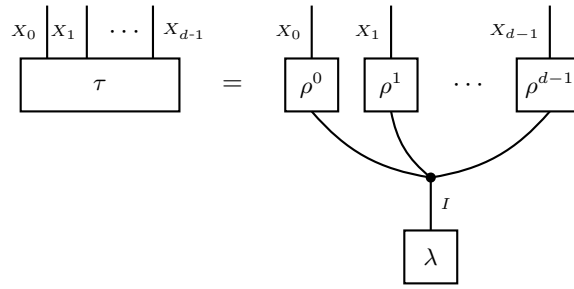


Figure 41: Tensor Network diagram of a generic CP decomposition (see Def. 89)

### 18.1.1 Directed Leg Cores

The constraint of directionality of the leg cores does not influence decomposability of a tensor, as we show next.

**Lemma 37.** *For any tensor  $\tau [X_{[d]}]$  we have*

$$\text{rank}(\tau) = \text{rank}^{\text{dir}}(\tau) .$$

*Proof.* Let there be a CP decomposition of  $\tau$  by

$$\tau [X_{[d]}] = \langle \{ \lambda[I] \} \cup \{ \rho^k [X_k, I] : k \in [d] \} \rangle [X_{[d]}] .$$

We then transform the scalar core to another core  $\tilde{\lambda}[I]$  with coordinates to  $i \in [n]$  by

$$\tilde{\lambda}[I = i] = \lambda[I = i] \cdot \prod_{k \in [d]} \langle \rho^k [X_k, I = i] \rangle [\emptyset] .$$

It follows for any  $i \in [n]$ , that

$$\lambda[I = i] \cdot \bigotimes_{k \in [d]} \rho^k [X_k, I = i] = \tilde{\lambda}[I = i] \cdot \bigotimes_{k \in [d]} \langle \rho^k [X_k, I = i] \rangle [X_k | \emptyset]$$

and thus

$$\tau [X_{[d]}] = \left\langle \{ \tilde{\lambda}[I] \cup \{ \langle \rho^k [X_k, I] \rangle [X_k, I | \emptyset] : k \in [d] \} \} \right\rangle [X_{[d]}] .$$

We have thus constructed a directed CP decomposition of same size  $n$  to an arbitrary CP decomposition and conclude that  $\text{rank}(\tau) = \text{rank}^{\text{dir}}(\tau)$ .  $\square$

### 18.1.2 Basis CP decompositions and the $\ell_0$ -norm

The slices of directed and boolean tensors with respect to incoming variables are basis tensors. We have thus called CP decomposition with the restriction of directed and boolean leg vectors basis CP decomposition. Based on this intuition, we can interpret basis CP decomposition by mappings to non-vanishing coordinates of the decomposed tensor. To start, let us define the number of nonzero coordinates of tensors by the  $\ell_0$ -norm.

**Definition 90.** The  $\ell_0$ -norm counts the nonzero coordinates of a tensor by

$$\ell_0(\tau) = \#\{x_0, \dots, x_{d-1} : \tau_{x_0, \dots, x_{d-1}} \neq 0\}.$$

The  $\ell_0$ -norm is not a norm, but at each tensor the limit of  $\ell_p$ -norms (which are norms for  $p \geq 1$ ) for  $p \rightarrow 0$ .

The  $\ell_0$  norm is the number of non-vanishing coordinates of a tensor. We understand the leg cores as the basis encoding of functions mapping to the slices of these coordinates given an enumeration. This is consistent with the previous analysis of Chapter 17, where we characterized boolean and directed cores by the encoding of associated functions. Based on this idea, we can prove, that any tensor has a directed and boolean CP decomposition with rank  $\ell_0(\tau)$ .

**Theorem 118.** For any tensor  $\tau [X_{[d]}]$  we have

$$\text{rank}^{\text{bas}}(\tau) = \ell_0(\tau).$$

*Proof.* Let us first show, that  $\text{rank}^{\text{bas}}(\tau) \leq \ell_0(\tau)$ . We find a map

$$D : [\ell_0(\tau)] \rightarrow \bigtimes_{k \in [d]} [m_k]$$

whose image is the set of non-vanishing coordinates of  $\tau [X_{[d]}]$ . Denoting its image coordinate maps by  $D_k$  we have

$$\tau [X_{[d]}] = \sum_{j \in [m]} \lambda[D((j))] \left( \bigotimes_{k \in [d]} \epsilon_{D_k(j)} [X_k] \right).$$

This is a basis CP decomposition of size  $\ell_0(\tau)$  and we thus have  $\text{rank}^{\text{bas}}(\tau) \leq \ell_0(\tau)$ .

Conversely, let us show  $\text{rank}^{\text{bas}}(\tau) \geq \ell_0(\tau)$ . Any basis CP decomposition of  $\tau$  with size  $r$  would have at most  $r$  coordinates different from zero and thus  $\ell_0(\tau) \leq r$ . Thus, there cannot be a CP decomposition with a dimension  $r \leq \ell_0(\tau)$ .  $\square$

The next theorem relates the basis CP decomposition with encodings of  $d$ -ary relations (see Def. 83).

**Theorem 119.** Any boolean tensor  $\tau [X_{[d]}] \in \bigotimes_{k \in [d]} \mathbb{R}^{m_k}$  is the encoding of a  $d$ -ary relation  $\mathcal{R} \subset \bigtimes_{k \in [d]} [m_k]$  with cardinality

$$|\mathcal{R}| = \text{rank}^{\text{bas}}(\tau).$$

*Proof.* We find a basis CP decomposition of  $\tau [X_{[d]}]$  with  $n = \text{rank}^{\text{bas}}(\tau)$ . Since  $\tau [X_{[d]}]$  is boolean, and since each  $i$  labels a disjoint non-vanishing coordinate (see proof of The. 118), the decomposition has a trivial scalar core  $\lambda[I] = \mathbb{I}[I]$ . It follows, that

$$\tau [X_{[d]}] = \sum_{i \in [n]} \left( \bigotimes_{k \in [d]} \rho^k [X_k, I = i] \right)$$

Since the CP decomposition is basis, the slice  $\rho^k [X_k, I = i]$  is for any  $k \in [d]$  and  $i \in [n]$  a basis vector. We then define

$$x_k^i = \epsilon^{-1}(\rho^k [X_k, I = i])$$

and notice, that for the relation

$$\mathcal{R} = \{x_{[d]}^i : i \in [n]\} \subset \bigtimes_{k \in [d]} [m_k]$$

we have

$$\epsilon_{\mathcal{R}} [X_{[d]}] = \sum_{i \in [n]} \left( \bigotimes_{k \in [d]} \rho^k [X_k, I = i] \right).$$

This coincides with the above CP decomposition of  $\tau [X_{[d]}]$  and the claim is established.  $\square$

**Remark 23** (Matrix Storage of basis CP decompositions). *The storage demand of any CP decomposition is at most linear in the size and the sum of its leg dimension. When we have a basis CP decomposition, this demand can be further improved. The basis vectors can be stored by its preimage of the one hot encoding  $\epsilon_{\cdot}$ , that is the number of the basis vector in  $[m]$ . This reduces the storage demand of each basis vector to the logarithms of the space dimension without the need of storing the full vector. More precisely, we can define a leg selecting variable  $L$  taking values in  $[d + 1]$  and store a basis CP decomposition of size  $n$  by the matrix*

$$M[I, L] \in \mathbb{R}^{m \times (d+1)}$$

defined for  $i \in [n]$  and  $k \in [d]$  by

$$M[I = i, L = k] = \begin{cases} \lambda[I = i] & \text{if } k = d \\ \epsilon^{-1}(\rho^k [X_k, I = i]) & \text{else} \end{cases}.$$

This is a common trick to store relational databases.

### 18.1.3 Basis+ CP decompositions and polynomials

The basis+ CP decompositions are closely related to monomial decompositions of a tensor, which we will define next.

**Definition 91.** *A monomial decomposition of a tensor  $\tau [X_{[d]}] \in \bigotimes_{k \in [d]} \mathbb{R}^{m_k}$  is a set  $\mathcal{M}$  of tuples  $(\lambda, A, x_A)$  where  $\lambda \in \mathbb{R}$ ,  $A \subset [d]$  and  $x_A \in \times_{k \in A} [m_k]$  such that*

$$\tau [X_{[d]}] = \sum_{(\lambda, A, x_A) \in \mathcal{M}} \lambda \cdot \langle \epsilon_{x_A} [X_A] \rangle [X_{[d]}]. \quad (46)$$

For any tensor  $\tau [X_{[d]}] \in \bigotimes_{k \in [d]} \mathbb{R}^{m_k}$  we define its polynomial sparsity of order  $r$  as

$$\text{rank}^r(\tau) = \min \left\{ |\mathcal{M}| : \tau [X_{[d]}] = \sum_{(\lambda, A, x_A) \in \mathcal{M}} \lambda \cdot \langle \epsilon_{x_A} [X_A] \rangle [X_{[d]}], \forall (\lambda, A, x_A) \in \mathcal{M} \ |A| \leq r \right\}$$

We refer to the terms in a decomposition (46) in Def. 91 as monomials of boolean features, which are enumerated by pairs  $(k, x_k)$  and indicate whether the variable  $X_k$  is in state  $x_k \in [m_k]$ . Each such boolean features is represented by the indicator

$$\mathbb{I}_{X_k=x_k} [X_k] = \epsilon_{x_k} [X_k].$$

The monomial of multiple such boolean features indicates, whether all variables labelled by a set  $A$  are in the state  $X_A$ . We have

$$\mathbb{I}_{\forall k \in A: X_k=x_k} [X_A] = \epsilon_{x_A} [X_A] = \bigotimes_{k \in A} \epsilon_{x_k} [X_k].$$

The states of the variables labeled by  $k \in [d]/A$  are not specified in the monomial and the indicators are trivially extended to

$$\langle \epsilon_{x_A} [X_A] \rangle [X_{[d]}] = \epsilon_{x_A} [X_A] \otimes \mathbb{I} [X_{[d]/A}].$$

Since we are working with boolean features, there is no need to consider higher-order powers of individual features, since for any  $n \in \mathbb{N}$ ,  $n \geq 1$  and any boolean value  $z \in \{0, 1\}$  we have  $z^n = z$ .

For some monomial orders  $r < d$  there are tensors  $\tau [X_{[d]}]$ , which do not have a monomial decomposition of order  $r$ . In that case the minimum is over an empty set and we define  $\text{rank}^r(\tau) = \infty$ . We characterize in the next theorem the set of tensors with monomial decompositions of order  $r$ .

**Theorem 120.** For any  $d, r$ , the set of tensors of  $d$  variables with leg dimension  $m$ , which have a monomial decomposition of order  $r$ , is a linear subspace  $V^{d,r}$  with dimension

$$\dim(V^{d,r}) \leq \sum_{s \in [r]} m^s \binom{d}{s}.$$

*Proof.* The set of tensors admitting a monomial decomposition of order  $r$  is closed under addition and scalar multiplication. Specifically, the sum of two such tensors retains a monomial decomposition, formed by concatenating their respective decompositions. Scalar multiplication can be performed by a rescaling of each scalar  $\lambda$  and therefore preserves the decomposition structure. Hence, these tensors form a linear subspace.

To bound the dimension of this subspace, we consider tensors of the form  $\langle \epsilon_{x_A} \rangle [X_{[d]}]$ . The number of such tensors is given by

$$\sum_{s \in [r]} m^s \binom{d}{s}.$$

Since any tensor with a monomial decomposition is a weighted sum of those, this provides an upper bound on the dimension.

We notice, that the set of slices is in general not linear independent, and therefore forms a frame instead of a linear basis Casazza et al. (2013). The number of elements in the frame is therefore in general a loose upper bound on the dimension.  $\square$

The. 120 states, that the tensors admitting a monomial decomposition of a small order build a low-dimensional subspace in the  $m^d$  dimensional space of tensors, since for  $r \ll d$  we have

$$\dim(V^{d,r}) \ll m^d.$$

If  $r \geq d$ , we always find a monomial decomposition by an enumeration of nonzero coordinates. In the next theorem, we show that in that case the  $\text{rank}^r(\tau)$  furthermore coincides with the basis+ CP rank  $\text{rank}^{\text{bas}+}(\tau)$ .

**Theorem 121.** For any tensor  $\tau [X_{[d]}] \in \bigotimes_{k \in [d]} \mathbb{R}^{m_k}$  we have

$$\text{rank}^d(\tau) = \text{rank}^{\text{bas}+}(\tau).$$

In case of two-dimensional legs, that is  $m_k = 2$  for all  $k \in [d]$ , we also have

$$\text{rank}^{\text{bin}}(\tau) = \text{rank}^d(\tau).$$

*Proof.* To proof the first claim, we construct a basis+ CP decomposition given a monomial decomposition and vice versa. To show  $\text{rank}^d(\tau) \geq \text{rank}^{\text{bas}+}(\tau)$ , let there be an arbitrary tensor  $\tau [X_{[d]}]$  with a monomial decomposition by  $\mathcal{M}$  with  $|\mathcal{M}| = m$  and let us enumerate the elements in  $\mathcal{M}$  by  $(\lambda^i, A^i, x_{A^i}^i)$  for  $i \in [n]$ . We define for each  $k \in [d]$  the tensors

$$\rho^k [I, X_k] = \left( \sum_{i \in [n] : k \in A} \epsilon_i [I] \otimes \epsilon_{x_k^i} [X_k] \right) + \left( \sum_{i \in [n] : k \notin A} \epsilon_i [I] \otimes \mathbb{I} [X_k] \right)$$

and

$$\lambda [I] = \sum_{i \in [n]} \lambda^i \cdot \epsilon_i [I]$$

and notice that

$$\begin{aligned} \tau [X_{[d]}] &= \sum_{i \in [n]} \lambda^i \cdot \langle \epsilon_{x_A^i} \rangle [X_{[d]}] \\ &= \sum_{i \in [n]} \left( \lambda [I = i] \cdot \bigotimes_{k \in [d]} \rho^k [I = i, X_k] \right) \\ &= \langle \{ \lambda [I] \} \cup \{ \rho^k [I, X_k] : k \in [d] \} \rangle [X_{[d]}]. \end{aligned}$$

By construction this is a basis+ CP decomposition with rank  $n$ . Since any monomial decomposition can be transformed into a basis+ CP decomposition with same rank we have

$$\text{rank}^d(\tau) \geq \text{rank}^{\text{bas}+}(\tau) .$$

To show  $\text{rank}^d(\tau) \leq \text{rank}^{\text{bas}+}(\tau)$ , let there now be a basis+ CP decomposition of an arbitrary  $\tau [X_{[d]}]$ . We define for each  $i \in [n]$

$$A^i = \{k \in [d] : \rho^k [I = i, X_k] \neq \mathbb{I} [X_k]\} \quad \text{and} \quad x_A^i = \{\epsilon^{-1}(\rho^k [I = i, X_k]) : k \in A\}$$

where by  $\epsilon^{-1}(\cdot)$  we denote the inverse of the one-hot encoding.

We notice that this is a monomial decomposition of  $\tau [X_{[d]}]$  to the tuple set

$$\mathcal{M} = \{(\lambda[I = i], A^i, x_A^i) : i \in [n]\} .$$

It follows from this that

$$\text{rank}^d(\tau) \leq \text{rank}^{\text{bas}+}(\tau)$$

and the first claim is shown.

The second claim follows from the observation, that the set of non-vanishing boolean vectors coincides with the set of one-hot encodings extended by the trivial vector. Thus, a CP decomposition with non-vanishing slices is boolean if and only if it is basis+. This establishes, that both ranks are equal, since a CP decomposition of minimal rank cannot contain non-vanishing slices.  $\square$

**Remark 24** (Sparse representation of propositional formulas). *When all leg dimensions of a boolean tensor  $\tau$  are 2, we can further interpret  $\tau$  as a logical formula. We can use the boolean CP decomposition of any tensor  $\tilde{\tau}$  with  $\mathbb{I}_{\neq 0}(\tilde{\tau}) = \tau$  as a CNF of  $\tau$ . Finding the sparsest CNF thus amounts to finding the  $\tilde{\tau}$  with minimal  $\text{rank}^d(\tilde{\tau})$  such that  $\mathbb{I}_{\neq 0}(\tilde{\tau}) = \tau$ .*

**Remark 25** (Matrix Storage of basis+ CP decompositions). *We can adapt the storage format of Remark 23 from basis to basis+ CP decompositions. To this end, let there be a basis+ CP decomposition of a tensor with scalar core  $\lambda[I]$  and leg cores  $\{\rho^k [X_k, I] : k \in [d]\}$ . We use a value  $z \in \mathbb{R}/\text{im}(\lambda)$  distinguished from the coordinates of the scalar core and define a matrix*

$$M^z [I, L]$$

where  $L$  takes values in  $[d]$ , coordinatewise as

$$M^z [I = i, L = k] = \begin{cases} \lambda[I = i] & \text{if } k = d \\ z & \text{if } \rho^k [X_k, I = i] = \mathbb{I} [X_k] \\ \epsilon^{-1}(\rho^k [X_k, I = i]) & \text{else} \end{cases} .$$

## 18.2 Constructive Bounds on CP Ranks

After having defined different CP decompositions, let us investigate bounds on their ranks, which proofs rely on explicit core constructions.

### 18.2.1 Cascade of ranks

We start by showing a cascade of bounds of CP ranks, when demanding different leg restrictions as in Def. 89.

**Theorem 122.** *For any tensor  $\tau [X_{[d]}] \in \bigotimes_{k \in [d]} \mathbb{R}^{m_k}$  we have*

$$\text{rank}(\tau) = \text{rank}^{\text{dir}}(\tau) \leq \text{rank}^{\text{bin}}(\tau) \leq \text{rank}^{\text{bas}+}(\tau) \leq \text{rank}^{\text{bas}}(\tau) .$$

*Proof.* The equality  $\text{rank}(\tau) = \text{rank}^{\text{dir}}(\tau)$  has been established in Lem. 37. The further inequalities follow by consecutive subset relations of the set of allowed leg slices in the respective CP decompositions. These imply, that any basis CP decomposition of a tensor  $\tau [X_{[d]}]$  is also a basis+ CP decomposition, further that any basis+ CP decomposition is also a boolean CP decomposition and that any boolean CP decomposition is trivially an unrestricted CP decomposition. Thus, the ranks are minima of enlarging sets and the claimed rank cascade is established.  $\square$

Let us notice, that the stated bounds are not tight in general. To give an example, let us consider the tensor  $\tau[X_{[d]}] = \mathbb{I}[X_{[d]}]$ , for which we have

$$\text{rank}^{\text{bas}}(\tau) = \ell_0(\tau) = \prod_{k \in [d]} m_k.$$

Since in the other restricted CP formats we can choose trivial slices to the leg cores, we have

$$\text{rank}^{\text{bas}+}(\tau) = 1 = \text{rank}^{\text{bin}}(\tau) = \text{rank}^{\text{dir}}(\tau) = \text{rank}(\tau).$$

The trivial tensor serves thus as an example, where the demand of the the storage format in Remark 23 has an exponential overhead compared to the storage format in Remark 25.

### 18.2.2 Operations on CP decompositions

When using CP decompositions of tensors in practice applications, such as those investigated in Part I and Part II, we have to perform numerical manipulations in the form of summations, contractions and normalizations of the represented tensors. Let us here investigate, how these operations influence the decomposition.

**Summation** We start with the sum of tensors in a CP decomposition, which can be captured by a concatenation of the slices.

**Theorem 123.** *For any collections of tensors  $\{\tau^l[X_{\mathcal{V}}] : l \in [p]\}$  with identical variables and scalars  $\lambda^l \in \mathbb{R}$  for  $l \in [p]$  we have*

$$\text{rank}\left(\sum_{l \in [p]} \lambda^l \cdot \tau^l\right) \leq \sum_{l \in [p]} \text{rank}(\tau^l).$$

*The bound still holds, when we replace on both sides  $\text{rank}(\cdot)$  by  $\text{rank}^{\text{bin}}(\cdot)$ , by  $\text{rank}^{\text{bas}}(\cdot)$  or by  $\text{rank}^{\text{bas}+}(\cdot)$ .*

*Proof.* Products with scalars do not change the rank, since they just rescale the core  $\lambda$ . The sum of CP decomposition is just the combination of all slices, thus the rank is at most additive.  $\square$

Let us notice, that the upper bound is loose in many applications. For example, if two slice tuples of two decomposed tensors agree on  $x_A, A$ , then their sum can be performed by a sum of the corresponding scalar.

**Contraction** We continue to show rank bounds for arbitrary contractions by the product of the ranks of contracted tensors.

**Theorem 124.** *For any tensor network  $\tau^{\mathcal{G}}[X_{\mathcal{V}}]$  on a graph  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ , we have for any subset  $\tilde{\mathcal{V}} \subset \mathcal{V}$*

$$\text{rank}(\langle \tau^{\mathcal{G}} \rangle [X_{\tilde{\mathcal{V}}}]) \leq \prod_{e \in \mathcal{E} : \tilde{\mathcal{V}} \cap e \neq \emptyset} \text{rank}(\tau^e).$$

*The bound still holds, when we replace on both sides  $\text{rank}(\cdot)$  by  $\text{rank}^{\text{bin}}(\cdot)$ , by  $\text{rank}^{\text{bas}}(\cdot)$  or by  $\text{rank}^{\text{bas}+}(\cdot)$ .*

Remarkably, in The. 124 the upper bound on the CP rank is build only by the ranks of the tensor cores, which have remaining open edges. We prepare for its proof by first showing the following lemmata.

**Lemma 38.** *For any tensors  $\tau^1[X_{\mathcal{V}^1}]$  and  $\tau^2[X_{\mathcal{V}^2}]$  and any set of variables  $\tilde{\mathcal{V}} \subset \mathcal{V}^1 \cup \mathcal{V}^2$  we have*

$$\text{rank}(\langle \tau^1, \tau^2 \rangle [X_{\tilde{\mathcal{V}}}]) \leq \text{rank}(\tau^1) \cdot \text{rank}(\tau^2).$$

*The bound still holds, when we replace on both sides  $\text{rank}(\cdot)$  by  $\text{rank}^{\text{bin}}(\cdot)$ , by  $\text{rank}^{\text{bas}}(\cdot)$  or by  $\text{rank}^{\text{bas}+}(\cdot)$ .*

*Proof.* Let there be CP decompositions of  $\tau^1[X_{\mathcal{V}^1}]$  and  $\tau^2[X_{\mathcal{V}^2}]$  by

$$\tau^1[X_{\mathcal{V}^1}] = \langle \{\lambda^1[I_1]\} \cup \{\rho^{1,k}[X_k, I_1] : k \in \mathcal{V}^1\} \rangle [X_{\mathcal{V}^1}]$$

and

$$\tau^2[X_{\mathcal{V}^2}] = \langle \{\lambda^2[I_2]\} \cup \{\rho^{2,l}[X_l, I_2] : l \in \mathcal{V}^2\} \rangle [X_{\mathcal{V}^2}].$$

By linearity of contractions we have

$$\begin{aligned}
 \langle \tau^1, \tau^2 \rangle [X_{\tilde{\mathcal{V}}}] &= \sum_{i_1 \in n_1} \sum_{i_2 \in n_2} \lambda^1[I_1 = i_1] \cdot \lambda^2[I_2 = i_2] \\
 &\quad \cdot \langle \{ \rho^{1,k} [X_k, I_1 = i_1] : k \in \mathcal{V}^1 \} \cup \{ \rho^{2,l} [X_l, I_2 = i_2] : l \in \mathcal{V}^2 \} \rangle [X_{\tilde{\mathcal{V}}}] \\
 &= \sum_{i_1 \in n_1} \sum_{i_2 \in n_2} \lambda^1[I_1 = i_1] \cdot \lambda^2[I_2 = i_2] \cdot \\
 &\quad \left\langle \{ \rho^{1,k} [X_k, I_1 = i_1] : k \in \mathcal{V}^1 / \tilde{\mathcal{V}} \} \cup \{ \rho^{2,l} [X_l, I_2 = i_2] : l \in \mathcal{V}^2 / \tilde{\mathcal{V}} \} \right\rangle [\emptyset] \cdot \\
 &\quad \bigotimes_{k \in \tilde{\mathcal{V}}} \rho^k [X_k, I_1 = i_1, I_2 = i_2] ,
 \end{aligned}$$

where we denote

$$\rho^k [X_k, I_1 = i_1, I_2 = i_2] = \begin{cases} \rho^{1,k} [X_k, I_1 = i_1] & \text{if } k \notin \mathcal{V}^2 \\ \rho^{2,k} [X_k, I_2 = i_2] & \text{if } k \notin \mathcal{V}^1 \\ \langle \rho^{1,k} [X_k, I_1 = i_1], \rho^{2,k} [X_k, I_2 = i_2] \rangle [X_k] & \text{else} \end{cases} .$$

Note, that since  $k \in \tilde{\mathcal{V}} \subset \mathcal{V}^1 \cup \mathcal{V}^2$ , these slices are well-defined. We build a new decomposition variable  $I$  enumerating the summands to indices  $[n_1] \times [n_2]$  and have thus found a CP decomposition of  $\langle \tau^1, \tau^2 \rangle [X_{\tilde{\mathcal{V}}}]$  of size  $n = n_1 \cdot n_2$ . This shows the claim in the case of  $\text{rank}(\cdot)$ .

When the CP decompositions of  $\tau^1$  and  $\tau^2$  are boolean, basis or basis+, then the property is preserved in the constructed CP decomposition, since the constructed slices  $\rho^k [X_k, I_1 = i_1, I_2 = i_2]$  are either copies of the leg cores or their contractions and the respective property is preserved in both cases. Thus, the constructive rank bounds hold also for  $\text{rank}^{\text{bin}}(\cdot)$ ,  $\text{rank}^{\text{bas}}(\cdot)$  and  $\text{rank}^{\text{bas+}}(\cdot)$ .  $\square$

When one core of the contracted tensor network does not contain variables which are left open, we can drastically sharpen the bound provided by Lem. 38 as we show next.

**Lemma 39.** *For any two tensors  $\tau^1 [X_{\mathcal{V}^1}]$ ,  $\tau^2 [X_{\mathcal{V}^2}]$  and any set  $\tilde{\mathcal{V}}$  with  $\tilde{\mathcal{V}} \cap \mathcal{V}^2 = \emptyset$  we have*

$$\text{rank}(\langle \tau^1, \tau^2 \rangle [X_{\tilde{\mathcal{V}}}]) \leq \text{rank}(\tau^1) .$$

*The bound still holds, when we replace on both sides  $\text{rank}(\cdot)$  by  $\text{rank}^{\text{bin}}(\cdot)$ , by  $\text{rank}^{\text{bas}}(\cdot)$  or by  $\text{rank}^{\text{bas+}}(\cdot)$ .*

*Proof.* As in the proof of Lem. 38 we assume a CP decomposition of  $\tau^1 [X_{\mathcal{V}^1}]$  and  $\tau^2 [X_{\mathcal{V}^2}]$  and use the linearity of contractions to get

$$\begin{aligned}
 \langle \tau^1, \tau^2 \rangle [X_{\tilde{\mathcal{V}}}] &= \sum_{i_1 \in n_1} \sum_{i_2 \in n_2} \lambda^1[I_1 = i_1] \cdot \lambda^2[I_2 = i_2] \cdot \\
 &\quad \left\langle \{ \rho^{1,k} [X_k, I_1 = i_1] : k \in \mathcal{V}^1 / \tilde{\mathcal{V}} \} \cup \{ \rho^{2,l} [X_l, I_2 = i_2] : l \in \mathcal{V}^2 / \tilde{\mathcal{V}} \} \right\rangle [\emptyset] \cdot \\
 &\quad \bigotimes_{k \in \tilde{\mathcal{V}}} \rho^{1,k} [X_k, I_1 = i_1] ,
 \end{aligned}$$

where we used that  $\tilde{\mathcal{V}} \cup \mathcal{V}^2 = \emptyset$ . By rearranging the sum of  $i_2$ , we have a CP decomposition with decomposition variable  $I_1$  and slices

$$\begin{aligned}
 \lambda[I_1 = i_1] &= \sum_{i_2 \in n_2} \lambda^1[I_1 = i_1] \cdot \lambda^2[I_2 = i_2] \cdot \\
 &\quad \left\langle \{ \rho^{1,k} [X_k, I_1 = i_1] : k \in \mathcal{V}^1 / \tilde{\mathcal{V}} \} \cup \{ \rho^{2,l} [X_l, I_2 = i_2] : l \in \mathcal{V}^2 / \tilde{\mathcal{V}} \} \right\rangle [\emptyset] .
 \end{aligned}$$

This shows the rank bound for  $\text{rank}(\cdot)$ . The properties of the CP decomposition are trivially inherited by the constructed decomposition, since the leg cores of the decomposition of  $\tau^1 [X_{\mathcal{V}^1}]$  are chosen. Thus, the rank bounds hold also for any other rank in the claim.  $\square$



*Proof of The. 124.* We partition the edges into the set  $\mathcal{E}^1 = \{e \in \mathcal{E} : e \cup \tilde{\mathcal{V}} \neq \emptyset\}$  and  $\mathcal{E}^2 = \{e \in \mathcal{E} : e \cup \tilde{\mathcal{V}} = \emptyset\}$ . We then have

$$\langle \tau^{\mathcal{G}} \rangle [X_{\tilde{\mathcal{V}}}] = \left\langle \langle \{\tau^e [X_e] : e \in \mathcal{E}^1\} \rangle \left[ X_{\bigcup_{e \in \mathcal{E}^1} e} \right], \langle \{\tau^e [X_e] : e \in \mathcal{E}^2\} \rangle \left[ X_{\bigcup_{e \in \mathcal{E}^2} e} \right] \right\rangle [X_{\tilde{\mathcal{V}}}] \quad (47)$$

By an iterative application of Lem. 38 when including the cores to  $e \in \mathcal{E}^1$  after each other to the contraction, we get the bound

$$\text{rank} \left( \langle \{\tau^e [X_e] : e \in \mathcal{E}^1\} \rangle \left[ X_{\bigcup_{e \in \mathcal{E}^1} e} \right] \right) \leq \prod_{e \in \mathcal{E}^1} \text{rank}(\tau^e).$$

With the decomposition (47) and Lem. 39 we then arrive at the claim

$$\langle \tau^{\mathcal{G}} \rangle [X_{\tilde{\mathcal{V}}}] \leq \prod_{e \in \mathcal{E}^1} \text{rank}(\tau^e).$$

Since the applied lemmata hold also for the restricted CP ranks in the claim, the derived bound is also for those valid.  $\square$

**Example 25** (Composition of formulas with connectives). *For any formula  $f$  we have  $1 - f = \neg f$ . The CP rank bound brings an increase by at most factor 2 when taking the contraction with  $\beta^\neg$  which has slice sparsity of 2. This is not optimal, since  $\neg f$  has at most an absolute slice sparsity increase of 1.*

*For any formulas  $f$  and  $h$  we have  $f \cdot h = f \wedge h$ . Here the CP rank bounds on contractions can also be further tightened.*

**Example 26** (Distributions of independent variables). *Independence means factorization, conditional independence means sum over factorizations. Again, the  $\ell_0$  norm is bounded by the product of the  $\ell_0$  norm of the factors.*

### 18.3 Sparse Encoding of Functions

We now state that the basis CP rank of basis encodings is equal to the cardinality of the domain. The basis CP format can therefore not provide a sparse representation when the factored system contains many categorical variables.

**Theorem 125.** *For any function*

$$q : \bigtimes_{k \in [d]} [m_k] \rightarrow \bigtimes_{l \in [r]} [m_l]$$

*between factored systems we have*

$$\text{rank}^{\text{bas}}(\beta^q) = \prod_{k \in [d]} m_k.$$

*Proof.* The bound follows from The. 118, using that  $\ell_0(\beta^q) = \prod_{k \in [d]} m_k$ .  $\square$

Let us further provide a construction scheme to find a basis CP decomposition of  $\beta^q$  of size  $\prod_{k \in [d]} m_k$ . We notice that

$$\beta^q [Y_{[p]}, X_{[d]}] = \sum_{x_{[d]} \in \bigtimes_{k \in [d]} [m_k]} \epsilon_{x_{[d]}} [X_{[d]}] \otimes \epsilon_{q(X_{[d]}=x_{[d]})} [Y_{[p]}].$$

We build for  $k \in [d]$  decomposition variables  $I_{[d]}$  with  $n_k = m_k$  and define leg cores

$$\rho^k [X_k, I_{[d]}] = \delta [X_k, I_k]$$

and for  $\tilde{k} \in [r]$  and  $i_{[d]}$

$$\rho^{\tilde{k}} [Y_{\tilde{k}}, I_{[d]} = i_{[d]}] = \epsilon_{q_{\tilde{k}}[X_{[d]}=i_{[d]}]} Y_{\tilde{k}}.$$

We then have with a trivial scalar core

$$\beta^q [Y_{[p]}, X_{[d]}] = \left\langle \{I_{[d]}\} \cup \{\rho^k [X_k, I_{[d]}] : k \in [d]\} \cup \{\rho^{\tilde{k}} [Y_{\tilde{k}}, I_{[d]}] : \tilde{k} \in [r]\} \right\rangle [Y_{[p]}, X_{[d]}].$$

This is a basis CP decomposition of size  $\prod_{k \in [d]} m_k$ .

In combination with The. 122, The. 125 also provides bounds on all other CP ranks defined in Def. 89. This is obvious, since basis leg slices are the most restrictive properties compared with boolean, directed or basis+.

We restate The. 110 as a basis CP decomposition bound.

**Theorem 126.** *Let  $q$  and be a function between factored systems*

$$q : [m] \rightarrow \bigtimes_{k \in [d]} [m_k]$$

and  $q_k$  as in The. 110. Then  $\beta^q [X_{[d]}, X]$  has a basis CP decomposition with decomposition index  $X$ , trivial slices  $\mathbb{I}[X]$  leg vectors  $\beta^{q_k} [X_k, X]$ , that is

$$\beta^q [X, X_{[d]}] = \langle \{\mathbb{I}[X]\} \cup \{\beta^{q_k} [X_k, X] : k \in [d]\} \rangle [X_{[d]}]$$

*Proof.* The claimed decomposition directly follows from The. 110, since the trivial scalar core  $\lambda[X] = \mathbb{I}[X]$  does not influence the contraction and can be omitted.  $\square$

Basis CP decompositions can be constructed by understanding the variable  $O_{\text{in}}$  of the basis encoding of a function  $q : \mathcal{U}^{\text{in}} \rightarrow \mathcal{U}^{\text{out}}$  as the slice selection variable.

**Example 27** (Empirical distributions, see The. 23). *Let there be a data map*

$$D : [m] \rightarrow \bigtimes_{k \in [d]} [m_k] .$$

We can use The. 126 to find a tensor network representation for  $\beta^D$  as

$$\beta^D [X_{[d]}, J] = \langle \{\beta^{D_k} [X_k, J] : k \in [d]\} \rangle [X_{[d]}, J] .$$

This representation is a basis CP decomposition, when adding trivial scalar core. This provides also a basis CP decomposition for the empirical distribution, since normalization can be done by setting a slice core to  $\frac{1}{m} \mathbb{I}[J]$ .

**Example 28.** *Exponential families The statistic has a CP decomposition with rank by the cardinality of states, that is*

$$\beta^S [Y_{[p]}, X_{[d]}] = \langle \{\beta^{S_l} [Y_l, X_{[d]}] : l \in [p]\} \rangle [Y_{[p]}, X_{[d]}] .$$

While The. 125 and The. 126 provide CP rank bounds based on the domain factored system, we can also show in the next theorem a bound using the structure of the image.

**Theorem 127.** *Let  $q : \mathcal{U}^{\text{in}} \rightarrow \mathcal{U}^{\text{out}}$  be an arbitrary function and let us consider for each  $y \in \text{im}(q)$  the indicator*

$$\mathbb{I}^{q=y} [O_{\text{in}}] = \begin{cases} 1 & \text{if } q(O_{\text{in}}) = y \\ 0 & \text{else} . \end{cases}$$

*The basis+ rank of the basis encoding of  $q$  then obeys the bound*

$$\text{rank}(\beta^q) \leq \sum_{y \in \text{im}(q)} \text{rank}(\mathbb{I}^{q=y}) .$$

*The bound still holds, when we replace on both sides  $\text{rank}(\cdot)$  by  $\text{rank}^{\text{bin}}(\cdot)$ , by  $\text{rank}^{\text{bas}}(\cdot)$  or by  $\text{rank}^{\text{bas}+}(\cdot)$ .*

*Proof.* We have

$$\beta^q [O_{\text{out}}, O_{\text{in}}] = \sum_{y \in \text{im}(q)} \epsilon_{I(y)} [O_{\text{out}}] \otimes \mathbb{I}^{q=y} [O_{\text{in}}] .$$

For any  $y \in \text{im}(q)$  it is obvious that

$$\text{rank}(\epsilon_{I(y)} [O_{\text{out}}] \otimes \mathbb{I}^{q=y} [O_{\text{in}}]) = \text{rank}(\mathbb{I}^{q=y} [O_{\text{in}}]) ,$$

which also holds true for the other bounds in the claim. We then use the summation bound of The. 123 to get

$$\begin{aligned} \text{rank}(\beta^q [O_{\text{out}}, O_{\text{in}}]) &\leq \sum_{y \in \text{im}(q)} \text{rank}(\epsilon_{I(y)} [O_{\text{out}}] \otimes \mathbb{I}^{q=y} [O_{\text{in}}]) \\ &\leq \sum_{y \in \text{im}(q)} \text{rank}(\mathbb{I}^{q=y} [O_{\text{in}}]) . \end{aligned}$$

Again, the bound still hold for the other ranks in the claim.  $\square$

The above claim still holds when replacing  $\text{rank}^{\text{bas}+}(\cdot)$  with the ranks  $\text{rank}^{\text{bas}}(\cdot)$  or  $\text{rank}^{\text{bin}}(\cdot)$ . For the rank  $\text{rank}^{\text{bas}}(\cdot)$  it leads to the bound of The. 125, since summing the number of non zero coordinators of the indicators is the cardinality of the domain.

**Example 29** (Propositional formulas). *Let us now illustrate how the above representation scheme can be leveraged for the sparse representation of propositional formulas. For an arbitrary propositional formula  $f$  we have  $\text{im}(f) \subset \{0, 1\}$  and the indicators*

$$\mathbb{I}^{f=1}[X_{[d]}] = f[X_{[d]}] \quad \text{and} \quad \mathbb{I}^{f=0}[X_{[d]}] = \neg f[X_{[d]}] = \mathbb{I}[X_{[d]}] - f[X_{[d]}] .$$

For the conjunction  $\wedge[X_0, X_1] = X_0 \wedge X_1$  we have

$$\beta^\wedge[X_0, X_1] = \epsilon_1[Y_\wedge] \otimes \epsilon_{1,1}[X_0, X_1] + \epsilon_0[Y_\wedge] \otimes (\mathbb{I}[X_0, X_1] - \epsilon_{1,1}[X_0, X_1])$$

and thus

$$\text{rank}^{\text{bas}+}(\beta^\wedge) \leq 3$$

while  $\text{rank}^{\text{bas}}(\beta^\wedge) = 4$ .

We can even generalize this observation to  $d$ -ary conjunctions  $\wedge[X_{[d]}] = X_0 \wedge \dots \wedge X_{d-1}$  (see Remark 3)

$$\wedge[X_{[d]}] = \bigotimes_{k \in [d]} \epsilon_1[X_k] \quad \text{and} \quad \neg \wedge[X_{[d]}] = \mathbb{I}[X_{[d]}] - \bigotimes_{k \in [d]} \epsilon_1[X_k]$$

and thus

$$\beta^\wedge[X_{[d]}] = \epsilon_1[Y_\wedge] \otimes \left( \bigotimes_{k \in [d]} \epsilon_1[X_k] \right) + \epsilon_0[Y_\wedge] \otimes \left( \mathbb{I}[X_{[d]}] - \bigotimes_{k \in [d]} \epsilon_1[X_k] \right)$$

Thus, while the basis CP rank is  $\text{rank}^{\text{bas}}(\beta^\wedge) = 2^d$ , the basis+ rank is bounded by 3, independently of  $d$ .

## 18.4 Optimization of sparse tensors

Let us now study the problem of searching for the maximal coordinate in a tensor represented by a monomial decomposition. Given a tensor  $\tau[X_{[d]}]$  we state this as the problem:

$$\text{argmax}_{x_{[d]} \in \times_{k \in [d]} [m_k]} \tau[X_{[d]} = x_{[d]}] \quad (\text{P}_\tau^{\text{max}})$$

Problem  $\text{P}_\tau^{\text{max}}$  can be reformulated as optimization over the standard simplex

$$\mathcal{M}_\wedge = \text{conv} \left( \epsilon_{x_{[d]}}[X_{[d]}] : x_{[d]} \in \times_{k \in [d]} [m_k] \right)$$

as

$$\text{argmax}_{\mu[L_{[d]}] \in \mathcal{M}_\wedge} \langle \mu, \tau \rangle [\emptyset] .$$

**Example 30** (Mode search in exponential families). *Given a statistic  $S$ , a canonical parameter  $\theta$  and a boolean base measure  $\nu$ , the mode search problem for the member  $\mathbb{P}^{S, \theta, \nu}$  of the exponential family  $\Gamma^{S, \nu}$  is*

$$\max_{x_{[d]} \in \times_{k \in [d]} [2] : \nu[X_{[d]} = x_{[d]}] = 1} \langle \sigma^S[X_{[d]} = x_{[d]}, L], \theta[L] \rangle [\emptyset] = \max_{\mu \in \mathcal{M}_{S, \nu}} \langle \mu[L], \theta[L] \rangle [\emptyset] .$$

Such mode search problems have appeared as generic MAP queries (see Chapter 6). In Chapter 12 we have discussed them for the specific cases of hybrid logic networks and grafting proposal distributions.

### 18.4.1 Unconstrained Binary Optimization

For leg dimensions  $m_k = 2$ , Problem  $\text{P}_\tau^{\text{max}}$  is known as the unconstrained binary optimization. Problem  $\text{P}_\tau^{\text{max}}$  is a Higher-Order Unconstrained Binary Optimization (HUBO), when  $\tau[X_{[d]}]$  has a when  $\tau$  has a monomial decomposition (see Def. 91) with  $|A^i| \leq r$  for all  $i \in [n]$ , that is when  $\text{rank}^r(\tau) < \infty$ .

**Definition 92.** Let  $\tau [X_{[d]}]$  be a tensor with a monomial decomposition  $\{(\lambda^i, A^i, x_{A^i}^i) : i \in [n]\}$ , where  $\max_{i \in [n]} |A^i| = r$ . Se then call Problem  $P_\tau^{\max}$  a  $r$ -Order Unconstrained Binary Optimization (HUBO), which we denote as

$$\operatorname{argmax}_{x_{[d]} \in \times_{k \in [d]} [2]} \sum_{i \in [n]} \lambda^i \left\langle \epsilon_{x_{A^i}^i} [X_{A^i}] \right\rangle [X_{[d]} = x_{[d]}] . \quad (P_\tau^{\text{HUBO}})$$

**Remark 26** (Leg dimensions larger than 2). We demanded leg dimensions  $m_k = 2$  to have boolean valued variables  $X_k$ , which is required to connect with the formalism of binary optimization. Categorical variables with larger dimensions can be represented by atomization variables, which are created by contractions with categorical constraint tensors (see Sect. 8.3.2).

The sparsity rank<sup>r</sup> ( $\tau$ ) is the minimal number of monomials, for which a weighted sum is equal to  $\tau$ . Thus we interpret Problem  $P_\tau^{\text{HUBO}}$  as searching for the maximum in a polynomial consistent of rank<sup>r</sup> ( $\tau$ ) monomial terms.

Problem  $P_\tau^{\text{HUBO}}$  is called Quadratic Unconstrained Binary Optimization problems, if  $r = 2$ . We can transform certain Higher-Order Unconstrained Binary Optimization (HUBO) problems into Quadratic Unconstrained Binary Optimization (QUBO) problems by introducing auxiliary variables. An example of such an transform is provided by the next lemma.

**Lemma 40.** For any  $x_0, \dots, x_{d-1} \in [2]$  and  $A \subset [d]$  we have

$$\left( \prod_{k \in A} x_k \right) \left( \prod_{k \notin A} (1 - x_k) \right) = \max_{z \in [2]} z \cdot 2 \cdot \left( \sum_{k \in A} x_k - |A| - \sum_{k \notin A} x_k + \frac{1}{2} \right) .$$

*Proof.* Only if  $x_k = 1$  for  $k \in A$  and  $x_k = 0$  else we have

$$\left( \sum_{k \in A} x_k - |A| - \sum_{k \notin A} x_k + \frac{1}{2} \right) \geq 0 .$$

In this case the maximum is taken for  $z = 1$  and we have

$$\max_{z \in [2]} z \cdot 2 \cdot \left( \sum_{k \in A} x_k - |A| - \sum_{k \notin A} x_k + \frac{1}{2} \right) = 1 = \left( \prod_{k \in A} x_k \right) \left( \prod_{k \notin A} (1 - x_k) \right) .$$

In all other cases, the maximum is taken for  $z = 0$  and thus vanishes, that is

$$\max_{z \in [2]} z \cdot 2 \cdot \left( \sum_{k \in A} x_k - |A| - \sum_{k \notin A} x_k + \frac{1}{2} \right) = 0 = \left( \prod_{k \in A} x_k \right) \left( \prod_{k \notin A} (1 - x_k) \right) .$$

Thus, the claim holds in all cases.  $\square$

#### 18.4.2 Integer Linear Programming

Let us now show how optimization problems can be represented as linear programming problems. To this end, we understand each index tuple  $x_{[d]} \in \times_{k \in [d]} [m_k]$  as a vector  $v_{x_{[d]}} [L] \in \mathbb{R}^d$  with coordinates

$$v_{x_{[d]}} [L = k] = x_k .$$

**Definition 93.** The integer linear program (ILP) of  $M[J, L] \in \mathbb{R}^{n \times d}$ ,  $b[J] \in \mathbb{R}^n$  and  $c \in \mathbb{R}^d$  is the problem

$$\operatorname{argmax}_{x_{[d]} \in \times_{k \in [d]} [m_k]} \langle c[L], v_{x_{[d]}} [L] \rangle [\emptyset] \quad \text{subject to} \quad \langle M[J, L], v_{x_{[d]}} [L] \rangle [\emptyset] \prec b[J] , \quad (P_{c, M, b}^{\text{ILP}})$$

where by  $\prec$  we denote partial ordering of tensors (see Def. 81).

We now show that any binary optimization problem of a tensor can be transformed into a integer linear program, given a monomial decomposition of the tensor  $\tau [X_{[d]}]$  by  $\mathcal{M} = \{(\lambda^i, A^i, x_{A^i}^i) : i \in [n]\}$ . For this we choose state indices by vectors

$$y_{[d+n]} = x_0, \dots, x_{d-1}, z_0, \dots, z_{n-1} \in \left( \times_{k \in [d]} [2] \right) \times \left( \times_{i \in [n]} [2] \right) ,$$

that is we added for each monomial an index  $z_i$ , which will represent the evaluations of the respective monomial.

We furthermore define a vector  $c^{\mathcal{M}}[L]$ , where  $L$  takes values in  $[d+n]$ , as

$$c^{\mathcal{M}}[L=l] = \begin{cases} \lambda^{l-d} & \text{if } l > d \\ 0 & \text{else} \end{cases}. \quad (48)$$

To construct a matrix  $M[J, L]$  and a vector  $b[J]$  to the monomial decomposition  $\mathcal{M}$ , we now introduce a variable  $J$  enumerating linear inequalities, which takes values in  $[m]$ , where

$$m = \sum_{i \in [n]} (|A^i| + 1).$$

We define for each  $i \in [n]$  an auxiliary number

$$m_i = \sum_{\tilde{i}=0}^i (|A^{\tilde{i}}| + 1)$$

and further enumerate the set  $A^i$  by a function  $I : [|A^i|] \rightarrow A^i$ .

We then construct a matrix  $M^{\mathcal{M}}[J, L]$ , where for  $l \in [d+n]$ ,  $i \in [n]$  and  $j \in [|A^i|]$  we have

$$M^{\mathcal{M}}[J = m_i + j, L = l] = \begin{cases} 1 - 2 \cdot x_{I(j)}^i & \text{if } j < |A^i|, j = l \text{ and } l = I(j) \\ 1 & \text{if } j < |A^i| \text{ and } l = d + i \\ -x_{I(j)}^i & \text{if } j = |A^i| \text{ and } l = I(j) \\ -1 & \text{if } j = |A^i| \text{ and } l = d + i \\ 0 & \text{else} \end{cases}. \quad (49)$$

Similarly, we define  $b^{\mathcal{M}}[J]$  as the vector which nonvanishing coordinates are for  $i \in [n]$  at

$$b^{\mathcal{M}}[J = m_i + j] = \begin{cases} 1 - x_{I(j)}^i & \text{if } j < |A^i| \\ -1 + |\{k \in A^i : x_k^i = 1\}| & \text{if } j = |A^i| \end{cases}. \quad (50)$$

Informally, we pose for each tuple  $(\lambda, A, x_A)$   $|A| + 1$  linear equations. The first  $|A|$  enforce, that the slice representing variable  $z$  is zero once a leg is 0. The last enforces that the slice representing variable is 1. We prove this claim more formally in the next theorem.

**Theorem 128.** *Given a monomial decomposition  $\mathcal{M} = \{(\lambda^i, A^i, x_{A^i}^i) : i \in [n]\}$  of a tensor  $\tau$ , let  $y^{ILP, \mathcal{M}}$  be a solution of the integer linear program defined by the matrix and vectors in equations (48), (49) and (50). Then we have*

$$y^{ILP, \mathcal{M}}|_{[d]} \in \operatorname{argmax}_{x_{[d]} \in \times_{k \in [d]} [2]} \tau[X_{[d]} = x_{[d]}].$$

where by  $y^{ILP, \mathcal{M}}|_{[d]}$  we denote the restriction of the index tuple  $y^{ILP, \mathcal{M}}$  to the first  $d$ .

*Proof.* We show that the linear constraints by

$$\langle M^{\mathcal{M}}[J, L], v_{y_{[d+n]}}[L] \rangle [\emptyset] \prec b^{\mathcal{M}}[J]$$

are satisfied for a vector  $y_{[d+n]} = (x_{[d]}, z_{[n]})$ , if and only if for all  $i \in [n]$  the product constraints

$$z_i = \left( \prod_{k \in A^i, x_k^i = 0} (1 - x_k) \right) \cdot \left( \prod_{k \in A^i, x_k^i = 1} x_k \right) \quad (51)$$

hold. We will see, that the linear constraints where  $J$  takes indices in  $m_i + [|A^i|]$  are equivalent to the upper bound on  $z_i$  and the constraint to  $J = m_i + |A^i|$  is equivalent to an lower bound on  $z_i$ . To show the upper bound, we notice that for any  $j \in [|A^i|]$  the constraint  $J = m_i + j$  is

$$z_i \leq \begin{cases} x_{I(j)} & \text{if } x_{I(j)}^i = 1 \\ (1 - x_{I(j)}) & \text{if } x_{I(j)}^i = 0 \end{cases}.$$

Thus, whenever a factor on the right side of (51) is 0, we have  $z_i = 0$  if the respective constraint is satisfied. We conclude, that

$$z_i \leq \left( \prod_{k \in A^i, x_k^i=0} (1 - x_k) \right) \cdot \left( \prod_{k \in A^i, x_k^i=1} x_k \right).$$

To show the lower bound, we have the constraint to  $J = m_i + |A^i|$  by

$$z_i \geq 1 - \left( \sum_{k \in A^i, x_k^i=0} x_k \right) + \left( \sum_{k \in A^i, x_k^i=1} (x_k - 1) \right).$$

The right side of this inequality is 1, if and only if all factors on the right side of (51) are 1, and less or equal to 0 else. Thus, whenever this constraint is satisfied, we have

$$z_i \geq \left( \prod_{k \in A^i, x_k^i=0} (1 - x_k) \right) \cdot \left( \prod_{k \in A^i, x_k^i=1} x_k \right).$$

In summary, the equation (51) holds, if and only if the constraints where  $J$  takes indices in  $m_i + [|A^i| + 1]$  are satisfied.

This characterization of the constraints implies, that for any  $x_{[d]} \in \times_{k \in [d]} [m_k]$  there is exactly one feasible index  $y_{[d+n]}$  with  $(y_{[d+n]})_{[d]} = x_{[d]}$ , and the objective takes for this index the value

$$\begin{aligned} \langle c[L], v_{y_{[d+n]}}[L] \rangle [\emptyset] &= \sum_{i \in [n]} \lambda^i \cdot z_i \\ &= \sum_{i \in [n]} \lambda^i \cdot \left( \prod_{k \in A^i, x_k^i=0} (1 - x_k^i) \right) \cdot \left( \prod_{k \in A^i, x_k^i=1} x_k^i \right) \\ &= \tau [X_{[d]} = x_{[d]}]. \end{aligned}$$

Therefore, any solution of the ILP reduced to the first  $d$  indices corresponding with the axis of  $\tau$ , is a solution of the binary optimization problem to  $\tau$ .  $\square$

In order to achieve a sparse linear program it is beneficial to use a monomial decomposition with small order and rank. Beside this sparsity, the matrix  $M^{\mathcal{M}}[J, L]$  is often  $\ell_0$ -sparse, and has thus an efficient representation in a basis CP format. More precisely we have by the above construction

$$\ell_0(M^{\mathcal{M}}[J, L]) \leq \sum_{i \in [n]} 3 \cdot |A^i| + 1.$$

## 19 Tensor Approximation

Often reasoning requires the execution of demanding contractions of tensors networks, or combinatorial search of maximum coordinates. We in this chapter investigate methods, to replace hard to be sampled tensor networks by approximating tensor networks, which then serve as a proxy in inference tasks.

### 19.1 Selection tensor networks for CP decompositions

In this section, we show that the set of tensors representable in specific CP formats coincides with the expressivity of tailored selection architectures (see Chapter 10). We first define a basis+ CP selecting tensor and then show its decomposition into a formula selecting neural network.

**Definition 94.** Given a set of categorical variables  $X_{[d]}$  with  $m = \max_{k \in [d]} m_k$ , a CP selecting tensor of maximal cardinality  $r$  is the tensor

$$\mathcal{H}_{\wedge, d, r}(X_{[d]}, L_{0,0}, \dots, L_{r-1,0}, L_{0,1}, \dots, L_{r-1,1})$$

with dimensions

$$p_{s,0} = m + 1, p_{s,1} = d \quad \text{for } s \in [r]$$

and coordinates

$$\begin{aligned} \mathcal{H}_{\wedge,d,r} (X_{[d]} = x_{[d]}, L_{0,0} = l_{0,0}, \dots, L_{r-1,0} = l_{r-1,0}, L_{0,1} = l_{0,1}, \dots, L_{r-1,1} = l_{r-1,1}) \\ = \begin{cases} 1 & \text{if } \forall k \in [d], s \in [r] : (l_{s,1} = k \wedge l_{s,0} \neq m) \Rightarrow (l_{s,0} = x_k) \\ 0 & \text{else} \end{cases} \end{aligned}$$

Intuitively, the selection variables  $L_{s,1}$  of  $\mathcal{H}_{\wedge,d,r}$  select a variable out of  $[d]$  to be included in  $A$  and the selection variables  $L_{s,0}$  select a corresponding state to that variable. As in Chapter 10 we refer to variables  $L_{s,1}$  as variable selectors and  $L_{s,0}$  as state selectors. When  $L_{s,0} = m$ , the slice is left trivial, that is the selected variable is effectively not included in  $A$ . This then allows to also represent slices where  $|A| < r$ . We in the following prove this more formally.

**Theorem 129.** *Let the non-vanishing indices of  $\theta [L_{[r] \times [m+1]}]$  denote by  $\{l_{[r] \times [m+1]}^i : i \in [n]\}$ . Let further  $\mathcal{M} \subset [n]$  be the set of agreeing selection indices, that is*

$$\mathcal{M} = \{i : i \in [n], \forall s, \tilde{s} \in [r] : (l_{s,1}^i = l_{\tilde{s},1}^i \wedge l_{s,1}^i \neq m \wedge l_{\tilde{s},1}^i \neq m) \Rightarrow (l_{s,1}^i = l_{\tilde{s},1}^i)\}$$

Then

$$\begin{aligned} \langle \mathcal{H}_{\wedge,d,r} (X_{[d]}, L_{[r] \times [m+1]}), \theta [L_{[r] \times [m+1]}] \rangle [X_{[d]}] \\ = \sum_{i \in \mathcal{M}} \theta [L_{[r] \times [m+1]} = l_{[r] \times [m+1]}^i] \langle \{\epsilon_{l_{s,0}} [X_{l_{s,1}}] : s \in [r], l_{s,0} \neq m\} \rangle [X_{[d]}] \end{aligned}$$

Further we have

$$\text{rank}^{\text{bas}+} (\langle \mathcal{H}_{\wedge,d,r} (X_{[d]}, L_{[r] \times [m+1]}), \theta [L_{[r] \times [m+1]}] \rangle [X_{[d]}]) \leq \ell_0(\theta) .$$

*Proof.* Let us notice, that whenever  $i \notin \mathcal{M}$  then

$$\mathcal{H}_{\wedge,d,r} (X_{[d]}, L_{[r] \times [m+1]} = l_{[r] \times [m+1]}^i) = 0 [X_{[d]}] ,$$

since the condition for non-vanishing coordinates

$$\forall k \in [d], s \in [r] : (l_{s,1} = k \wedge l_{s,0} \neq m) \Rightarrow (l_{s,0} = x_k)$$

in Def. 94 cannot be satisfied for any  $x_{[d]}$ . For  $i \notin \mathcal{M}$  we have

$$\mathcal{H}_{\wedge,d,r} (X_{[d]}, L_{[r] \times [m+1]} = l_{[r] \times [m+1]}^i) = \langle \{\epsilon_{l_{s,0}} [X_{l_{s,1}}] : s \in [r], l_{s,0} \neq m\} \rangle [X_{[d]}] .$$

We now use these insights and linearity of contraction to get

$$\begin{aligned} \langle \mathcal{H}_{\wedge,d,r} (X_{[d]}, L_{[r] \times [m+1]}), \theta [L_{[r] \times [m+1]}] \rangle [X_{[d]}] \\ = \sum_{i \in [n]} \theta [L_{[r] \times [m+1]} = l_{[r] \times [m+1]}^i] \mathcal{H}_{\wedge,d,r} (X_{[d]}, L_{[r] \times [m+1]} = l_{[r] \times [m+1]}^i) \\ = \sum_{i \in \mathcal{M}} \theta [L_{[r] \times [m+1]} = l_{[r] \times [m+1]}^i] \langle \{\epsilon_{l_{s,0}} [X_{l_{s,1}}] : s \in [r], l_{s,0} \neq m\} \rangle [X_{[d]}] \end{aligned}$$

This shows the decomposition claim. We notice, that this is a basis+ CP decomposition of size  $\mathcal{M}$  and thus

$$\text{rank}^{\text{bas}+} (\langle \mathcal{H}_{\wedge,d,r} (X_{[d]}, L_{[r] \times [m+1]}), \theta [L_{[r] \times [m+1]}] \rangle [X_{[d]}]) \leq |\mathcal{M}| \leq \ell_0(\theta) .$$

□

Towards a practice usage of this representation scheme for basis+ CP decompositions, let us show that the CP selecting tensor coincides with a formula selecting network.

**Lemma 41.** *The CP selection tensor coincides with a formula selecting neural network with neurons (see Figure 19.1):*

- unary state selecting neurons enumerated by  $s$ , selecting one of the  $X_{[d]}$  with the variable  $L_{s,1}$  and selecting a state, extended by a possible choice of trivial legs  $\mathbb{I}$
- $r$ -ary output neuron fixed to the  $\wedge$  connective.

*Proof.* This can be easily checked on each coordinate. □

If  $m_k = 2$  for  $k \in [d]$ , the state selector chooses between the connectives  $\{\neg, \text{Id}, \text{True}\}$ . We can in that case understand each slice by a logical term.

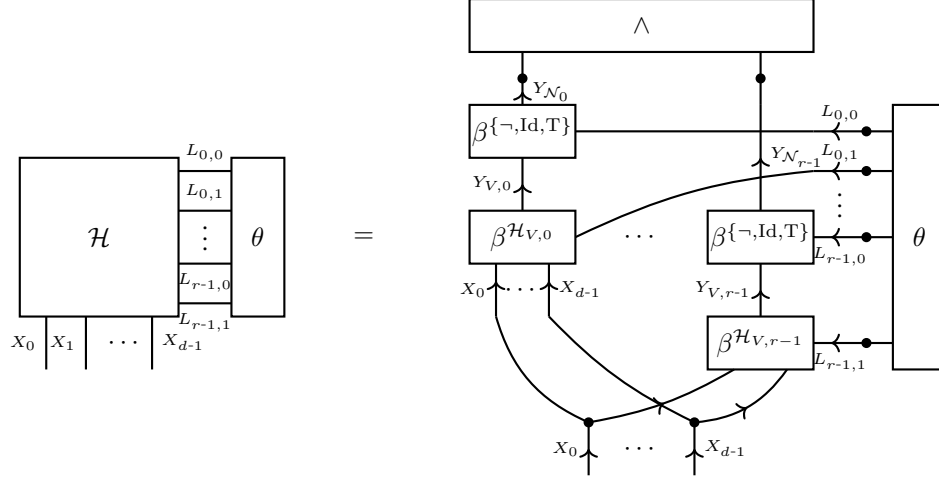


Figure 42: Representation of a basis+ Tensor by the contraction of a parameter tensor  $\theta$  with a CP selecting architecture  $\mathcal{H}$ , which has a decomposition as a formula selecting neural network (see Lem. 41). The nonzero coordinates of  $\theta$  represent slices of the CP decomposition.

### 19.1.1 Applications

One application is as a parametrization scheme in the approximation of a tensor by a slice-sparse tensor, see Chapter 19. The approximated parameter can then be used as a proxy energy to be maximized. When choosing  $r = 2$ , the approximating tensor contains only quadratic slices, which then poses a QUBO problem.

**Remark 27** (Extension to arbitrary CP formats). *Select at each input neuron a specific leg. For finite number of legs, as it is the case in the binary, basis and basis+ formats, we can enumerate all possibilities by the selection variable. For the basis+ format, in case of binary leg dimensions, we here exemplified the approach, by enumerating the three possibilities  $\epsilon_0, \epsilon_1, \mathbb{I}[1]$ . This approach, however, fails as a generic representation of the directed format, since the directed legs are continuous and there therefore are infinite choosable legs.*

### 19.2 Approximation of Energy tensors

The Hilbert-Schmidt norm of a tensor is the contraction of the coordinatewise transform with the square function

$$\|\tau[X_{[d]}]\|_2 = \sqrt{\langle (\tau[X_{[d]}])^2 | \emptyset \rangle}.$$

Approximation involving a selection architecture  $\mathcal{H}$  is the problem

$$\operatorname{argmin}_{\theta \in \Gamma^g} \|\phi - \langle \sigma^{\mathcal{H}}, \theta \rangle [X_{[d]}]\|^2.$$

Direct approximation is then the choice of the minterm statistic

$$\operatorname{argmin}_{\theta \in \Gamma^g} \|\phi[X_{[d]}] - \theta[X_{[d]}\|^2.$$

In a tensor network diagram we depict this as

$$\operatorname{argmin}_{\theta \in \Gamma^g} \left\| \begin{array}{c} \sigma^{\mathcal{F}} \\ \vdots \\ \theta \end{array} \begin{array}{c} L_{n-1} \\ \vdots \\ L_0 \end{array} \begin{array}{c} X_0 \\ \vdots \\ X_{d-1} \end{array} - \begin{array}{c} Y \\ \vdots \\ X_0 \\ \vdots \\ X_{d-1} \end{array} \right\|^2$$

**Example 31** (Approximate based on a slice sparsity selecting architecture). *Use a term selecting neural network (conjunction neuron on  $d$  unary neurons selecting a variable and Id,  $\neg$ , True as connective selector. Demand the parameter tensor  $\theta$  to be in a basis CP format, then each slice of the parameter tensor corresponds with the slice of the energy. The use the approximation for MAP search. Same construction possible for probability tensors, but often more involved to instantiate them as tensor network.*



### 19.3 Transformation of Maximum Search to Risk Minimization

By the squares risk trick, maximum coordinate searches involving contractions with boolean tensors can be turned into squares risk minimization problems. This trick can be applied in MAP inference of MLN and the proposal distribution.

#### 19.3.1 Weighted Squares Loss Trick

**Lemma 42.** *Let  $\tau$  be a boolean tensor, that is  $\text{im}(\tau) \subset \{0, 1\}$ . Then*

$$\tau[X_{[d]}] = \mathbb{I}[X_{[d]}] - (\tau[X_{[d]}] - \mathbb{I}[X_{[d]}])^2$$

where  $\mathbb{I}$  is a tensor with same shape as  $\tau$  and all coordinates being 1.

*Proof.* Since for each  $x_{[d]} \in \times_{k \in [d]} [m_k]$  we have  $\tau[X_{[d]} = x_{[d]}] \in \{0, 1\}$ , it holds that

$$\tau[X_{[d]} = x_{[d]}] = 1 - (\tau[X_{[d]} = x_{[d]}] - 1)^2$$

and thus in coordinatewise calculus

$$\tau[X_{[d]}] = \mathbb{I}[X_{[d]}] - (\tau[X_{[d]}] - \mathbb{I}[X_{[d]}])^2.$$

□

We apply this property to reformulate optimization problems over boolean tensors into weighted least squares problems.

**Theorem 130** (Weighted Squares Loss Trick). *Let  $\Gamma$  be a set of boolean tensors in  $\bigotimes_{k \in [d]} \mathbb{R}^{m_k}$  and  $I \in \bigotimes_{k \in [d]} \mathbb{R}^{m_k}$  arbitrary. Then we have*

$$\text{argmax}_{\tau \in \Gamma} \langle I, \tau \rangle [\emptyset] = \text{argmin}_{\tau \in \Gamma} \langle I, (\tau[X_{[d]}] - \mathbb{I}[X_{[d]}])^2 \rangle [\emptyset] \quad (52)$$

*Proof.* Using the Lemma above,  $\tau$  is identical to  $\mathbb{I}[X_{[d]}] - (\tau[X_{[d]}] - \mathbb{I}[X_{[d]}])^2$  and we get

$$\langle I, \tau \rangle [\emptyset] = \langle I, \mathbb{I}[X_{[d]}] \rangle [\emptyset] - \langle I, (\tau[X_{[d]}] - \mathbb{I}[X_{[d]}])^2 \rangle [\emptyset]$$

Since the first term does not depend on  $\tau$ , it can be dropped in the maximization problem. The  $(-1)$  factor then turns the maximization into a minimization problem. □

The. 130 reformulates maximization of binary tensors with respect to an angle to another tensor into minimization of a squares risk. This squares risk trick is especially useful when combining it with a relaxation of  $\Gamma$  to differentially parametrizable sets, since then common squares risk solvers can be applied. We will call  $I$  in the The. 130 importance tensor, since it manipulates the relevance of each coordinate in the squares loss.

As a result, we interpret the objective

$$\langle I, (\tau[X_{[d]}] - \mathbb{I}[X_{[d]}])^2 \rangle [\emptyset]$$

as a weighted squares loss.

**Example 32** (Proposal distribution maxima). *The Problem  $\text{P}_{D, \mathcal{F}, \mathcal{H}}^{\text{grad}}$  of finding the maximal coordinate of the proposal distribution can thus be turned into*

$$\begin{aligned} & \text{argmax}_{l_{[n]}} \left\langle (\mathbb{P}^D - \tilde{\mathbb{P}}), \mathcal{H} \right\rangle [L_{[n]} = l_{[n]}] \\ &= \text{argmin}_{l_{[n]}} \left\langle (\mathbb{P}^D - \tilde{\mathbb{P}}), (\langle \mathcal{H}, \epsilon_{l_{[n]}} [L_{[n]}] \rangle [X_{[d]}] - \mathbb{I}[X_{[d]}])^2 \right\rangle [\emptyset]. \end{aligned}$$

#### 19.3.2 Problem of the trivial tensor

By the above we motivated least squares problems on the set of one-hot encoded states. One is tempted to extend this set to  $\Gamma^{\mathcal{S}^G, \mathbb{I}}$  for efficient solutions by alternating algorithms.

However, for any hypergraph  $\mathcal{G}$  we have  $\mathbb{I}[X_{[d]}] \in \Gamma^{\mathcal{S}^G, \mathbb{I}}$ . In many situations (e.g. disjoint model sets supported at positive data) the objective is more in favor at the trivial tensor than at the one-hot encoding. As a result, we do not solve the previously posed one-hot encoding problem, when allowing such an hypothesis embedding.

**Example 33** (Fitting a boolean tensor by a formula tensor). *Given a tensor  $\tau$ , we want to find a formula  $f \in \mathcal{F}$  such that it approximates  $\tau$ .*

*If  $\tau$  is a binary tensor, we understand it as a formula and want to find an  $f$  such that its number of worlds is maximal, that is solve the problem*

$$\operatorname{argmax}_{f \in \mathcal{F}} \langle f \Leftrightarrow \tau \rangle [\emptyset] .$$

*We can use the squares risk trick and get an equivalent problem*

$$\operatorname{argmin}_{f \in \mathcal{F}} \| \langle f \Leftrightarrow \tau \rangle [X_{[d]}] - \mathbb{I} [X_{[d]}] \|^2 .$$

*We have since  $\tau$  and  $f$  are boolean*

$$\| \langle f \Leftrightarrow \tau \rangle [X_{[d]}] - \mathbb{I} [X_{[d]}] \|^2 = \| \langle f \rangle [X_{[d]}] - \tau [X_{[d]}] \|^2$$

*Now, when representing  $\mathcal{F}$  in a formula selecting architecture we have*

$$\operatorname{argmin}_{\theta \in \Gamma_1} \| \langle \mathcal{H}, \theta \rangle [X_{[d]}] - \tau [X_{[d]}] \|^2 .$$

*where  $\Gamma_1$  is the set of basis tensors.*

*When we extend  $\Gamma_1$  to a set including the trivial tensor  $\mathbb{I} [L]$ , when the formulas  $f$  are pairwise disjoint and  $\tau [X_{[d]}] = \mathbb{I} [X_{[d]}]$ , then the solution would be  $\mathbb{I} [L]$ .*

#### 19.4 Alternating Solution of Least Squares Problems

When the parameter tensor  $\theta$  is only restricted to have a decomposition as a tensor network on  $\mathcal{G}$ , we can iteratively update each core. The resulting algorithm is called Alternating Least Squares (ALS) (see Algorithm 11).

---

##### Algorithm 11 Alternating Least Squares (ALS)

---

**Require:** Target tensor  $Y[X_{[d]}]$ , hypergraph  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$  with  $[d] \subset \mathcal{V}$  specifying the approximation format

**Ensure:** Approximation of  $Y[X_{[d]}]$  by a tensor network  $\langle \{\tau^e [X_e] : e \in \mathcal{E}\} \rangle [X_{[d]}]$

---

**for**  $e \in \mathcal{E}$  **do**

    Set  $\tau^e [X_e]$  to a random element in  $\bigotimes_{k \in e} \mathbb{R}^{m_k}$

**end for**

**while** Stopping criterion is not met **do**

**for**  $e \in \mathcal{E}$  **do**

        Set  $\tau^e [X_e]$  to a solution of the local problem, that is

$$\tau^e [X_e] \leftarrow \operatorname{argmin}_{\tau^e [X_e]} \langle I, \langle \langle \mathcal{H}, \theta \rangle [X_{[d]}] - Y[X_{[d]}] \rangle^2 \rangle [\emptyset]$$

**end for**

**end while**

**return**  $\{\tau^e [X_e] : e \in \mathcal{E}\}$

---

##### 19.4.1 Choice of Representation Format

The choice of the hypergraph  $\mathcal{G}$  used for approximation bears a tradeoff between expressivity and complexity in sampling. Hidden variables, that is variables only present in  $\mathcal{G}$ , but not in the sensing matrix, increase the expressivity, especially when assigning large dimensions to them. When there are no hidden variables, the maximum of  $\theta$  can be found by maximum calibration through a message passing algorithm, since no hidden variable has to be marginalized.

#### 19.5 Regularization and Compressed Sensing

When regularizing the least squares problem by enforcing the sparsity of  $\theta$ , we arrive at the compressed sensing problem

$$\operatorname{argmin}_{\theta [L]} \ell_0(\theta) \quad \text{subject to} \quad \| \langle \sigma^S, \theta \rangle [X_{[d]}] - \phi [X_{[d]}] \|_2 \leq \eta \quad (53)$$

Here, the sensing matrix is the selection tensor.

**Example 34** (Formula fitting to an example). *Choosing the best formula fitting data (see Example 33) is the problem*

$$\operatorname{argmin}_{\theta[L]: \ell_0(\theta)=1} \left\| \langle I, \sigma^S, \theta \rangle [X_{[d]}] - Y \right\|_2 \quad (54)$$

where  $I$  has nonzero entries at marked coordinates and  $Y$  stores in Boolean coordinates whether the marked coordinates are positive or negative examples. *When the number of positive and negative examples are identical, we can linearly transform the objective to that of a grafting instance, where the current model is the empirical distribution of negative examples and the data consists of the positive examples.*

The sparse tensor solving the problem then has a small number of nonzero coordinates and the selection tensor can be restricted to those. As a consequence, inference can be performed more efficiently.

The algorithmic solution of these problems can be done by greedy algorithms, thresholding based algorithms or optimization based algorithms Foucart and Rauhut (2013).

Guarantees for the success of the algorithms depend on the properties of the sensing matrices. Here the sensing matrices are deterministic, since constructed as selection tensors, and concentration based approaches towards probabilistic bounds on these properties (see Goeßmann (2021)) are not applicable.

**Example 35** (Sensing matrix for propositional Formulas). *Let there be a set  $\mathcal{F}$  of formulas, then we have*

$$\langle \sigma^{\mathcal{F}} [X_{[d]}, L_{\text{in}}], \sigma^{\mathcal{F}} [X_{[d]}, L_{\text{out}}] \rangle [L_{\text{in}} = l_{\text{in}}, L_{\text{out}} = l_{\text{out}}] = \langle f_{l_{\text{in}}}, f_{l_{\text{out}}} \rangle [\emptyset].$$

*If the formulas have disjoint model sets then*

$$\langle \sigma^{\mathcal{F}} [X_{[d]}, L_{\text{in}}], \sigma^{\mathcal{F}} [X_{[d]}, L_{\text{in}}] \rangle [L_{\text{in}} = l_{\text{in}}, L_{\text{out}} = l_{\text{out}}] = \begin{cases} \langle f_{l_{\text{in}}} \rangle [\emptyset] & \text{if } l_{\text{in}} = l_{\text{out}} \\ 0 & \text{else} \end{cases}.$$

*In that case, the sensing matrix is a restricted isometry, in the sense that the norm of any mapped vector is its norm multiplied by a factor between the smallest and the largest  $\langle f_{l_{\text{in}}} \rangle [\emptyset]$ .*

**Example 36** (Sensing matrix for slice selection networks). *For the slice selection network*

$$\begin{aligned} & \langle \mathcal{H}_{\wedge, d, r} (X_{[d]}, L_{\text{in}}), \mathcal{H}_{\wedge, d, r} (X_{[d]}, L_{\text{out}}) \rangle [L_{\text{in}} = l_{\text{in}}, L_{\text{out}} = l_{\text{out}}] \\ &= \begin{cases} 0 & \text{if for a } \tilde{k} \in A^{l_{\text{in}}} \cap A^{l_{\text{out}}} \text{ we have } x_{\tilde{k}}^{l_{\text{in}}} \neq x_{\tilde{k}}^{l_{\text{out}}} \\ \prod_{\tilde{k} \notin A^{l_{\text{in}}} \cup A^{l_{\text{out}}}} m_{\tilde{k}} & \text{else} \end{cases}. \end{aligned}$$

*Given a fixed  $l_{\text{in}}$ , the maximum value in the respective slice is thus taken at  $l_{\text{in}} = l_{\text{out}}$ .*

## 19.6 Discussion and Outlook

The slice selection network described here have been tailored to the CP format.

We can extend the slice selection network to parametrize more general tensor network formats than the CP format. Here the graph structure of the format itself can be optimized, by the usage of variable selectors. The activation selectors, generalizing the connective selection, are then choices of specific cores in the format. Thus, each neuron represents an hypercore  $\tau^e [X_e]$ , where  $e \subset \mathcal{V}$  is selected by variable selectors and the values of the tensor  $\tau^e$  by activation selectors.

## 20 Message Passing

In this chapter we introduce local contraction passed along tensor clusters to calculate global contractions exactly or approximatively. These message passing schemes provide tradeoffs between efficiency increases and exactness of the global contraction.

We use the CP decompositions to investigate the asymptotic behavior of the message passing algorithms.

*The application of message passing schemata to calculate contractions are motivated by commutations of contractions. We first show this property and then provide message passing schemata.*

### 20.1 Commutation of Contractions

We show in the next theorem, that a contractions can be performed by contracting a subnetwork first and then further contracting the result with the rest.

**Theorem 131** (Commutativity of Contractions). *Let  $\tau^{\mathcal{G}}$  be a tensor network on a hypergraph  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ . Let us now split the  $\mathcal{G}$  into two graphs  $\mathcal{G}_1 = (\mathcal{V}^1, \mathcal{E}_1)$  and  $\mathcal{G}_2 = (\mathcal{V}^2, \mathcal{E}_2)$ , such that  $\mathcal{E}_1 \dot{\cup} \mathcal{E}_2 = \mathcal{E}$ ,  $\mathcal{V}^1 \cup \mathcal{V}^2 = \mathcal{V}$  and all nodes in  $\mathcal{V}^2$  are contained in an hyperedge of  $\mathcal{E}_2$ . We then have for any  $\tilde{\mathcal{V}} \subset \mathcal{V}$*

$$\langle \tau^{\mathcal{G}} \rangle [X_{\tilde{\mathcal{V}}}] = \left\langle \tau^{\mathcal{G}_1} [X_{\mathcal{V}^1}] \cup \{ \langle \tau^{\mathcal{G}_2} \rangle [X_{\mathcal{V}^2 \cap (\mathcal{V}^1 \cup \tilde{\mathcal{V}})}] \} \right\rangle [X_{\tilde{\mathcal{V}}}] .$$

*Proof.* For any index  $x_{\tilde{\mathcal{V}}}$  we show that

$$\langle \tau^{\mathcal{G}} \rangle [X_{\tilde{\mathcal{V}}} = x_{\tilde{\mathcal{V}}}] = \left\langle \tau^{\mathcal{G}_1} \cup \{ \langle \tau^{\mathcal{G}_2} \rangle [X_{\mathcal{V}^2 \cap (\mathcal{V}^1 \cup \tilde{\mathcal{V}})}] \} \right\rangle [X_{\tilde{\mathcal{V}}} = x_{\tilde{\mathcal{V}}}] .$$

By definition we have

$$\begin{aligned} \langle \tau^{\mathcal{G}} \rangle [X_{\tilde{\mathcal{V}}} = x_{\tilde{\mathcal{V}}}] &= \sum_{x_{\mathcal{V}/\tilde{\mathcal{V}}}} \prod_{e \in \mathcal{E}} \tau^e [X_e = x_e] \\ &= \sum_{x_{\mathcal{V}/\tilde{\mathcal{V}}}} \left( \prod_{e \in \mathcal{E}_1} \tau^e [X_e = x_e] \right) \cdot \left( \prod_{e \in \mathcal{E}_2} \tau^e [X_e = x_e] \right) \\ &= \sum_{x_{\mathcal{V}^1/\tilde{\mathcal{V}}}} \sum_{x_{\mathcal{V}^2/(\tilde{\mathcal{V}} \cup \mathcal{V}^1)}} \left( \prod_{e \in \mathcal{E}_1} \tau^e [X_e = x_e] \right) \cdot \left( \prod_{e \in \mathcal{E}_2} \tau^e [X_e = x_e] \right) \\ &= \sum_{x_{\mathcal{V}^1/\tilde{\mathcal{V}}}} \left( \prod_{e \in \mathcal{E}_1} \tau^e [X_e = x_e] \right) \cdot \left( \sum_{x_{\mathcal{V}^2/(\tilde{\mathcal{V}} \cup \mathcal{V}^1)}} \prod_{e \in \mathcal{E}_2} \tau^e [X_e = x_e] \right) . \end{aligned}$$

When contracting the variables  $X_{\mathcal{V}^2/(\tilde{\mathcal{V}} \cup \mathcal{V}^1)}$  on  $\tau^{\mathcal{G}_2}$ , the variables  $X_{\mathcal{V}^2 \cap (\tilde{\mathcal{V}} \cup \mathcal{V}^1)}$  are left open. We therefore have for any  $x_{\mathcal{V}^2 \cap (\tilde{\mathcal{V}} \cup \mathcal{V}^1)}$

$$\langle \tau^{\mathcal{G}_2} \rangle [X_{\mathcal{V}^2 \cap (\tilde{\mathcal{V}} \cup \mathcal{V}^1)} = x_{\mathcal{V}^2 \cap (\tilde{\mathcal{V}} \cup \mathcal{V}^1)}] = \left( \sum_{x_{\mathcal{V}^2/(\tilde{\mathcal{V}} \cup \mathcal{V}^1)}} \prod_{e \in \mathcal{E}_2} \tau^e [X_e = x_e] \right) .$$

It follows with the above, that

$$\begin{aligned} \langle \tau^{\mathcal{G}} \rangle [X_{\tilde{\mathcal{V}}} = x_{\tilde{\mathcal{V}}}] &= \sum_{x_{\mathcal{V}^1/\tilde{\mathcal{V}}}} \left( \prod_{e \in \mathcal{E}_1} \tau^e [X_e = x_e] \right) \cdot \langle \tau^{\mathcal{G}_2} \rangle [X_{\mathcal{V}^2 \cap (\tilde{\mathcal{V}} \cup \mathcal{V}^1)} = x_{\mathcal{V}^2 \cap (\tilde{\mathcal{V}} \cup \mathcal{V}^1)}] \\ &= \left\langle \tau^{\mathcal{G}_1} \cup \{ \langle \tau^{\mathcal{G}_2} \rangle [X_{\mathcal{V}^2 \cap (\mathcal{V}^1 \cup \tilde{\mathcal{V}})}] \} \right\rangle [X_{\tilde{\mathcal{V}}} = x_{\tilde{\mathcal{V}}}] . \end{aligned}$$

□

We can interpret the inner contraction  $\langle \tau^{\mathcal{G}_2} \rangle [X_{\mathcal{V}^2 \cap (\mathcal{V}^1 \cup \tilde{\mathcal{V}})}]$  as a message, sent from  $\mathcal{G}_2$  to  $\mathcal{G}_1$ . Based on this intuition, we will define message passing schemes in the next section.

## 20.2 Exact Contractions

We apply Theorem 131 to split a contraction into subcontractions, which are consecutively performed.

Contractions can be performed partially, and the result passed to the rest of the network as a message.

### 20.2.1 Construction of Cluster Graphs

Let us first introduce with the cluster graph a mechanism to coarse grain the hypergraph capturing a tensor network.

**Definition 95** (Cluster Graph). *Given a tensor network  $\tau^{\mathcal{G}}$  a cluster partition is a partition of the tensor network into  $n$  clusters, by a function*

$$\alpha : \mathcal{E} \rightarrow [n] .$$

*The clusters are with tensors decorated edge sets  $C_i = \{e : \alpha(e) = i\}$  with variables  $\mathcal{V}_i = \bigcup_{e \in C_i} e$ .*

*We say, that the cluster graph satisfies the running intersection property, when for any clusters  $C_i$  and  $C_j$  and any  $v \in \mathcal{V}_i \cup \mathcal{V}_j$  there is a path between  $C_i$  and  $C_j$  with  $v \in \mathcal{V}_k$  for any cluster  $C_k$  along the path.*

Given a cluster graph to a tensor network, we can execute any global contraction by a contraction of local contraction to each cluster.

**Theorem 132.** *Given a tensor network  $\tau^G$  and a cluster graph. We then define for each cluster the node set*

$$\tilde{\mathcal{V}}_i = \bigcup_{j \neq i} \mathcal{V}_j$$

and have

$$\langle \tau^G \rangle [X_{\tilde{\mathcal{V}}}] = \left\langle \{ \langle \tau^{C_i} \rangle [X_{\mathcal{V}_i \cap (\tilde{\mathcal{V}}_i \cup \tilde{\mathcal{V}})}] : i \in [n] \} \right\rangle [X_{\tilde{\mathcal{V}}}] .$$

*Proof.* By Theorem 131 applied for each cluster seen as a subgraph.  $\square$

### 20.2.2 Message Passing to calculate contractions

Having a hypergraph  $\mathcal{G}$ , we iteratively apply Theorem 131 and call the  $\mathcal{G}_2$  a cluster. When iterating until  $\mathcal{G}$  is empty, we get a cluster graph, where all tensors are assigned to a cluster.

When the cluster are a polytree, that is a union of disjoint trees, we define messages between neighbored clusters  $C_i$  and  $C_j$  with  $C_j \prec C_i$  by the contractions

$$\delta_{j \rightarrow i} [X_{\mathcal{V}^i \cap \mathcal{V}^j}] = \left\langle \{ \delta_{\tilde{j} \rightarrow j} [X_{\mathcal{V}^{\tilde{j}} \cap \mathcal{V}^j}] : C_{\tilde{j}} \prec C_j \} \cup \tau^{C_j} \right\rangle [X_{\mathcal{V}^i \cap \mathcal{V}^j}] .$$

We note, that the messages are well defined by these recursive equations, exactly when the cluster graph is a polytree.

When the cluster graph is a tree, we can choose a root cluster and order the clusters by the topological order  $\prec$ .

**Lemma 43.** *When the cluster graph is a tree satisfying the running intersection property, we have for neighbored clusters  $C_i$  and  $C_j$  with  $C_j \prec C_i$*

$$\delta_{i \rightarrow \tilde{i}} [X_{\mathcal{V}^i \cap \mathcal{V}^{\tilde{i}}}] = \left\langle \{ \tau^{C_j} : C_j \prec C_i \} \right\rangle [X_{\mathcal{V}^i \cap \mathcal{V}^{\tilde{i}}}] .$$

*Proof.* By induction over the cardinality  $n$  of the preceding clusters.

$n = 1$  : For a single preceding cluster the statement holds trivial, since the preceding cluster is the cluster itself.

$n + 1 \rightarrow n$  : Let us now assume, that the statement holds for up to  $n$  preceding clusters, and let there be  $n + 1$  preceding clusters. We build another cluster graph for the cores different from  $C_i$ , by assigning each cluster  $C_{\tilde{j}}$  to the neighbor  $C_j$  where  $j \in N(i)$ , for which

$$C_{\tilde{j}} \prec C_j .$$

We use The. 132 on this constructed cluster graph and get

$$\left\langle \{ \tau^{C_{\tilde{j}}} [X_{\mathcal{V}^{\tilde{j}}}] : \tilde{j} \neq i \} \right\rangle [X_{\mathcal{V}^i}] = \left\langle \left\{ \langle \tau^{C_{\tilde{j}}} [X_{\mathcal{V}^{\tilde{j}}}] : \tilde{j} \prec j \rangle [X_{\tilde{\mathcal{V}}^j}] : j \in N(i) \right\} \right\rangle [X_{\mathcal{V}^i}]$$

Here by  $\tilde{\mathcal{V}}^j$  we denote the intersection of

$$\tilde{\mathcal{V}}^j = \left( \bigcup_{\tilde{j} \prec j} \mathcal{V}^{\tilde{j}} \right) \cap \left( \bigcup_{\tilde{j} \not\prec j} \mathcal{V}^{\tilde{j}} \right)$$

By the running intersection property, we have  $\mathcal{V}^j \cap \mathcal{V}^i = \tilde{\mathcal{V}}^j$ .

We further have for any  $j \in N(i)$  that

$$|\{ \tilde{j} \prec j \}| \leq n .$$

We can therefore apply the assumption of the induction and get

$$\left\langle \{ \tau^{C_{\tilde{j}}} [X_{\mathcal{V}^{\tilde{j}}}] : \tilde{j} \prec j \} \right\rangle [X_{\tilde{\mathcal{V}}^j}] = \delta_{j \rightarrow i} [X_{\mathcal{V}^j \cap \mathcal{V}^i}]$$

With the above, we arrive at

$$\delta_{i \rightarrow \tilde{i}} [X_{\mathcal{V}^i \cap \mathcal{V}^{\tilde{i}}}] = \left\langle \{ \tau^{C_j} : C_j \prec C_i \} \right\rangle [X_{\mathcal{V}^i \cap \mathcal{V}^{\tilde{i}}}] .$$

$\square$

**Theorem 133.** *When the cluster graph is a tree satisfying the running intersection property, then we have for each cluster  $C_i$  with neighbors  $N(i)$*

$$\langle \tau^{\mathcal{G}} \rangle [X_{V_i}] = \langle \{ \delta_{j \rightarrow i} [X_{V_i \cap V_j}] : j \in N(i) \} \cup \{ \tau^{C_i} \} \rangle [X_{V_i}]. \quad (55)$$

*Proof.* We use the topological order  $\prec$  of the clusters by the tree, when choosing a root by cluster  $C_i$ .

The claim then follows from The. 132 and Lem. 43. □

While we have defined message passing along the topological order of a graph, we can also define messages against the topological order, that is

$$\delta_{j \leftarrow i} = \langle \{ \delta_{i \leftarrow j} : C_i \prec C_j \} \cup \tau^{C_i} \rangle [X_{V_i \cap V_j}]$$

To this end, we can get a similar statement for nodes, which are not the roots of the cluster tree. The constructions at each cluster can then be computed batchwise, based on message passed along a topological order and against.

These message passing schemes can be derived from Lagrangian parameters given a local consistency polytope Wainwright and Jordan (2008).

### 20.2.3 Variable Elimination Cluster Graphs

**Remark 28** (Construction of Cluster Graphs by Variable Elimination). *Following an elimination order of the colors, mark those tensors containing the colors, which have not been marked before, as the cluster. A clique tree can be constructed by these cluster; when iterating through the clusters and either connect them to previous disconnected clusters or leave the current cluster disconnected. Add the disconnected clusters with the current cluster in case there are overlaps of their open colors. If the disconnected cluster added has more open colors,*

### 20.2.4 Bethe Cluster Graphs

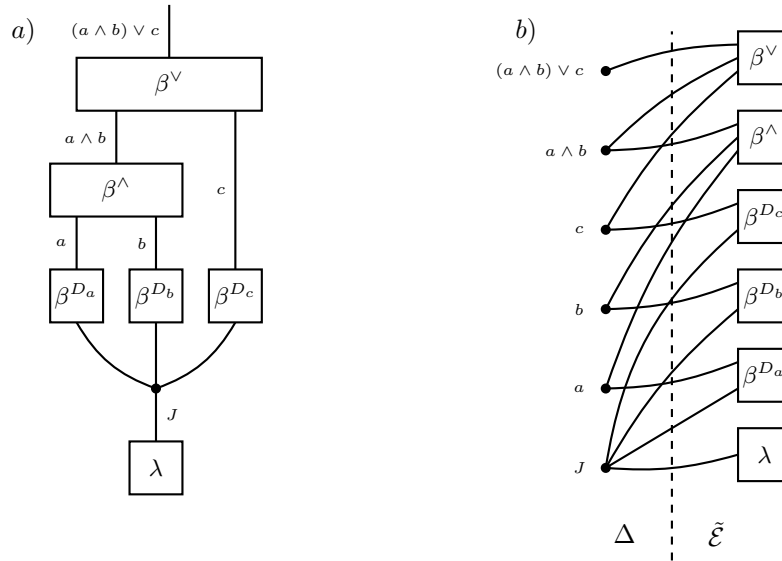


Figure 43: Example of a Bethe Cluster Graph. a) Example of a Tensor Network  $\tau^{\mathcal{G}}$ , which represents the by  $\lambda$  averaged evaluation of the formula  $(a \wedge b) \vee c$  on data  $D$ . b) Corresponding Bethe Cluster Hypergraph, which dual is bipartite by the sets  $\Delta$  and  $\tilde{\mathcal{E}}$ .

By adding delta tensors to each node  $v \in \mathcal{V}$  and defining its leg variables by  $v^e$  for  $e \in \mathcal{E}$ . We mark each such delta tensor by a cluster in  $\Delta^{\mathcal{G}}$ , as defined in the following (see also Figure 43).

**Definition 96.** *Given a tensor network  $\tau^{\mathcal{G}}$  on a decorated hypergraph  $\mathcal{G}$ , we define the Bethe Cluster Hypergraph  $\tilde{\mathcal{G}}$  as  $(\tilde{\mathcal{V}}, \tilde{\mathcal{E}} \cup \Delta^{\mathcal{G}})$  where we have*

- *Recolored Edges*  $\tilde{\mathcal{E}} = \{\tilde{e} : e \in \mathcal{E}\}$  where  $\tilde{e} = \{v^e : v \in e\}$ , which decoration tensor has same coordinates as  $\tau^e$
- *Nodes*  $\tilde{\mathcal{V}} = \bigcup_{e \in \mathcal{E}} \tilde{e}$
- *Delta Edges*  $\Delta^{\mathcal{G}} = \{\{v^e : e \ni v\} : v \in \mathcal{V}\}$ , each of which decorated by a delta tensor  $\delta^{\{v^e : e \ni v\}}$

By Lem. 27 this construction does not change contractions.

The dual is bipartite, since any variable appears exactly in one cluster in  $\tilde{\mathcal{E}}$  and in one cluster of  $\Delta^{\mathcal{G}}$ . This further makes the dual of the Bethe Cluster Hypergraph a proper graph (i.e. edges consistent of node pairs).

### 20.2.5 Computational Complexity

**Tree-width here:** By building a cluster tree to any partition, simply by including variables for the running intersection property. The maximum number of variables at a cluster is the tree-width and provides a complexity bound for the local contractions.

Naive execution of  $\langle \tau^{\mathcal{G}} \rangle [\tilde{\mathcal{V}}]$ :  $\prod_{v \in \mathcal{V}} m_v$  many products are built and summed up. When splitting contractions into local subcontractions, the product can be turned into sums with tremendous decrease in complexity.

### 20.3 Boolean Message Passing

Instead of the exact calculation of a contraction, let us now investigate schemes to sparsify the tensors before a contraction. To this end, we first show underlying properties of contractions enabling these schemes.

#### 20.3.1 Monotonicity of tensor contraction

To state the next theorem we use the nonzero function  $\mathbb{I}_{\neq 0} : \mathbb{R} \rightarrow [2]$  by  $\mathbb{I}_{\neq 0}(x) = 1$  if  $x \neq 0$  and  $\mathbb{I}_{\neq 0}(x) = 0$  else. Applied coordinatewise on tensors it marks the nonzero coordinates by 1.

We show that adding boolean tensor cores to an contraction orders the results by the partial ordering introduced in Def. 81.

**Theorem 134** (Monotonicity of Tensor Contractions). *Let  $\tau^{\mathcal{G}}, \tau^{\tilde{\mathcal{G}}}$  be tensor network of non-negative tensors and  $X_{\tilde{\mathcal{V}}}$  an arbitrary set of random variables. Then we have*

$$\mathbb{I}_{\neq 0} \left( \langle \tau^{\mathcal{G}} \cup \tau^{\tilde{\mathcal{G}}} \rangle [X_{\tilde{\mathcal{V}}}] \right) \prec \mathbb{I}_{\neq 0} \left( \langle \tau^{\mathcal{G}} \rangle [X_{\tilde{\mathcal{V}}}] \right).$$

*Proof.* It suffices to show that for any  $x_{\tilde{\mathcal{V}}}$  with

$$\mathbb{I}_{\neq 0} \left( \langle \tau^{\mathcal{G}} \cup \tau^{\tilde{\mathcal{G}}} \rangle [X_{\tilde{\mathcal{V}}} = x_{\tilde{\mathcal{V}}}] \right) = 1$$

we also have

$$\mathbb{I}_{\neq 0} \left( \langle \tau^{\mathcal{G}} \rangle [X_{\tilde{\mathcal{V}}} = x_{\tilde{\mathcal{V}}}] \right) = 1.$$

For any  $x_{\tilde{\mathcal{V}}}$  satisfying the first equation we find an extension  $x_{\mathcal{V}}$  to all variables of the tensor networks such that

$$\langle \tau^{\mathcal{G}} \cup \tau^{\tilde{\mathcal{G}}} \rangle [X_{\mathcal{V}} = x_{\mathcal{V}}] > 0$$

and it follows that

$$\langle \tau^{\mathcal{G}} \rangle [X_{\mathcal{V}} = x_{\mathcal{V}}] > 0 \quad \text{and} \quad \langle \tau^{\tilde{\mathcal{G}}} \rangle [X_{\mathcal{V}} = x_{\mathcal{V}}] > 0.$$

But this already implies, that

$$\mathbb{I}_{\neq 0} \left( \langle \tau^{\mathcal{G}} \rangle [X_{\tilde{\mathcal{V}}} = x_{\tilde{\mathcal{V}}}] \right) = 1.$$

□

### 20.3.2 Invariance of adding subcontractions

Let us now state an equivalence of the contraction, when we add the result of the same contraction. This property was used in the proof of The. 54.

**Theorem 135** (Invariance under adding subcontractions). *Let  $\tau^{\mathcal{G}}$  be a tensor network of non-negative tensors with variables  $X_{\mathcal{V}}$  and let  $\tau^{\tilde{\mathcal{G}}}$  be a subset. Then we have for any subset  $X_{\tilde{\mathcal{V}}}$  of  $X_{\mathcal{V}}$*

$$\left\langle \tau^{\mathcal{G}} \cup \{\mathbb{I}_{\neq 0} \left( \left\langle \tau^{\tilde{\mathcal{G}}} \right\rangle [X_{\tilde{\mathcal{V}}}] \right) \} \right\rangle [X_{\mathcal{V}}] = \langle \tau^{\mathcal{G}} \rangle [X_{\mathcal{V}}] .$$

*Proof.* For any  $x_{\mathcal{V}}$  with

$$\langle \tau^{\mathcal{G}} \rangle [X_{\mathcal{V}} = x_{\mathcal{V}}] = 0$$

we also have

$$\left\langle \tau^{\mathcal{G}} \cup \{\mathbb{I}_{\neq 0} \left( \left\langle \tau^{\tilde{\mathcal{G}}} \right\rangle [X_{\tilde{\mathcal{V}}}] \right) \} \right\rangle [X_{\mathcal{V}} = x_{\mathcal{V}}] = 0 .$$

For any  $x_{\mathcal{V}}$  with

$$\langle \tau^{\mathcal{G}} \rangle [X_{\mathcal{V}} = x_{\mathcal{V}}] \neq 0$$

we have for the reduction  $x_{\tilde{\mathcal{V}}}$  of the index  $x_{\mathcal{V}}$  that

$$\left\langle \tau^{\tilde{\mathcal{G}}} \right\rangle [X_{\tilde{\mathcal{V}}} = x_{\tilde{\mathcal{V}}}] \neq 0$$

and thus

$$\begin{aligned} \left\langle \tau^{\mathcal{G}} \cup \{\mathbb{I}_{\neq 0} \left( \left\langle \tau^{\tilde{\mathcal{G}}} \right\rangle [X_{\tilde{\mathcal{V}}}] \right) \} \right\rangle [X_{\mathcal{V}} = x_{\mathcal{V}}] &= \langle \tau^{\mathcal{G}} \rangle [X_{\mathcal{V}} = x_{\mathcal{V}}] \cdot \mathbb{I}_{\neq 0} \left( \left\langle \tau^{\tilde{\mathcal{G}}} \right\rangle [X_{\tilde{\mathcal{V}}}] \right) [X_{\tilde{\mathcal{V}}} = x_{\tilde{\mathcal{V}}}] \\ &= \langle \tau^{\mathcal{G}} \rangle [X_{\mathcal{V}} = x_{\mathcal{V}}] . \end{aligned}$$

□

**Remark 29.** *Similar statements hold, when dropping the non-negativity assumption on the, but demanding that all variables are left open.*

### 20.3.3 Basis Calculus as message passing scheme

Message Passing of directed and boolean message by basis encoding of functions can be interpreted as function evaluation. Each subfunction evaluation is passed in its one-hot encoding.

This is because any basis encoding of a function, the decomposition

$$\beta^q = \sum_{y \in \text{im}(q)} \left( \sum_{i: q(i)=y} \epsilon_i \right) \otimes \epsilon_y$$

is a SVD of the matricification of  $\beta^q$  with respect to incoming and outgoing legs.

Passing a message  $\epsilon_i$  in direction thus gives the message  $\epsilon_{q(i)}$ .

Note, that this is exact, whenever the graph is directed and acyclic. We do not need acyclicity of the underlying undirected graph.

**Remark 30** (Basis Calculus as Message Passing). *Given a tensor network of directed and binary tensor cores, each representing a function  $q_e$  depending on variables  $e^{\text{in}}$ . When there are not directed cycles, we define the compositions of  $q_e$  to be the function  $q$  from the nodes  $\mathcal{V}^1$  not appearing as incoming nodes to the nodes  $\mathcal{V}^2$  not appearing as outgoing nodes in an edge. Choosing arbitrary  $x_v \in [m_v]$  for  $v \in \mathcal{V}^1$  we have*

$$\left\langle \{ \beta^{q_e} [X_{e^{\text{out}}}, X_{e^{\text{in}}}] : e = (e^{\text{out}}, e^{\text{in}}) \in \mathcal{E} \} \right\rangle [\mathcal{V}^2] = \epsilon_{q(x_v : v \in \mathcal{V}^1)} .$$

### 20.3.4 Application

This properties can be applied as a sparsification of tensors before the execution of a contraction. The Knowledge Propagation Algorithm 8 produces in the knowledge cores conditions on non-vanishing coordinates. Thus, the knowledge cores can be locally contracted with the tensors, as a sparsification before performing a global contraction.



## 20.4 Discussion

Computing contractions by message passing is known to the graphical model community as belief propagation. There, the objective is the calculation of marginal probabilities of Markov Networks, which involve contractions of the corresponding factor tensors.

**Remark 31** (Approximate Message Passing Schemes). *When the cluster graphs are not trees, we cannot find a topological order of the clusters any more. Messages can still be defined implicitly by received neighbored messages, but the equivalence with global contractions cannot be established in general.*

Such algorithms are known in the graphical model community as loopy belief propagation.

When queries share same parts, can perform their contraction using dynamic programming. For conditional probability queries, which variables are the clusters of a cluster tree, this results in belief propagation.

## A Implementation in the `tnreason` package

We here document the implementation of the discussed concepts in the python package `tnreason`, in the version 2.0.0 `tnreason` is an abbreviation of **t**ensor **n**etwork **r**easoning, by which we emphasize the capabilities of this package to represent and answer reasoning tasks by tensor network contractions.

The package can be installed either by cloning the repository

<https://github.com/EnexaProject/enexa-tensor-reasoning>

or by

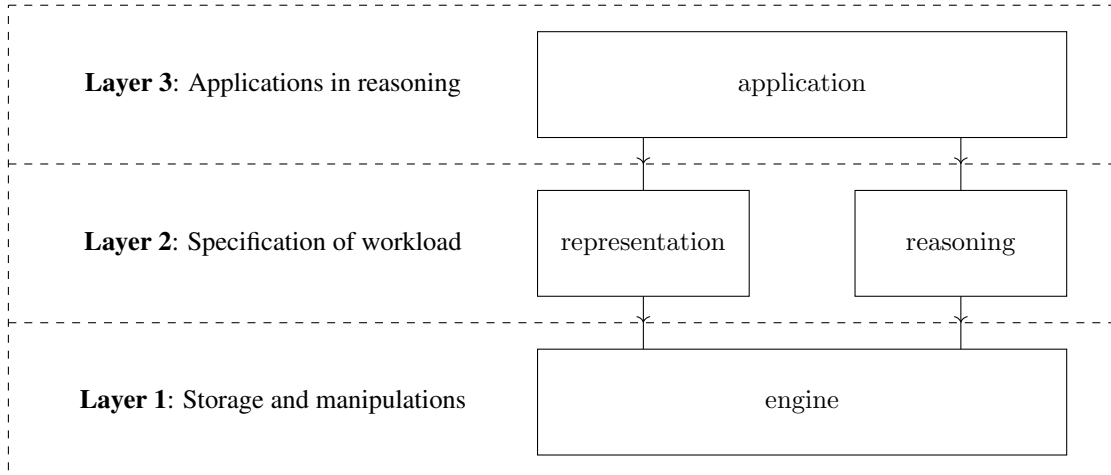
```
!pip install tnreason==2.0.0
```

### A.1 Architecture

`tnreason` is structured in four subpackages and three layers

- Layer 1: Storage and numerical manipulations, by subpackage `engine`, "Tensor Networks" -> building "tn" of `tnreason`
- Layer 2: Specification of workload, subpackage representation specific for storage, subpackage reasoning specific for manipulations
- Layer 3: Applications in reasoning, by subpackage `application`, "Reasoning" -> building "reason" of `tnreason`

We sketch this structure by



## A.2 Implementation of basic notation

First of all, we explain how the basis notation explained in Chapter 3 is reflected in the implementation.

### A.2.1 Categorical Variables and Representations

Categorical Variables are identified by strings, which then appear as colors of the corresponding tensor axes. Their dimension is stored in `shapeDicts`, but most practically these shapes are stored in the tensors in which variables appear. Suffixes in the color string (defined in `representation . suffixes`) denote the type of the variable:

- Distributed variables with color suffix "`_dV`":  $X$ .
- Computed variables with color suffix "`_cV`":  $Y$ .
- Selection variables with color suffix "`_sV`":  $L$ .
- Term variables with color suffix "`_tV`":  $O$ .

### A.2.2 Tensors

**Tensors** are objects of classes inheriting `engine.TensorCore` with main attributes

- `values`: Storing the coordinates of the tensors (individual realization for different cores)
- `colors`: List of the variables  $[Y_f, X_0, X_1]$
- `name`: Reflecting the notation such as  $\beta^f$
- `shape`: Storing the dimension of each appearing variable, as a list of integers with the same length as `colors`.

Suffixes in the name string (defined in `representation . suffixes`) highlight the origin and purpose of the tensor. Cores are named with suffixes based on their functionality

- Computation core with name suffix "`_cC`": They represent the computation of a function in basis calculus, and are directed cores. Their colors are  $[\text{headColors}] + [\text{inputColors}]$ , where  $[\text{inputColors}]$  are either distributed variables or, if having a composition of formulas. When the function is a selection augmentation of other functions, selection colors are listed in the end of  $[\text{inputColors}]$ .
- Activation core with name suffix "`_aC`": two-dimensional vectors representing of the activation core to a formula

Both the cores and the colors are further refined by infixes before the suffices to denote specific instantiations.

- "`_s`": Involving a selection variable
- "`_e`": Storing evidence about a variable
- "`_h`": Head of a function, typically the variable computed at a activation selector
- "`_f`": Function selection variables
- "`_p`"+"`_i`": Variable selection for argument at position  $i$
- "`_d`": Involving data (data cores and colors)

Further infixes are strings denoting atom names and neuron names.

Exploiting efficient representation tricks we further have the tensor name suffices:

- "`_atoC`": Atomization core, for sparse representation of categorical constraints
- "`_yselC`": Variable selection core: For sparse representation of variable selectors

**Initialization** Tensors are instantiated by

```
engine.getCore(coreType)(values, colors, name, shape)
```

where `coreType` is a string further specifying a specific implementation of tensors (see for more detail Sect. A.3). The default tensor implementation `NumpyCore` is chosen, when `coreType` is not specified.

One-hot encodings are specific tensors created in `representation`.

### A.2.3 Contractions

**Tensor networks**  $\tau^{\mathcal{G}} = \{\tau^e [X_e] : e \in \mathcal{E}\}$  are stored as dictionaries of tensors, where the keys coincide with the names of the corresponding tensors. The edges of the hypergraph  $\mathcal{G}$  used in the definition of tensor networks in Def. 8 are here labeled by the names of the tensors and the affected variables by the list of colors, being an attribute of each tensor.

**Contractions** are implemented in the subpackage `engine`, orienting on Def. 10. Reflected in the notation

$$\langle \tau^{\mathcal{G}} \rangle [X_{\tilde{\mathcal{V}}}]$$

a contraction is defined by

- Tensor Network  $\tau^{\mathcal{G}}$ , specified by a dictionary of tensor names as keys and valued by tensor cores.
- Open Variables  $\tilde{\mathcal{V}}$ , specified by a list of colors to the variables.

Contraction calls are implemented as

```
engine.contract(contractionMethod, coreDict, openColors, dimensionDict, evidenceColorDict)
```

where the arguments are

- `contractionMethod`: str, chooses one of the contraction providers. The default contraction method `NumpyEinsumis` chosen, when
- `coreDict`: Dictionary of TensorCores (of the above formats), representing the Tensor Network  $\tau^{\mathcal{G}}$
- `openColors`: List of str, each str identifying a color, that is a variable to be left open in the contraction
- `dimensionDict`: Dict valued by int and keys by str, storing dimensions to each variable. This is of optional usage, when a color in `openColors` does not appear in the `coreDict`.
- `evidenceColorDict`: Dict valued by int and keys by str, indicating sliced variables

Coordinates of tensors can be retrieved by

$$\langle \tau [X_{\mathcal{V}}] \rangle [X_{\tilde{\mathcal{V}}} = x_{\tilde{\mathcal{V}}}] .$$

We implement this by leaving `openColors` empty and passing  $x_{\tilde{\mathcal{V}}}$  as the `evidenceColorDict`, as a dictionary with keys by the **str** colors to the variables and values by the corresponding **int** indices.

Graphical illustrations can be generated by

```
engine.draw_factor_graph(coreDict)
```

where `coreDict` is a tensor network to be visualized.

### A.2.4 Function encoding schemes

Encoding schemes are implemented in the subpackage `representation`.

## A.3 Subpackage engine

The `engine` subpackage is for the storage and numerical manipulation of tensors and tensor networks. We organize the subpackage as the lowest layer of `tnreason`, specializing in storage of Tensor Networks and performing the contractions.

### A.3.1 Basis+ CP Decompositions storing values

**Specification of basis+ elementary tensors** We orient on basis+ sparse tensor decomposition in the initialization of tensor cores, as discussed in detail in Chapter 18. The basis+ elementary tensors have basis+ CP rank of 1 and admit a decomposition as (see Def. 91)

$$\tau [X_{\mathcal{V}}] = \lambda \cdot \langle \epsilon_{x_A} [X_A] \rangle [X_{\mathcal{V}}] .$$

Elementary basis+ tensors are in `tnreason` stored by a tuple

(value , posDict)

where posDict specifies the values to the variables, which do not have a trivial leg vector, and value a scalar scaling the basis vector. Comparing with the notation of Chapter 18, the keys of posDict correspond with  $A$ , the values of posDict with  $x_A$  and value corresponds with  $\lambda$ .

**Elementary Iterators** The initialization, coordinate retrieval and conversion operations of all tensor cores are oriented on basis+ CP Decompositions (46) of tensors. A tensor

$$\tau [X_V] = \sum_{(\lambda, A, x_A) \in \mathcal{M}} \lambda \cdot \langle \epsilon_{x_A} [X_A] \rangle [X_V] .$$

corresponds with an iterator over tuples (value , posDict), each specifying a basis+ elementary tensor in the sum.

**Core Arithmetics** When subscribing an instance exampleCore of engine.TensorCore by

exampleCore[posDict] = value

a basis+ elementary tensor specified by (value , posDict) is added to its values, that is

$$\tau [X_{[d]}] \leftarrow \tau [X_{[d]}] + \langle \epsilon_{x_A} [X_A] \rangle [X_{[d]}] .$$

The linear structure of tensors spaces are more further reflected in sums of engine.TensorCore instances, which are implemented with the same coreType, as

summed = exampleCore1 + exampleCore2

and scalar multiplication, where a scalar value of type **int** or **float**

multiplied = value \* exampleCore

Both operations are performed as manipulations of the tensors values. Contraction of two engine.TensorCore instances are performed by

contracted = exampleCore1.contract\_with(exampleCore2)

and are used in corewise contraction, where contractionMethod="CorewiseContractor".

**Initialization** Any instance of engine.TensorCore is initialized as a vanishing tensor  $\mathbb{I}[X_V]$ , when values is not specified. The values are then assigned by iteration over a sliceIterator over (value , posDict) tuples specifying elementary basis+ tensors, where the CP rank is the length of the iterator This initialization is by applied in the method

engine.create\_from\_slice\_iterator(shape , colors , sliceIterator , coreType , name)

where shape, colors , coreType, name are used in the call of an empty core by engine.get\_core and sliceIterator used to iterative add the basis+ elementary tensors to create the tensor.

**Storage of basis+ CP decompositions** The implemented tensor classes derived from engine.TensorCore differ in their implementation of values. Motivated from basis+ CP decompositions, most classes rely on a data base storing the (value , posDict) tuples. An overview over the derived classes is provided in Figure 44. Here PandasCore, TentrCore and PolynomialCore support sparse basis+ CP decompositions, by utilizing pandas.DataFrame, tentr . hypertrie and **list** as storage data base. These are implementations of the matrix representation of Remark 25. The NumpyCore class on the other hand is based relies on arrays as numpy.array as storage solution, which corresponds with the demand that each posDict contains all colors of the tensor. Effectively, this amounts to restricting to basis CP decomposition, which demanded ranks are always larger than basis+ CP decompositions (see The. 122).

### A.3.2 Contractions

The supported contraction methods are listed in Figure 45.

**Einstein Summation** is a syntax of specifying the contractions of arrays. Different possibilities are available to optimize over possible contraction paths, for example in numpy by the numpy.einsum\_path.

coreType	Used Package	Storage of values	Sparse basis+ CP support
NumpyCore	numpy	numpy.array	No
PandasCore	pandas	pandas.DataFrame	Yes
TentrisCore	tentris Bigerl et al. (2020)	tentris . hypertrie	Yes
PolynomialCore	--	list of (value , posDict)	Yes

Figure 44: Derived classes from engine .TensorCore, differing in the implemented storage of values .

contractionMethod (str)	Package	Applied procedure
"NumpyEinsum"	numpy	Einstein summation numpy.einsum
"TentrisEinsum"	tentris Bigerl et al. (2020)	Einstein summation tentris .einsum
"PgmpyVariableEliminator"	pgmpy	Variable Elimination of pgmpy.DiscreteFactor
"CorewiseContractor"	--	Contraction using core . contract_with ()

Figure 45: Implemented contraction methods in tnreason .

**Variable Elimination** contracts along a junction tree, build by variable elimination. We here use an implementation in the pgmpy package.

**Corewise Contraction** uses the `contract_with` method of tensor cores to contract in a given order.

#### A.4 Subpackage representation

The representation subpackage consists in a collection of core creation methods. We arrange the representation subpackage into the second layer of the tnreason architecture, since it specifies tensor cores which formats are specified in engine .

**Coordinate Calculus** Coordinatewise transformations (see Def. 78) are supported by

```
engine.coordinatewise_transform(coresList , transformFunction)
```

where `coresList` is a list of  $p$  tensors with identical variables and `transformFunction` a function  $h : \mathbb{R}^p \rightarrow \mathbb{R}$ .

**Basis Calculus** Basis encodings (see Def. 84) of functions  $q : \times_{k \in [d]} [m_k] \rightarrow \times_{l \in [r]} [m_l]$  are created by

```
engine.create_relational_encoding_from_lambda(inshape , outshape , incolors , outcolor)
```

where `indicesToIndices` is a lambda-function representing  $q$  and `inshape`, `outshape`, `incolors`, `outcolors` specify the input and output variables. Let us notice, that this procedure produces sums of  $\prod_{k \in [d]} m_k$  basis tensors corresponding with a basis CP decompositions. More involved initialization procedures based on basis+ elementary tensors calling `engine.create_from_iterator` might result in sparser representations.

**Propositional Connectives** are represented by strings. Figure 46 lists the supported logical connectives, which are implemented in `representation.basisplus_calculus`. If the **str** to the connective starts with "n", then the negated connective is encoded.

**Wolfram codes** provide a classification scheme of propositional formulae by natural numbers, which is supported in the script language. The Wolfram code has been designed for the classification of cellular automaton rules Wolfram (1983) and popularized in the book Wolfram (2002). Along this, the coordinate encodings of  $d$ -ary connectives  $\circ$  are flattened and interpreted as a binary number, which is transformed into a decimal number and represented as a string  $S(\circ)$ . To be more precise, to each  $d$ -ary connective its Wolfram code is calculated by

$$N(\circ) = \sum_{x \in [2^d]} 2^{2^d - x - 1} \cdot \circ(x).$$

$\circ$	$S(\circ)$	Notes	$\text{rank}^{\text{bas}+}(\circ)$	$\text{rank}^{\text{bas}+}(\beta^\circ)$
$\wedge$	"and"		1	3
$\vee$	"or"		2	3
$\Rightarrow$	"imp"	last variable as head, others premises	2	3
$\oplus$	"xor"	implemented as the negation of "eq", i.e. "neq"	3	5
$\Leftrightarrow$	"eq"		2	5
$X_k$	"pas" + "k"	kth atom	1	2
$\neg$	"not"	negation of the first argument, i.e. "npas0"	1	2

Figure 46: Supported connectives in the script language. The arity of all connectives is not restricted. We notice, that the basis+ rank  $\text{rank}^{\text{bas}+}(\beta^\circ)$  is independent of the arity and in most cases less than the naive bound of  $2^d$ .

featureType	Purpose	Canonical Parameter	Activation Core
"SingleSoftFeature"	$S_l$ of a statistic	Scalar ( <b>int</b> or <b>float</b> )	$\exp[\theta \cdot I_{\text{im}(S_l)}(Y_l)]$
"SoftPartitionFeature"	Partition statistics (see Def. 88) $\mathcal{S}$	$\theta[L]$	$\exp[\theta[L]]$
"HardPartitionFeature"	Partition statistics (see Def. 88)	Boolean tensor $\theta[L]$	$\theta[L]$

Figure 47: Features derived from representation .ComputedFeature, which are implemented in representation .features .

In the script language, the connective is represented by the string concatenation of the arity and the Wolfram code as

$$S(\circ) = "d" + " " + "N(\circ)".$$

#### A.4.1 Computation Activation Networks

**Features** are generation procedures of activation cores given canonical parameters. They are initialized by

```
ComputedFeature(featureColors, affectedComputationCores, shape, name)
```

where

- `featureColors` is a list of **str** variable colors, which are assigned to a created activation core
- `affectedComputationCores` is a list of **str** names of computation cores, which are required to compute the feature colors
- `shape` is a list of **int** dimensions to the variables

Mean parameters to features are computed by

```
exampleFeature.compute_meanParam(environmentMean)
```

where `environmentMean` is the contraction of a tensor network with open colors by the `featureColors`. They are further capable of computing local changes to canonical parameter in order to match mean parameters by

```
exampleFeature.local_update(environmentMean, meanParam)
```

To customize their purposes, individual feature classes are derived from `representation.ComputedFeature`, as listed in Figure 47.

**Computation Activation Networks** are the most general models representable in `tnreason`. They generalize distributions such as those computable by a statistic (see Def. 26) as well as constraint satisfaction problems (see Def. 48). Computation Activation Networks are initialized by

```
representation.ComputationActivationNetwork(featureDict, computationCoreDict, baseL
```

where

- `featureDict` is a dictionary of features, where the keys correspond with the names of the features
- `computationCoreDict` is a tensor network of computation cores, which needs to contain all names of computation cores required to compute all feature colors
- `baseMeasureCoreDict` is a tensor network representing of  $\nu [X_V]$
- `canParamDict` is a dictionary of canonical parameters (see Figure 47)

Computation Activation Networks can be instantiated as cores by `exampleCANetwork.create_cores()` or as an energy dictionary by `exampleCANetwork.create_cores()`. Energy dictionaries are stored as dictionaries with values by tuples of value, tensorNetwork, representing a by value weighted sum of the tensorNetwork.

**Example 37** (Representation of a member of an exponential family). To represent  $\mathbb{P}^{\mathcal{S}, \theta, \nu}$  we initialize a

- `featureDict` as a dictionary of representation. `SingleSoftFeature` to each feature  $\mathcal{S}_l$
- `computationCoreDict`  $\tau^G$  of computation cores such that

$$\tau^G [Y_{[p]}, X_V] = \beta^{\mathcal{S}} [Y_{[p]}, X_V] .$$

Decompositions and redundancies between coordinates  $\mathcal{S}_l$  can be exploited to find an efficient tensor network.

- `baseMeasureCoreDict`  $\tau^{\tilde{G}}$  to represent the base measure  $\nu$
- `canParamDict` of canonical parameters, valued by the coordinates  $\theta [L = l]$

## A.5 Subpackage reasoning

The reasoning subpackage implements contraction-based reasoning algorithm on representation.ComputationActivationNetworks. As the representation subpackage it is arranged in the second layer of the tnreason architecture, since it specifies the manipulation of tensor networks in the engine subpackage.

### A.5.1 Sampling

Sampling is performed by MCMC methods calling local sampling methods, which are derived classes from reasoning.SampleCoreBase. The energy-based algorithms execute reasoning tasks solely on energy dictionaries, which are created by reasoning.ComputationActivationNetworks. reasoning.EnergyBasedGibbs

### A.5.2 Variational Inference

An overview over the variational inference methods in presented in Figure 48.

**Forward mappings** are implemented by reasoning.ForwardContractor as contraction of all cores, and in reasoning.ExpectationPropagator as a message-passing approach.

**Backward mappings** (see Sect. 6.8) are implemented by reasoning.BackwardAlternator as alternating algorithms iteratively updating the canonical parameters to single features (see Algorithm 5).

**Mean field methods** are approximation methods of energy tensors by tractable exponential families (see Sect. 6.7). Given an energy dictionary (e.g. created by `exampleCANetwork.get_energy_dict()`) the naive mean field method is implemented as

```
reasoning.NaiveMeanField(energyDict)
```

and the more general markov network based mean field method by

```
reasoning.GenericMeanField(energyDict, edgeColorDict)
```

where `edgeColorDict` specifies the graph of the approximating markov network.

### A.5.3 Optimization

Optimization is a reasoning task of finding a maximal coordinate given a tensor network. The supported methods are implemented in reasoning.optimization\_handling and listed in Figure 49.

<b>inferenceMethod</b> (str)	<b>Applied procedure</b>	<b>Dependency</b>
"ForwardContractor"	Contraction of all cores keeping the feature colors open	--
"BackwardAlternator"	Iterative local updates to match the mean parameters	Forward inferer to iteratively update the mean parameters
"ExpectationPropagator"	Iterative updates of messages between feature clusters	Forward and backward inferer used for the computation of messages

Figure 48: Implemented inference methods in `tnreason` .

<b>optimizationMethod</b> (str)	<b>Package</b>	<b>Applied procedure</b>
"numpyArgMax"	numpy	Transformation into a numpy core and solution by <code>numpy.argmax</code>
"gurobi"	gurobipy	Transformation into an ILP and solution by <code>gurobipy.optimize</code>
"gibbsSample"	--	Simulated annealing based on gibbs sampling
"meanFieldSample"	--	Mean field approximation combined with "gibbsSample"

Figure 49: Implemented optimization methods in `tnreason` .

## A.6 Subpackage application

With the application subpackage we provide an interface for reasoning workload. It builds a third layer, since it used representation to represent knowledge by tensor networks and reasoning in the execution of reasoning tasks. A user-friendly high-level syntax of script language (logical formulas or neuro-symbolic architectures) for the specification of tensor networks creation, such as propositional formulas or categorical constraints, is introduced. Given a specification of a formula  $f$  in script language  $S(\cdot)$ , the task amounts to building a semantic representation based on the syntactic specification.

### A.6.1 Representation of formulas

Propositional formulas  $f$  are represented in three schemes:

- Syntactical representation by a script language  $S(f)$  as nested lists (see Sect. A.6.2).
- Syntactical representation by a **str** specifying a color to the categorical variables  $Y_f$ .
- Representation of formulas by tensor networks being contracted to  $\beta^f[Y_f, X_V]$

Conversions of the formats:

- $S(f)$  to color by

```
application.get_formula_color(S(f))
```

Here the nested lists are turned in a string by concatenating all elements of a list with "\_" and adding "[" and "]" at the beginning and end of each list.

- $S(f)$  to tensor network

```
application.create_raw_cores(S(f))
```

This creates the connective cores for the semantic representation of  $\beta^f$ . We encode them by iterative calls of `engine.create_from_iterators` .

### A.6.2 Script Language

To specify propositional sentences, neuro-symbolic architectures and Markov Logic Networks, we developed a script language.



$d$ -ary connective	"and"   "or"   "imp"   "xor"   "eq"   "not"   ... " $d$ " + "_" + " $N$ ", where $N < 2^{2^d-1}$
Atomic Formula	Set of strings not in Connectives
Complex Formula	Atomic Formula   [ $d$ -ary connective, $d$ Complex Formulas]

Figure 50: Backus-Naur form of the grammar producing the nested list expressions. The string connectives are either appearing in the list of Figure 46, or represented by a Wolfram code.

**Atomic Formulas** are represented by arbitrary strings, which are not used for the representation of connectives. We further avoid the symbols {"(", ")", "\_"} in the names of atoms, to not confuse them with colors of categorical variables.

**Composed Formulas** are represented by nested lists, where each sublist is either specifying an atomic formula (if string) or another composed formula. For example, a formula  $f_1 \circ f_2$  is represented by

$$S(f_1 \circ f_2) = [S(\circ), S(f_1), S(f_2)]$$

where we apply the conventions

- Connectives are at the 0th position in each list
- Further entries are either atoms as strings or encoded formulas itself

The nested lists follows a grammar, which is provided in Figure 50 in its Backus-Naur form.

**Example 38** (Encoding of the Wet Street example). *For example we have*

- *Atomic variable Rained by*  
 $S(\text{Rained}) = \text{"Rained"}$
- *Negative literal  $\neg$ Rained by*  
 $S(\neg \text{Rained}) = [\text{"not"}, \text{"Rained"}]$
- *Horn clause  $(\text{Rained} \Rightarrow \text{Wet})$  by*  
 $S(\text{Rained} \Rightarrow \text{Wet}) = [\text{"imp"}, \text{"Rained"}, \text{"Wet"}]$
- *Knowledge Base  $(\neg \text{Rained}) \wedge (\text{Rained} \Rightarrow \text{Wet})$  by*  
 $S(\neg \text{Rained}) \wedge (\text{Rained} \Rightarrow \text{Wet}) = [\text{"and"}, [\text{"not"}, \text{"Rained"}], [\text{"imp"}, \text{"Rained"}, \text{"Wet"}]]$

### Knowledge Bases

We distinguish here formulas, with propositional logic interpretation and formulas which have a soft logic interpretation. The formulas with hard interpretation are called facts in a knowledge base  $\mathcal{KB}$  and encoded by dictionaries

$$\{\text{key}(f) : S(f) \text{ for } f \in \mathcal{KB}\}$$

### Markov Logic Networks

The formulas with soft interpretation are called weighted formulas and encoded by  $\exp[\theta_f \cdot f]$ . We thus require a specification of the weights, which we do by adding  $\theta_f$  as a float or an int to the list  $S(f)$ . We then store Markov Logic Networks by dictionaries

$$\{\text{key}(f) : S(f) + [\theta_f] \text{ for } f \in \mathcal{F}\}$$

### Neuro-Symbolic Architecture by Nested Lists

To specify neuro-symbolic architectures in terms of formula selecting maps, as has been the subject of Chapter 10 we further exploit the nested list structure of encoding propositional logics. We replace, in each hierarchy of the nested structure each entry by a list of possible choices. In this way, we reinterpret the list index as the choice indices  $l$  introduced for connective and formula selections (see Def. 53 and 56). More formally, the production rules are formalized in Figure 51 by the extension of the Backus-Naur form in Figure 50.

$d$ -ary connectives	$[d\text{-ary connective}] \mid [d\text{-ary connective}] + d\text{-ary connectives}$
Dependency Choice	Atomic Formula $\mid$ Neuron
Dependency Choices	$[\text{Dependency Choice}] \mid [\text{Dependency Choice}] + \text{Dependency Choices}$
Neuron	$[d\text{-ary connectives}, d \text{ Dependency Choices}]$

Figure 51: Extension of the grammar in Backus-Naur form in Figure 50 to describe selections of functions.

**Connective selectors** (see Def. 53) are encoded by the list

$$S(\circ) = [S(\circ_0), \dots, S(\circ_{p-1})]$$

and a formula selector (see Def. 56) by

$$S(\mathcal{H}) = [S(\circ_0), \dots, S(\circ_{p-1})]$$

**A logical neuron** of order  $n$  (see Def. 57), defined by a connective selector  $\circ$ , and a formula selector  $\mathcal{H}_k$  on each argument  $k \in [n]$ , is encoded by

$$S(\mathcal{N}) = [S(\circ), S(\mathcal{H}_0), \dots, S(\mathcal{H}_{n-1})]$$

Only the unary  $n = 1$  and the  $n = 2$  cases are supported.

The resulting nested lists indices have an alternating interpretation at each level compared with the elements of each list. That is, when  $S(\mathcal{N})$  is the encoding of a neuron, then any element  $x \in S(\mathcal{N})$  represents a list of choices. When  $x$  is not the first element, then each choice is either the encoding  $S(X)$  of an atomic formula, or another neuron.

**A neural architecture**  $\mathcal{A}$  is then represented in the dictionary

$$S(\mathcal{A}) = \{\text{key}(\mathcal{N}) : S(\mathcal{N}) \text{ for } f \in \mathcal{A}\}$$

where  $\text{key}(\mathcal{N})$  is a string, which can be used in the formula selections of other neurons.

It is important that the directed graph of neurons induced by the choice possibilities is acyclic, to ensure a well-defined architecture.

**Example 39** (Neuro-Symbolic Architecture for the Wet Street). *Following the wet street example, we can define a neuron by*

$$S(\mathcal{N}) = [["imp", "eq"], ["Wet", "Sprinkler"], ["Street"]]$$

from which the formulas

$$\begin{aligned} &["imp", "Wet", "Street"] \\ &["eq", "Wet", "Street"] \\ &["imp", "Sprinkler", "Street"] \\ &["eq", "Sprinkler", "Street"] \end{aligned}$$

can be chosen. Combining this neuron with further neurons, e.g. by the architecture

$$S(\mathcal{A}) = \{ \text{"neur1": } [["imp", "eq"], ["neur2"], ["Street"]], \\ \text{"neur2": } [["lnot", "id"], ["Wet", "Sprinkler"], ["Street"]] \}$$

the expressivity increases. In this case, the further neuron provides the flexibility of the first atoms to be replaced by its negation.

### A.6.3 Distributions

Distributions are procedures to specify representation .ComputationActivationNetworks and are derived from the base class representation . DistributionBase . Each distribution needs to have a routine

```
exampleDistribution.create_caNetwork()
```

creating the corresponding network.

**Markov Networks** are distributions, which do not have computation cores. The positive coordinates of the factors are represented by `representation . SoftPartitionFactors` and the support of the factors `representation . HardPartitionFactors`.

**Empirical Distributions** are special instances of Markov Networks specifying distributions of sample data. We represent the values as a CP Format of data cores as specified in Sect. 6.3.1. They are initialized by

```
application . get_empirical_distribution (sampleDf , atomColumns , interpretation , dimension)
```

where

- `sampleDf` is a `pandas.DataFrame` specifying the data
- `atomColumns` is a list of column names in `sampleDf` to be extracted as variables of the distribution.
- `interpretation` is either *"atomic"* or *"categorical"*, specifying whether the entries in `sampleDf` are interpreted as uncertainties in the interval  $[0, 1]$ , or as assignments to

Here the partition function is the number of samples used in the creation of the empirical distribution.

**HybridKnowledgeBases** are probability distributions, which are specified by propositional formulas in the script language.

```
application . HybridKnowledgeBase
```

They are initialized with arguments

- **facts**: Dictionary of propositional formulas stored as  $S(f)$  representing hard logical constraints
- **weightedFormulas**: Dictionary of propositional formulas stored as  $S(f)+[\theta_f]$  representing soft logical constraints
- **evidence**: Dictionary of atomic formulas, where key are the formulas in string representation and values the certainty in  $[0, 1]$  (float or int) of the atom being true
- **categoricalConstraints**: Dictionary of categorical constrained, which values are lists of atomic formulas stored as strings  $S(X)$

#### A.6.4 Inference

To simplify deductive inference on models a class

```
application . InferenceProvider
```

taking a `representation . ComputationActivationNetwork` or `application . Distribution` has been implemented.

**Probabilistic queries** as specified Def. 35) by

```
. query (variableList , evidenceDict)
```

**Mode queries** by

```
. exact_map_query ()
```

**Entailment** from the distribution (Def. 47) is decided by

```
. ask (queryFormula , evidenceDict)
```

where `queryFormula` is the formula  $f$  to be tested for entailment in the representation  $S(f)$ .

**Samples** can be drawn by

```
. draw_samples (sampleNum , variableList , annealingPattern)
```

based on Gibbs sampling, where

- `sampleNum` (int) gives the number of samples to be drawn

- `variableList` (list of str) defines the variables to be represented by the samples (default: all atoms in the distribution)
- `annealingPattern` specifies an annealing pattern

### A.6.5 Learning

To learn instances of `application.HybridKnowledgeBase` on data the class

`application.HybridLearner`

is initialized with the arguments

- `knowledgeBase`: Distribution representing a current model to be improved
- `specDict`: A neuro-symbolic architecture encoded in a dictionary of neurons

**Formula Selecting Neural Networks** (Def. 57) are specified to define a proposal distribution. They are encoded by creating all formula selecting neurons, each involving a

- a connective selection map (Def. 53)
- variable selection cores (Def. 54) to each argument using the decomposition of The. 56.

Each selection variable of each neuron comes with a control variable with suffix `"_sV"`.

**Structure Learning** is performed by

`application.HybridLearner.propose_candidate()`

where a proposal distribution is instantiated and then sampled given the specified inference method.

**Weight Estimation** is performed by

`application.HybridLearner.infer_weights_on_data(empDistribution)`

where `empDistribution` is used to infer the mean parameters to be matched by the canonical parameters of `knowledgeBase.weightedForm`

## B Glossary

### B.1 Tensors

Small greek letters are reserved for the notation of tensors:

Notation	Name	Reference
$\alpha^{l, \theta[L=l]}[Y_l]$	Activation Core	The. 11
$\beta^q[Y_q, X_{[d]}]$	Basis Encoding of a function $q$	Def. 84
$\delta^{[d], m}[X_{[d]}]$	Diracs delta	Example 19
$\epsilon_{x_{[d]}}[X_{[d]}]$	One-hot Encoding	Def. 6
$\eta^{\mathcal{S}, \mathbb{P}^*, D}[L]$	Noise tensor	Def. 62
$\theta[L]$	Canonical Parameter	Def. 24
$\kappa^e[X_e]$	Knowledge core	Def. 50
$\lambda[I]$	Scalar core in CP decompositions	Def. 89
$\mu[L]$	Mean Parameter	Def. 32
$\nu[X_{[d]}]$	Boolean base measure	Sect. 5.1.2
$\rho^k[X_k, I]$	Leg core in CP decompositions	Def. 89
$\sigma^q[X_{[d]}, L]$	Selection Encoding of a vector-valued function $q$	Def. 15
$\tau^{\mathcal{G}}[X_{[d]}]$	Tensor network of tensors	Def. 8
$\tau[X_{[d]}]$	Energy tensor	Def. 24
$\phi[X_{[d]}]$	Coordinate Encoding of a function $q$	Def. 77, often abbreviated by $q$
$\chi^q[X_{[d]}]$		

Sets of tensors are represented by large greek letters:

Notation	Name	Reference
$\Gamma^{\mathcal{S}, \nu}$	Exponential family	Def. 24
$\Lambda^{\mathcal{S}, \mathcal{G}}$	Sets of by $\mathcal{S}$ and $\mathcal{G}$ computable distributions	Def. 26
$\mathcal{M}_{\mathcal{S}, \nu}$	Polytope of mean parameters	Def. 32

In the implementation in `tnreason`, we distinguish between computation and activation cores. The coarse roles are the computation of a function using basis calculus and the activation of the prepared variable to shape a probability distribution.

Name	Notation	String Suffix
Computation Core	$\beta^{\cdot}$	"_cC"
Activation Core	$\alpha^{\cdot}, \kappa^{\cdot}$	"_aC"

## B.2 Variables

Variables are denoted by large latin letters, their indices by the corresponding small letters. We distinguish between the coarse types of variables:

Name	Notation	String Suffix
Distributed Variable	$X.$	"_dV"
Computed Variable	$Y.$	"_cV"
Selection Variable	$L.$	"_sV"
Term Variable	$O.$	"_tV"

### B.3 Maps

We have in this work encountered different maps, which have been encoded as tensors. Note, that in order to ease the notation, when not specified otherwise the coordinate encoding  $\chi$  has been used.

Notation	Name	Domain	Range	Reference
$f$	propositional formula	$\times_{k \in [d]} [2]$	$\{0, 1\}$	Def. 43
$\mathcal{H}_C$	$d$ -ary connective selecting map	$\times_{k \in [d]} [2]$	$\times_{l \in [p_C]} [2]$	Def. 53
$\mathcal{H}_V$	Variable selecting map	$\times_{l \in [p_V]} [2]$	$\times_{l \in [p_V]} [2]$	Def. 54
$\mathcal{H}_S$	State selection map	$[m]$	$\times_{k \in [d]} [2]$	Def. 55
$\mathcal{KB}$	Knowledge Base (conjunction of formulas)	$\times_{k \in [d]} [2]$	$\{0, 1\}$	
$\mathbb{P} [X_{[d]}]$	Probability distribution	$\times_{k \in [d]} [m_k]$	$[0, 1]$	Def. 16
$q$	Function between states of factored representations	$\times_{k \in [d]} [m_k]$	$\times_{l \in [r]} [m_l]$	

### B.4 Contraction equations

We here provide a summary for the application of contractions and normalization in probabilistic and logical reasoning.

Concept	Contraction Equation	Reference
Marginal probability	$\mathbb{P} [X_0] = \langle \mathbb{P} \rangle [X_0]$	Def. 19
Conditional probability	$\mathbb{P} [X_0   X_1] = \langle \mathbb{P} \rangle [X_0   X_1]$	Def. 20
Markov Network Distribution	$\mathbb{P}^{\tau^G} = \langle \tau^G \rangle [\mathcal{V}   \emptyset]$	Def. 27
Partition Function	$\mathcal{Z} (\tau^G) = \langle \tau^G \rangle [\emptyset]$	Def. 27
Independence of $X_0$ and $X_1$	$\langle \mathbb{P} \rangle [X_0, X_1] = \langle \mathbb{P} \rangle [X_0] \otimes \langle \mathbb{P} \rangle [X_1]$	Def. 21, The. 6
Independence of $X_0$ and $X_1$ conditioned on $X_2$	$\langle \mathbb{P} \rangle [X_0, X_1   X_2] = \langle \mathbb{P} \rangle [X_0   X_2] \otimes \langle \mathbb{P} \rangle [X_1   X_2]$	Def. 22, The. 7

## C \*

### Bibliography

Martín Abadi, Ashish Agarwal, Paul Barham, Eugene Brevdo, Zhifeng Chen, Craig Citro, Greg S. Corrado, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Ian Goodfellow, Andrew Harp, Geoffrey Irving, Michael Isard, Yangqing Jia, Rafal Jozefowicz, Lukasz Kaiser, Manjunath Kudlur, Josh Levenberg, Dan Mane, Rajat Monga, Sherry Moore, Derek Murray, Chris Olah, Mike Schuster, Jonathon Shlens, Benoit Steiner, Ilya Sutskever, Kunal Talwar, Paul Tucker, Vincent Vanhoucke, Vijay Vasudevan, Fernanda Viegas, Oriol Vinyals, Pete Warden, Martin Wattenberg, Martin Wicke, Yuan Yu, and Xiaoqiang Zheng. TensorFlow: Large-Scale Machine Learning on Heterogeneous Distributed Systems, March 2016. URL <http://arxiv.org/abs/1603.04467>. arXiv:1603.04467 [cs].

Christian Agerbeck and Mikael Hansen. A Multi-Agent Approach to Solving NP-Complete Problems. 2008. URL <https://www.semanticscholar.org/paper/A-Multi-Agent-Approach-to-Solving-NP-Complete-Agerbeck-Hansen/3762bf7893da14839e06ae000b9e04d63dac8af4>.

Grigoris Antoniou, Paul Groth, Frank Van Harmelen, and Rinke Hoekstra. *A Semantic Web Primer, third edition*. The MIT Press, Cambridge (Mass.), third edition edition, August 2012. ISBN 978-0-262-01828-9.

Artur S. Avila Garcez and Gerson Zaverucha. The Connectionist Inductive Learning and Logic Programming System. *Applied Intelligence*, 11(1):59–77, July 1999. ISSN 1573-7497. doi: 10.1023/A:1008328630915. URL <https://doi.org/10.1023/A:1008328630915>.

- Samy Badreddine, Artur d'Avila Garcez, Luciano Serafini, and Michael Spranger. Logic Tensor Networks. *Artificial Intelligence*, 303:103649, February 2022. ISSN 0004-3702. doi: 10.1016/j.artint.2021.103649. URL <https://www.sciencedirect.com/science/article/pii/S0004370221002009>.
- Ivana Balazevic, Carl Allen, and Timothy Hospedales. Tucker: Tensor Factorization for Knowledge Graph Completion. *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 5184–5193, 2019. doi: 10.18653/v1/D19-1522. URL <https://www.aclweb.org/anthology/D19-1522>. Conference Name: Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP) Place: Hong Kong, China Publisher: Association for Computational Linguistics.
- Albert-László Barabási. *Network Science*. Cambridge University Press, Cambridge, illustrated edition edition, July 2016. ISBN 978-1-107-07626-6.
- Richard E. Bellman. *Adaptive Control Processes*. Princeton University Press, New Jersey, 1961. ISBN 978-1-4008-7466-8. Publication Title: Adaptive Control Processes.
- Gregory Beylkin and Martin J. Mohlenkamp. Algorithms for Numerical Analysis in High Dimensions. *SIAM Journal on Scientific Computing*, 26(6):2133–2159, January 2005. ISSN 1064-8275, 1095-7197. doi: 10.1137/040604959.
- Alexander Bigerl, Felix Conrads, Charlotte Behning, Mohamed Ahmed Sherif, Muhammad Saleem, and Axel-Cyrille Ngonga Ngomo. Tentriss – A Tensor-Based Triple Store. In Jeff Z. Pan, Valentina Tamma, Claudia d'Amato, Krzysztof Janowicz, Bo Fu, Axel Polleres, Oshani Seneviratne, and Lalana Kagal, editors, *The Semantic Web – ISWC 2020*, volume 12506, pages 56–73. Springer International Publishing, Cham, 2020. ISBN 978-3-030-62418-7 978-3-030-62419-4. doi: 10.1007/978-3-030-62419-4\_4. URL [https://link.springer.com/10.1007/978-3-030-62419-4\\_4](https://link.springer.com/10.1007/978-3-030-62419-4_4). Series Title: Lecture Notes in Computer Science.
- Peter G. Casazza, Gitta Kutyniok, and Friedrich Philipp. Introduction to Finite Frame Theory. In Peter G. Casazza and Gitta Kutyniok, editors, *Finite Frames: Theory and Applications*, pages 1–53. Birkhäuser, Boston, 2013. ISBN 978-0-8176-8373-3. doi: 10.1007/978-0-8176-8373-3\_1. URL [https://doi.org/10.1007/978-0-8176-8373-3\\_1](https://doi.org/10.1007/978-0-8176-8373-3_1).
- Andrzej Cichocki. Era of Big Data Processing: A New Approach via Tensor Networks and Tensor Decompositions. *arXiv:1403.2048 [cs]*, March 2014. arXiv: 1403.2048.
- Andrzej Cichocki, Danilo Mandic, Lieven De Lathauwer, Guoxu Zhou, Qibin Zhao, Cesar Caiafa, and HUY ANH PHAN. Tensor Decompositions for Signal Processing Applications: From two-way to multiway component analysis. *IEEE Signal Processing Magazine*, 32(2):145–163, March 2015. ISSN 1558-0792. doi: 10.1109/MSP.2013.2297439. Conference Name: IEEE Signal Processing Magazine.
- P. Clifford and J. M. Hammersley. Markov fields on finite graphs and lattices. *Unpublished*, 1971. URL <https://ora.ox.ac.uk/objects/uuid:4ea849da-1511-4578-bb88-6a8d02f457a6>. Publisher: University of Oxford.
- William Cohen, Fan Yang, and Kathryn Rivard Mazaitis. TensorLog: A Probabilistic Database Implemented Using Deep-Learning Infrastructure. *Journal of Artificial Intelligence Research*, 67:285–325, February 2020. ISSN 1076-9757. doi: 10.1613/jair.1.11944. URL <https://jair.org/index.php/jair/article/view/11944>.
- Vin de Silva and Lek-Heng Lim. Tensor Rank and the Ill-Posedness of the Best Low-Rank Approximation Problem. *SIAM Journal on Matrix Analysis and Applications*, 30(3):1084–1127, January 2008. ISSN 0895-4798, 1095-7162. doi: 10.1137/06066518X.
- Morris H. DeGroot. *Probability and Statistics*. PEARSON INDIA, January 2016. ISBN 978-93-325-7387-1.
- Caglar Demir and Axel-Cyrille Ngonga Ngomo. DRILL- Deep Reinforcement Learning for Refinement Operators in ALC. *CoRR*, abs/2106.15373, 2021. URL <https://ris.uni-paderborn.de/record/25217>.
- Mike Espig, Wolfgang Hackbusch, Thorsten Rohwedder, and Reinhold Schneider. Variational calculus with sums of elementary tensors of fixed rank. *Numerische Mathematik*, 122(3):469–488, November 2012. ISSN 0945-3245. doi: 10.1007/s00211-012-0464-x.
- Antonio Falco and Wolfgang Hackbusch. On Minimal Subspaces in Tensor Representations. *Foundations of Computational Mathematics*, 12:765–803, December 2012. doi: 10.1007/s10208-012-9136-6.
- Simon Foucart and Holger Rauhut. *A Mathematical Introduction to Compressive Sensing*. Applied and Numerical Harmonic Analysis. Birkhäuser Basel, 2013. ISBN 978-0-8176-4947-0. doi: 10.1007/978-0-8176-4948-7. URL <https://www.springer.com/de/book/9780817649470>.
- Python Software Foundation. Python Language Reference, version 3.13.2, April 2025. URL <https://docs.python.org/3/>.

- Luis Antonio Galárraga, Christina Teflioudi, Katja Hose, and Fabian Suchanek. AMIE: association rule mining under incomplete evidence in ontological knowledge bases. In *Proceedings of the 22nd international conference on World Wide Web*, pages 413–422, Rio de Janeiro Brazil, May 2013. ACM. ISBN 978-1-4503-2035-1. doi: 10.1145/2488388.2488425. URL <https://dl.acm.org/doi/10.1145/2488388.2488425>.
- Varun Ganapathi, David Vickrey, John Duchi, and Daphne Koller. Constrained approximate maximum entropy learning of Markov random fields. In *Proceedings of the Twenty-Fourth Conference on Uncertainty in Artificial Intelligence*, UAI’08, pages 196–203, Arlington, Virginia, USA, July 2008. AUAI Press. ISBN 978-0-9749039-4-1.
- Artur d’Avila Garcez, Marco Gori, Luis C. Lamb, Luciano Serafini, Michael Spranger, and Son N. Tran. Neural-Symbolic Computing: An Effective Methodology for Principled Integration of Machine Learning and Reasoning, May 2019. URL <http://arxiv.org/abs/1905.06088>. arXiv:1905.06088 [cs].
- Patrick Gelß, Stefan Klus, Jens Eisert, and Christof Schütte. Multidimensional Approximation of Nonlinear Dynamical Systems. *Journal of Computational and Nonlinear Dynamics*, 14(6):061006–061006–12, April 2019. ISSN 1555-1415. doi: 10.1115/1.4043148.
- Lise Getoor and Ben Taskar. *Introduction to Statistical Relational Learning*. MIT Press, September 2019. ISBN 978-0-262-53868-8.
- Rafael Gillmann. *0/1-Polytopes: Typical and Extremal Properties*. PhD thesis, February 2007. URL <https://depositonce.tu-berlin.de/items/urn:nbn:de:kobv:83-opus-14695>.
- Vincenzo Nicosia Giovanni Russo Vito Latora. *Complex Networks: Principles, Methods and Applications. With 58 exercises*. Cambridge University Press, Cambridge, United Kingdom ; New York, NY, September 2017. ISBN 978-1-107-10318-4.
- Ivan Glasser, Ryan Sweke, Nicola Pancotti, Jens Eisert, and Ignacio Cirac. Expressive power of tensor-network factorizations for probabilistic modeling. *Advances in Neural Information Processing Systems*, 32, 2019.
- Alex Goeßmann, Ingo Roth, Gitta Kutyniok, Michael Götte, Ryan Sweke, and Jens Eisert. Tensor network approaches for data-driven identification of non-linear dynamical laws. In *Advances in Neural Information Processing Systems - First Workshop on Quantum Tensor Networks in Machine Learning*, page 21, 2020.
- Alex Christoph Goeßmann. *Uniform Concentration of Tensor and Neural Networks: An Approach towards Recovery Guarantees*. PhD Thesis, Technische Universität Berlin, Berlin, 2021. URL <https://depositonce.tu-berlin.de/handle/11303/15990>. Accepted: 2021-12-30T15:00:58Z.
- Lars Grasedyck. Hierarchical Singular Value Decomposition of Tensors. *SIAM J. Matrix Analysis Applications*, 31: 2029–2054, January 2010. doi: 10.1137/090764189.
- W. Hackbusch and S. Kühn. A New Scheme for the Tensor Representation. *Journal of Fourier Analysis and Applications*, 15(5):706–722, October 2009. ISSN 1531-5851.
- Wolfgang Hackbusch. *Tensor Spaces and Numerical Tensor Calculus*. Springer Series in Computational Mathematics. Springer-Verlag, Berlin Heidelberg, 2012. ISBN 978-3-642-28026-9. doi: 10.1007/978-3-642-28027-6.
- Jean-Baptiste Hiriart-Urruty and Claude Lemarechal. *Convex Analysis and Minimization Algorithms II: Advanced Theory and Bundle Methods*. Springer, Berlin, Heidelberg, 1993rd edition edition, October 1993. ISBN 978-3-540-56852-0.
- Frank L. Hitchcock. The Expression of a Tensor or a Polyadic as a Sum of Products. *Journal of Mathematics and Physics*, 6(1-4):164–189, 1927. ISSN 1467-9590. doi: <https://doi.org/10.1002/sapm192761164>.
- Sepp Hochreiter. Toward a broad AI. *Communications of the ACM*, 65(4):56–57, January 2022. ISSN 0001-0782, 1557-7317. doi: 10.1145/3512715. URL <https://dl.acm.org/doi/10.1145/3512715>.
- Aidan Hogan, Eva Blomqvist, Michael Cochez, Claudia d’Amato, Gerard de Melo, Claudio Gutierrez, Sabrina Kirrane, Jose Emilio Labra Gayo, Roberto Navigli, and Sebastian Neumaier. *Knowledge Graphs*. Springer, Cham, 1st edition edition, November 2021. ISBN 978-3-031-00790-3.
- Sebastian Holtz, Thorsten Rohwedder, and Reinhold Schneider. On manifolds of tensors of fixed TT-rank. *Numerische Mathematik*, 120(4):701–731, April 2012. ISSN 0029-599X, 0945-3245. doi: 10.1007/s00211-011-0419-7. URL <http://link.springer.com/10.1007/s00211-011-0419-7>.
- Norm Jouppi, George Kurian, Sheng Li, Peter Ma, Rahul Nagarajan, Lifeng Nai, Nishant Patil, Suvinay Subramanian, Andy Swing, Brian Towles, Clifford Young, Xiang Zhou, Zongwei Zhou, and David A Patterson. TPU v4: An Optically Reconfigurable Supercomputer for Machine Learning with Hardware Support for Embeddings. In *Proceedings of the 50th Annual International Symposium on Computer Architecture*, ISCA ’23, pages 1–14, New York, NY, USA, June 2023. Association for Computing Machinery. ISBN 979-8-4007-0095-8. doi: 10.1145/3579371.3589350. URL <https://dl.acm.org/doi/10.1145/3579371.3589350>.



- Tamara G. Kolda and Brett W. Bader. Tensor Decompositions and Applications. *SIAM Review*, 51(3):455–500, August 2009. ISSN 0036-1445, 1095-7200. doi: 10.1137/07070111X.
- Daphne Koller and Nir Friedman. *Probabilistic Graphical Models: Principles and Techniques*. The MIT Press, Cambridge, Mass., 1. edition edition, July 2009. ISBN 978-0-262-01319-2.
- N’Dah Jean Kouagou, Stefan Heindorf, Caglar Demir, and Axel-Cyrille Ngonga Ngomo. Neural Class Expression Synthesis, December 2022. URL <http://arxiv.org/abs/2111.08486>. arXiv:2111.08486 [cs].
- N’Dah Jean Kouagou, Stefan Heindorf, Caglar Demir, and Axel-Cyrille Ngonga Ngomo. Neural Class Expression Synthesis. In Catia Pesquita, Ernesto Jimenez-Ruiz, Jamie McCusker, Daniel Faria, Mauro Dragoni, Anastasia Dimou, Raphael Troncy, and Sven Hertling, editors, *The Semantic Web*, volume 13870, pages 209–226. Springer Nature Switzerland, Cham, 2023. ISBN 978-3-031-33454-2 978-3-031-33455-9. doi: 10.1007/978-3-031-33455-9\_13. URL [https://link.springer.com/10.1007/978-3-031-33455-9\\_13](https://link.springer.com/10.1007/978-3-031-33455-9_13). Series Title: Lecture Notes in Computer Science.
- J. Landsberg. *Tensors: Geometry and Applications*, volume 128 of *Graduate Studies in Mathematics*. American Mathematical Society, December 2011. ISBN 978-0-8218-6907-9 978-0-8218-8481-2 978-0-8218-8483-6 978-1-4704-0923-4.
- Jens Lehmann, Sören Auer, Lorenz Bühmann, and Sebastian Tramp. Class expression learning for ontology engineering. *Journal of Web Semantics*, 9(1):71–81, March 2011. ISSN 1570-8268. doi: 10.1016/j.websem.2011.01.001. URL <https://www.sciencedirect.com/science/article/pii/S1570826811000023>.
- David J. C. MacKay. *Information Theory, Inference and Learning Algorithms*. Cambridge University Press, Cambridge, illustrated edition edition, September 2003. ISBN 978-0-521-64298-9.
- Alan K. Mackworth. Consistency in networks of relations. *Artificial Intelligence*, 8(1):99–118, February 1977. ISSN 0004-3702. doi: 10.1016/0004-3702(77)90007-8. URL <https://www.sciencedirect.com/science/article/pii/0004370277900078>.
- Giuseppe Marra, Sebastijan Dumančić, Robin Manhaeve, and Luc De Raedt. From statistical relational to neurosymbolic artificial intelligence: A survey. *Artificial Intelligence*, 328:104062, March 2024. ISSN 0004-3702. doi: 10.1016/j.artint.2023.104062. URL <https://www.sciencedirect.com/science/article/pii/S0004370223002084>.
- John McCarthy. Programs with Common Sense. In *Proceedings of the Teddington Conference on the Mechanization of Thought Processes*, pages 75–91. Her Majesty’s Stationary Office, London, 1959. URL <http://www-formal.stanford.edu/jmc/mcc59.html>.
- Stephen Muggleton and Luc De Raedt. Inductive logic programming: Theory and methods. *The Journal of Logic Programming*, 19:629–679, 1994. URL <https://www.sciencedirect.com/science/article/pii/0743106694900353>. Publisher: Elsevier.
- Kevin P. Murphy. *Probabilistic Machine Learning: An Introduction*. The MIT Press, Cambridge, Massachusetts London, England, March 2022. ISBN 978-0-262-04682-4.
- Maximilian Nickel, Volker Tresp, and Hans-Peter Kriegel. A three-way model for collective learning on multi-relational data. In *Proceedings of the 28th International Conference on International Conference on Machine Learning*, ICML’11, pages 809–816, Madison, WI, USA, June 2011. Omnipress. ISBN 978-1-4503-0619-5.
- Maximilian Nickel, Kevin Murphy, Volker Tresp, and Evgeniy Gabrilovich. A Review of Relational Machine Learning for Knowledge Graphs. *Proceedings of the IEEE*, 104(1):11–33, January 2016. ISSN 0018-9219, 1558-2256. doi: 10.1109/JPROC.2015.2483592. URL <https://ieeexplore.ieee.org/document/7358050/>.
- Goran S. Nikolić, Bojan R. Dimitrijević, Tatjana R. Nikolić, and Mile K. Stojcev. A Survey of Three Types of Processing Units: CPU, GPU and TPU. In *2022 57th International Scientific Conference on Information, Communication and Energy Systems and Technologies (ICEST)*, pages 1–6, June 2022. doi: 10.1109/ICEST55168.2022.9828625. URL <https://ieeexplore.ieee.org/document/9828625>.
- Román Orús. Tensor networks for complex quantum systems. *Nature Reviews Physics*, 1(9):538–550, September 2019. ISSN 2522-5820. doi: 10.1038/s42254-019-0086-7.
- I. V. Oseledets and E. E. Tyrtshnikov. Breaking the Curse of Dimensionality, Or How to Use SVD in Many Dimensions. *SIAM Journal on Scientific Computing*, 31(5):3744–3759, January 2009. ISSN 1064-8275. doi: 10.1137/090748330. Publisher: Society for Industrial and Applied Mathematics.
- Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Köpf, Edward Yang, Zach DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. PyTorch:

- An Imperative Style, High-Performance Deep Learning Library, December 2019. URL <http://arxiv.org/abs/1912.01703>. arXiv:1912.01703 [cs].
- Judea Pearl. *Probabilistic Reasoning in Intelligent Systems: Networks of Plausible Inference*. Morgan Kaufmann, s.l., September 1988. ISBN 978-1-55860-479-7.
- Judea Pearl. *Causality: Models, Reasoning and Inference*. Ausgezeichnet: ACM Turing Award for Transforming Artificial Intelligence 2011. Cambridge University Press, Cambridge New York, NY Port Melbourne New Delhi Singapore, 2 edition, November 2009. ISBN 978-0-521-89560-6.
- Roger Penrose. *Spinors and Space-Time: Volume 1, Two-Spinor Calculus and Relativistic Fields*. Cambridge University Press, Cambridge, February 1987. ISBN 978-0-521-33707-6.
- D. Perez-Garcia, F. Verstraete, M. M. Wolf, and J. I. Cirac. Matrix product state representations. *Quantum Information & Computation*, 7(5):401–430, July 2007. ISSN 1533-7146.
- Matthew Richardson and Pedro Domingos. Markov logic networks. *Machine Learning*, 62(1-2):107–136, February 2006. ISSN 0885-6125, 1573-0565. doi: 10.1007/s10994-006-5833-1. URL <http://link.springer.com/10.1007/s10994-006-5833-1>.
- Elina Robeva and Anna Seigal. Duality of graphical models and tensor networks. *Information and Inference: A Journal of the IMA*, 8(2):273–288, June 2019. ISSN 2049-8772. doi: 10.1093/imaiai/iy009.
- Ralph Tyrell Rockafellar. *Convex Analysis*. Princeton University Press, Princeton, reprint edition edition, January 1997. ISBN 978-0-691-01586-6.
- Sebastian Rudolph. Foundations of Description Logics. In Axel Polleres, Claudia d’Amato, Marcelo Arenas, Siegfried Handschuh, Paula Kroner, Sascha Ossowski, and Peter Patel-Schneider, editors, *Reasoning Web. Semantic Technologies for the Web of Data: 7th International Summer School 2011, Galway, Ireland, August 23-27, 2011, Tutorial Lectures*, pages 76–136. Springer, Berlin, Heidelberg, 2011. ISBN 978-3-642-23032-5. doi: 10.1007/978-3-642-23032-5\_2. URL [https://doi.org/10.1007/978-3-642-23032-5\\_2](https://doi.org/10.1007/978-3-642-23032-5_2).
- Stuart Russell and Peter Norvig. *Artificial Intelligence: A Modern Approach, Global Edition: A Modern Approach, Global Edition*. Pearson, Boston, 4 edition, May 2021. ISBN 978-1-292-40113-3.
- Chiaki Sakama, Katsumi Inoue, and Taisuke Sato. Linear Algebraic Characterization of Logic Programs. In Gang Li, Yong Ge, Zili Zhang, Zhi Jin, and Michael Blumenstein, editors, *Knowledge Science, Engineering and Management*, volume 10412, pages 520–533. Springer International Publishing, Cham, 2017. ISBN 978-3-319-63557-6 978-3-319-63558-3. doi: 10.1007/978-3-319-63558-3\_44. URL [http://link.springer.com/10.1007/978-3-319-63558-3\\_44](http://link.springer.com/10.1007/978-3-319-63558-3_44). Series Title: Lecture Notes in Computer Science.
- Aaron Sander, Maximilian Fröhlich, Martin Eigel, Jens Eisert, Patrick Gelß, Michael Hintermüller, Richard M. Milbradt, Robert Wille, and Christian B. Mendl. Large-scale stochastic simulation of open quantum systems, January 2025. URL <http://arxiv.org/abs/2501.17913>. arXiv:2501.17913 [quant-ph].
- Md Kamruzzaman Sarker, Lu Zhou, Aaron Eberhart, and Pascal Hitzler. Neuro-symbolic artificial intelligence: Current trends. *AI Communications*, 34(3):197–209, March 2022. ISSN 18758452, 09217126. doi: 10.3233/AIC-210084. URL <https://www.medra.org/servlet/aliasResolver?alias=iospress&doi=10.3233/AIC-210084>.
- Taisuke Sato. A linear algebraic approach to datalog evaluation. *Theory and Practice of Logic Programming*, 17(3):244–265, May 2017. ISSN 1471-0684, 1475-3081. doi: 10.1017/S1471068417000023. URL <https://www.cambridge.org/core/journals/theory-and-practice-of-logic-programming/article/abs/linear-algebraic-approach-to-datalog-evaluation/CED3EEB903D9D8A16843CFC5AC4D577>. Publisher: Cambridge University Press.
- Luciano Serafini and Artur S. d’Avila Garcez. Learning and Reasoning with Logic Tensor Networks. In *AI\*IA 2016 Advances in Artificial Intelligence: XVth International Conference of the Italian Association for Artificial Intelligence, Genova, Italy, November 29 – December 1, 2016, Proceedings*, pages 334–348, Berlin, Heidelberg, November 2016. Springer-Verlag. ISBN 978-3-319-49129-5. doi: 10.1007/978-3-319-49130-1\_25. URL [https://doi.org/10.1007/978-3-319-49130-1\\_25](https://doi.org/10.1007/978-3-319-49130-1_25).
- Shalev-Schwartz, Shai and Ben-David, Shai. *Understanding Machine Learning: From Theory to Algorithms*. Cambridge University Press, New York, NY, USA, July 2014. ISBN 978-1-107-05713-5.
- C. E. Shannon. A Mathematical Theory of Communication. *Bell System Technical Journal*, 27(3):379–423, 1948. ISSN 1538-7305. doi: 10.1002/j.1538-7305.1948.tb01338.x. URL <https://onlinelibrary.wiley.com/doi/abs/10.1002/j.1538-7305.1948.tb01338.x>. \_eprint: <https://onlinelibrary.wiley.com/doi/pdf/10.1002/j.1538-7305.1948.tb01338.x>.

- Helmut Simonis. Sudoku as a constraint problem. In *CP Workshop on modeling and reformulating Constraint Satisfaction Problems*, volume 12, pages 13–27. Citeseer Sitges, Spain, 2005. URL [https://ai.dmi.unibas.ch/\\_files/teaching/fs21/ai/material/ai26-simonis-cp2005ws.pdf](https://ai.dmi.unibas.ch/_files/teaching/fs21/ai/material/ai26-simonis-cp2005ws.pdf).
- Edwin Stoudenmire and David J Schwab. Supervised Learning with Tensor Networks. In D. D. Lee, M. Sugiyama, U. V. Luxburg, I. Guyon, and R. Garnett, editors, *Advances in Neural Information Processing Systems 29*, pages 4799–4807. Curran Associates, Inc., 2016.
- Michel Talagrand. *Upper and Lower Bounds for Stochastic Processes: Modern Methods and Classical Problems*. Springer, Berlin, Heidelberg, 2014. ISBN 978-3-642-54074-5. doi: 10.1007/978-3-642-54075-2.
- Geoffrey G. Towell and Jude W. Shavlik. Knowledge-based artificial neural networks. *Artificial Intelligence*, 70(1):119–165, October 1994. ISSN 0004-3702. doi: 10.1016/0004-3702(94)90105-8. URL <https://www.sciencedirect.com/science/article/pii/0004370294901058>.
- Théo Trouillon and Maximilian Nickel. Complex and Holographic Embeddings of Knowledge Graphs: A Comparison, July 2017. URL <http://arxiv.org/abs/1707.01475>. arXiv:1707.01475 [cs, stat].
- Efthimis Tsilionis, Alexander Artikis, and Georgios Paliouras. A Tensor-Based Formalization of the Event Calculus. In *Proceedings of the Thirty-Third International Joint Conference on Artificial Intelligence*, pages 3584–3592, Jeju, South Korea, August 2024. International Joint Conferences on Artificial Intelligence Organization. ISBN 978-1-956792-04-1. doi: 10.24963/ijcai.2024/397. URL <https://www.ijcai.org/proceedings/2024/397>.
- Roman Vershynin. *High-Dimensional Probability: An Introduction with Applications in Data Science*. Cambridge University Press, New York, NY, 1st edition edition, September 2018. ISBN 978-1-108-41519-4.
- Martin J. Wainwright. *High-Dimensional Statistics: A Non-Asymptotic Viewpoint*. Cambridge Series in Statistical and Probabilistic Mathematics. Cambridge University Press, Cambridge, 2019. ISBN 978-1-108-49802-9. doi: 10.1017/9781108627771.
- Martin J. Wainwright and Michael Irwin Jordan. *Graphical Models, Exponential Families, and Variational Inference*. Now Publishers Inc, 2008. ISBN 978-1-60198-184-4.
- Stephen Wolfram. Statistical mechanics of cellular automata. *Reviews of Modern Physics*, 55(3):601–644, July 1983. doi: 10.1103/RevModPhys.55.601. URL <https://link.aps.org/doi/10.1103/RevModPhys.55.601>. Publisher: American Physical Society.
- Stephen Wolfram. *A New Kind of Science*. Wolfram Media, Champaign (Ill.), illustrated edition edition, May 2002. ISBN 978-1-57955-008-0.
- Bishan Yang, Wen-tau Yih, Xiaodong He, Jianfeng Gao, and Li Deng. Embedding Entities and Relations for Learning and Inference in Knowledge Bases. arXiv, August 2015. URL <http://arxiv.org/abs/1412.6575>. arXiv:1412.6575 [cs].
- Günter M. Ziegler. Lectures on 0/1-Polytopes. In Gil Kalai and Günter M. Ziegler, editors, *Polytopes — Combinatorics and Computation*, pages 1–41. Birkhäuser, Basel, 2000. ISBN 978-3-0348-8438-9. doi: 10.1007/978-3-0348-8438-9\_1. URL [https://doi.org/10.1007/978-3-0348-8438-9\\_1](https://doi.org/10.1007/978-3-0348-8438-9_1).
- Günter M. Ziegler. *Lectures on Polytopes*. Springer, New York, 1995th edition edition, October 2013. ISBN 978-0-387-94365-7.