
THE TENSOR NETWORK APPROACH TO EFFICIENT AND EXPLAINABLE AI

RESEARCH NOTES IN THE ENEXA PROJECT

Alex Goessmann, DATEV eG

March 6, 2025

ABSTRACT

Tensors appear naturally in factored representations of systems, when storing information about a systems state. Since the curse of dimensionality prevents feasible generic representations and reasoning, logical and probabilistic reasoning focuses on tradeoffs between efficiency and generality. In this work we present these tradeoffs based on the tensor network formalism and formulate feasible reasoning algorithms based on 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 parts we investigate in more detail schemes to exploit tensor network contractions for calculus.

Contents

1	Introduction	9
1.1	Application of the tensor network approach	10
1.1.1	Graphical Models	10
1.1.2	Neuro-Symbolic AI	10
1.2	Representation schemes of systems	10
1.2.1	Ontological Commitments	10
1.2.2	Logical and Probabilistic Approaches	10
1.3	Literature	11
1.3.1	Broad topics	11
1.3.2	Tensors Representations	12
1.4	Outline	12
2	Notation and Basic Concepts	12
2.1	Categorical Variables and Representations	13
2.2	Tensors	13
2.3	One-hot encodings	13
2.4	Contractions	14
2.4.1	Graphical Illustrations	14
2.4.2	Tensor Product	15

2.4.3	Generic Contractions	16
2.4.4	Decompositions	18
2.5	Properties of Tensors	18
2.6	Encoding schemes for functions	18
2.6.1	Relational encodings	18
2.6.2	Tensor-valued functions	19
I	Decomposition and Inference of Factored Representations	19
3	Probability Distributions	19
3.1	Tensor Representation of Distributions	20
3.1.1	Base measures	20
3.2	Marginal Distribution	21
3.3	Conditional Probabilities	22
3.4	Bayes Theorem and the Chain Rule	23
3.5	Independent Variables	23
3.6	Graphical Models	25
3.6.1	Markov Networks	25
3.6.2	Bayesian Networks	26
3.6.3	Bayesian Networks as Markov Networks	28
3.6.4	Example: Hidden Markov Models	29
3.7	Exponential Families	29
3.7.1	Tensor Network Representation	30
3.7.2	Mean Parameters	31
3.7.3	Examples	33
3.8	Empirical Distributions	34
3.9	Discussion and Outlook	35
4	Probabilistic Reasoning	35
4.1	Queries	36
4.1.1	Querying by functions	36
4.1.2	MAP Queries	36
4.1.3	Answering queries by energy contractions	37
4.2	Sampling based on queries	37
4.2.1	Exact Methods	37
4.2.2	Approximate Methods	38
4.2.3	Simulated Annealing	38
4.3	Maximum Likelihood Estimation	38
4.3.1	Likelihood and Loss	38
4.3.2	Entropic Interpretation	39

4.4	Forward Mapping in Exponential Families	40
4.4.1	Variational Formulation	40
4.4.2	Boundary of convex polytopes	42
4.4.3	Mode Search by annealing	43
4.4.4	Mean Field Method	43
4.4.5	Structured Variational Approximation	44
4.5	Backward Mapping in Exponential Families	46
4.5.1	Variational Formulation	47
4.5.2	Interpretation by Maximum Likelihood Estimation	47
4.5.3	Connection with Maximum Entropy	48
4.5.4	Alternating Algorithms to Approximate the Backward Map	48
4.6	Discussion	50
5	Propositional Logics	50
5.1	Encoding of Booleans	50
5.1.1	Booleans as categorical variables	50
5.1.2	One-hot Encoding	50
5.2	Semantics of Propositional Formulas	50
5.2.1	Formulas	51
5.2.2	Relational encoding of formulas	51
5.3	Syntax of Propositional Formulas	52
5.3.1	Atomic Formulas	52
5.3.2	Syntactical combination of formulas	53
5.3.3	Syntactical decomposition of formulas	55
5.4	Discussion and Outlook	57
6	Logical Inference	57
6.1	Entailment in Propositional Logics	57
6.1.1	Deciding Entailment by contractions	57
6.1.2	Contraction Knowledge Base	58
6.1.3	Sparse Representation of a Knowledge Base	58
6.2	Formulas as Random Variables	58
6.2.1	Conditioning on the atoms	59
6.2.2	Conditioning on the formula	59
6.2.3	Probability of a function given a Knowledge Base	60
6.2.4	Knowledge Bases as Base Measures for Probability Distributions	61
6.2.5	Deciding entailment on Markov Networks	62
6.3	Deciding Entailment by partial ordering	62
6.3.1	Monotonicity of Entailment	63
6.4	Deciding Entailment by local contractions	63

6.4.1	Knowledge Propagation	63
II	Neuro-Symbolic Learning	64
7	Formula Selecting Networks	64
7.1	Construction schemes	65
7.1.1	Connective Selecting Tensors	65
7.1.2	Variable Selecting Tensor Network	65
7.2	State Selecting Tensors	66
7.3	Composition of formula selecting maps	67
7.3.1	Formula Selecting Neuron	67
7.3.2	Formula Selecting Neural Network	68
7.4	Application of Formula Selecting Networks	69
7.4.1	Representation of selection encodings	69
7.4.2	Efficient Representation of Formulas	70
7.4.3	Batch contraction of parametrized formulas	70
7.4.4	Average contraction of parametrized formulas	70
7.5	Examples of formula selecting neural networks	71
7.5.1	Correlation	71
7.5.2	Conjunctive and Disjunctive Normal Forms	71
7.6	Extension to variables of larger dimension	72
8	Logic Network Representation	72
8.1	Markov Logic Networks	72
8.1.1	Markov Logic Networks as Exponential Families	72
8.1.2	Tensor Network Representation	73
8.1.3	Energy tensors	74
8.1.4	Expressivity	75
8.1.5	Distribution of independent variables	76
8.1.6	Boltzmann machines	77
8.2	Hard Logic Networks	77
8.2.1	The limit of hard logic	78
8.2.2	Tensor Network Representation	79
8.2.3	Polynomial Representation	79
8.2.4	Categorical Constraints	80
8.3	Hybrid Logic Network	81
8.3.1	Tensor Network Representation	83
8.3.2	Reasoning Properties	84
8.3.3	Expressivity	84
8.4	Applications	85

9	Logic Network Inference	85
9.1	Mean parameters of Hybrid Logic Networks	85
9.1.1	Extreme points by hard logic networks	86
9.1.2	Mean parameters in the interior	86
9.1.3	Mean parameters outside the interior	86
9.1.4	Expressivity of Hybrid Logic Networks	88
9.1.5	Case of tree computation networks	89
9.2	Entropic Motivation of unconstrained Parameter Estimation	89
9.2.1	Maximum Likelihood in Hybrid Logic Networks	89
9.2.2	Maximum Entropy in Hybrid Logic Networks	90
9.3	Alternating Algorithms to Approximate the Backward Map	90
9.4	Forward and backward mappings in closed form	92
9.4.1	Maxterms and Minterms	92
9.4.2	Atomic formulas	93
9.5	Constrained parameter estimation in the minterm family	94
9.5.1	Parameter Estimation	94
9.5.2	Structure Learning	95
9.6	Greedy Structure Learning	95
9.6.1	Greedy formula inclusions	95
9.6.2	Gain Heuristic	95
9.6.3	Gradient heuristic and the proposal distribution	97
9.6.4	Iterations	98
9.7	Proposal distribution	98
9.7.1	Mean parameter polytope	99
9.8	Discussion	99
10	Probabilistic Success Guarantees	99
10.1	Fluctuations	99
10.1.1	Fluctuation of the empirical distribution	99
10.1.2	Fluctuation tensors	100
10.1.3	Widths of random tensors	100
10.2	Fluctuations in Markov Logic Networks	100
10.2.1	Energy tensor in proposal distributions	101
10.2.2	Naive Exponential Family	101
10.3	Expected and Empirical Risk Minimization	102
10.3.1	Maximum Likelihood Estimation on random data	102
10.3.2	Solution of the Expected Problem	102
10.3.3	Recovery Guarantee based on Widths	103
10.4	Guarantees for Mode of the Proposal Distribution	104
10.5	Guarantees for Parameter Estimation	104

11 Statistical Models of First Order Logic (FOL)	105
11.1 World Tensors	105
11.1.1 Special Case of Propositional Logics	106
11.1.2 Enumeration of worlds	106
11.1.3 Random World Tensors	106
11.2 Formulas in first order logics	106
11.2.1 Syntactical Decomposition of quantifier-free formulas	107
11.2.2 Quantifiers	107
11.2.3 Basis CP Decomposition	108
11.2.4 Example: Queries	109
11.3 Representation of Knowledge Graphs	109
11.3.1 Representation as unary and binary predicates	110
11.3.2 Representation as ternary predicate	110
11.3.3 SPARQL Queries	110
11.4 Probabilistic Relational Models	111
11.4.1 Markov Logic Networks in FOL	111
11.4.2 Importance Formula	112
11.4.3 Decomposition of the log likelihood	113
11.4.4 Interpretation as Likelihood of Propositional Dataset	113
11.4.5 Approximation by Independent Samples	115
11.4.6 Extraction of Samples from FOL Knowledge Bases	115
11.4.7 Representation by Tensor Networks	115
11.4.8 Basis CP Decomposition of extracted data	116
11.4.9 Representation with auxiliary term variables	117
11.4.10 Generic Tensor Network Decomposition of Extracted Data	117
11.4.11 Design of the Formulas	117
11.5 Generation of FOL worlds	117
11.5.1 Samples by single objects	117
11.6 Samples by pairs of objects	118
11.7 Example: Generation of Knowledge Graphs	118
11.7.1 Samples by single resources	118
11.7.2 Samples by pairs of resources	119
11.7.3 General orthogonal projection approach	119
11.8 Discussion	120
 III Contraction Calculus	 120
12 Coordinate Calculus	120
12.1 One-hot encodings as basis	120

12.2	Coordinatewise Transforms	121
12.3	Differentiation of Contraction	121
12.4	Selection Encodings	123
12.4.1	Tensors as linear maps	123
12.4.2	Selection encodings	123
12.5	Discussion	124
13	Directed Tensor Calculus	125
13.1	Directed Tensors	125
13.1.1	Normation	126
13.1.2	Contraction of Directed Tensors	126
14	Contraction Equations	127
14.1	Contraction equations in logical and probabilistic reasoning	127
14.2	Normation Equations	127
14.3	Proof of Hammersley-Clifford Theorem	128
14.4	Commutation of Contractions	130
14.5	Support of Contractions	131
15	Basis Calculus	132
15.1	Encoding of Subsets and Relations	132
15.1.1	Higher order relations	133
15.2	Encoding of Functions	133
15.2.1	Relational Encoding of Functions	133
15.3	Calculus of relational encodings	134
15.3.1	Function Evaluation	134
15.3.2	Composition of function	135
15.3.3	Compositions with real functions	136
15.3.4	Decomposition in case of structured images	137
15.4	Effective Coordinate Calculus	138
15.5	Applications in Machine Learning	139
16	Contraction Message Passing	139
16.1	Exact Contractions	139
16.1.1	Construction of Cluster Graphs	140
16.1.2	Message Passing to calculate contractions	140
16.1.3	Variable Elimination Cluster Graphs	141
16.1.4	Bethe Cluster Graphs	141
16.1.5	Computational Complexity	141
16.2	Approximate Contractions	142
16.2.1	Exact Message Passing for Directed and Binary Contractions	142

16.2.2	Case of Matrices	142
16.2.3	Case of Tensors	143
16.3	Basis Calculus	143
16.4	Applications	143
17	Sparse Tensor Representations	143
17.1	CP Formats	143
17.1.1	Directed Leg Cores	144
17.1.2	Basis Leg Cores	144
17.1.3	Basis+ Leg Cores	145
17.2	Constructive Bounds on CP Ranks	147
17.2.1	Format Transformations	147
17.2.2	Summation of CP Decompositions	148
17.2.3	Contractions of CP Decompositions	148
17.2.4	Normations of CP Decompositions	149
17.2.5	Sparse Encoding of Functions	149
17.2.6	Construction by averaging the incoming legs	150
17.3	Representation by slice selection architectures	150
17.3.1	Applications	151
17.4	Optimization of sparse tensors	151
17.4.1	Mode search in exponential families	151
17.4.2	Higher-Order Unconstrained Binary Optimization (HUBO)	152
17.4.3	Quadratic Unconstrained Binary Optimization (QUBO)	152
17.4.4	Integer Linear Programming	153
17.5	Value subspaces of functions	155
17.5.1	Propositional Formula Subspaces	155
17.5.2	Formula Decomposition as a Subspace Choice	155
18	Reasoning by Tensor Approximation	155
18.1	Approximation of Energy tensors	155
18.1.1	Direct Approximation	155
18.1.2	Approximation involving Selection Architectures	156
18.2	Transformation of Maximum Search to Risk Minimization	156
18.2.1	Weighted Squares Loss Trick	156
18.2.2	Problem of the trivial tensor	157
18.3	Alternating Solution of Least Squares Problems	157
18.3.1	Choice of Representation Format	157
18.4	Regularization and Compressed Sensing	158
19	Uniform Concentration of Random Contractions	158
19.1	Naive bounds given binomial coordinates	158

19.1.1	Basis Vectors	159
19.1.2	Sphere	159
19.1.3	Bounds based on the sub-gaussian norm	160
19.2	Chaining bounds given binomial coordinates	160
19.2.1	Generic Width Bounds	162
A	Implementation in the tnreason package	162
A.1	Architecture	162
A.2	Subpackage engine	163
A.2.1	Contraction Calculus	163
A.2.2	Cores and Contractions	163
A.3	Subpackage encoding	165
A.3.1	Script Language	165
A.3.2	Core Nomenclature	167
A.3.3	Color Nomenclature	167
A.3.4	Relational encoding of formulas	167
A.3.5	Representation of MLNs	168
A.3.6	Formula Selecting Maps	168
A.4	Subpackage algorithms	168
A.4.1	Alternating Least Squares	168
A.4.2	Gibbs Sampling	168
A.4.3	Knowledge Propagation	169
A.4.4	Energy-based Algorithms	169
A.5	Subpackage knowledge	169
A.5.1	Distributions	169
A.5.2	Inference	169
A.5.3	Parameter Estimation	170
A.5.4	Structure Learning	170

1 Introduction

Reasoning based on tensors has the advantage to

- Provide a unifying approach to logical and probabilistic reasoning
- Define necessary workload of algorithms by contractions

The storage of numeric information about the states of a system with multiple variables is naturally a tensor. Tensors appear naturally in

- Logics: Boolean tensors indicating models (propositional case) and interpretation tensors in first order logics
- Probability theory: Truth tables, which are tensors of probabilities for joint distributions of categorical variables.

The curse of dimensionality renders the representations of tensors by a list of basis elements infeasible. Therefore, sparse representation formats by tensor networks have been investigated.

Tensor network decompositions as representation schemes appear in

- Logics: Conjunctions of formulas are Hadamard products of the tensor representation of formulas (Coordinate Calculus/ Effective Calculus)
- Probability theory: Graphical models are tensor networks of the factors
- Data bases: Relations encoded by lists as storage of nonvanishing coordinates of a relation encoding

Tensor network contractions as reasoning schemes appear in

- Logics: Model counts, used for satisfiability decisions and entailment
- Probability theory: Marginal probability distributions, extended to conditional probability distributions through normalizations

1.1 Application of the tensor network approach

1.1.1 Graphical Models

Further sparsity schemes for factors, by exploiting structure inside them.

1.1.2 Neuro-Symbolic AI

Tensor Approaches to Neuro-Symbolic AI

- TensorLog Cohen et al. (2020)
- Badreddine et al. (2022) representation of logic using tensor networks and automated differentiation to optimize.

In Deep Neural Networks, functions between the input layer and the output layer are decomposed into neurons. Typical neurons are linear transforms with an activation function.

Sparsity means restriction to functions, which are decomposability into a small number of neurons. Approximations of generic functions (see the universal approximation theorems) would require When restricting to functions based on a fixed architecture, the sparsity of functions

1.2 Representation schemes of systems

Following the book Russell and Norvig (2021): Atomic, Factored and Structured Representations, which are framed as ontological commitments.

Tensors have frequent appearance in machine learning. We show in this chapter how they are natural choices to encode states the states of factored systems, which are systems described by a collection of categorical variables.

1.2.1 Ontological Commitments

Differing ontological commitments are made by atomic systems, where we just enumerate the state of a system and understand them as assignments to a single variable. We can always transform a factored representation of a system to an atomic one, just by enumerating the tuples of states of the factored system and interpreting them as a single variable. 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. In structured systems the numbers of variables can differ depending on the state the system is in. Limitation to factored representations: Dealing with structured representations would mean differing order of the representing tensors.

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.

1.2.2 Logical and Probabilistic Approaches

Besides ontological commitments in the choice of a representation scheme, one also needs to make epistemologic commitments, by defining what properties we can reason about. We can encode different properties of the system in

the axes of the tensor representation. In logical approaches the properties of states are Boolean values representing, for example, whether a state is consistent with specific constraints. Probabilistic approaches represent on the coordinates of the tensors real numbers in $[0, 1]$ which can be interpreted as probability of a state. Compared with logical approaches to reasoning, probabilistic approaches thus bear a more expressive modelling.

Tensor Networks can represent both probabilistic and logical reasoning on factored systems, since in the tensor space we can

- Represent probability distributions by storing probabilistic values in each coordinate
- Represent set of states by sums over different one-hot encodings (this enables logical calculus)

Both probability and logic provide a human-understandable interface to machine learning. As we will describe in the following sections, they can be combined in one formalism providing efficient reasoning. This formalism of tensors along their network decompositions and contractions bears the potential of parallel computations exploited in the AI-dedicated soft and hardware.

Probability represents the uncertainty of states. The categorical variables are called random variables and their joint distribution is represented by a probability tensor. Humans can 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.

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 mathematical structures to represent properties of states in factored systems. The potential is then based in employing efficient multilinear algorithms based on tensor network alternations 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.

1.3 Literature

1.3.1 Broad topics

Inductive Logic Programming:

- ILP is a classical task Muggleton and De Raedt (1994)
- Amie Galárraga et al. (2013) is a method of learning Horn clauses using a refinement operator.
- Class Expression Learning Lehmann et al. (2011) is a more recent approach which has recently seen further popularity in combination with reinforcement learning Demir and Ngonga Ngomo (2021) and neural networks Kouagou et al. (2022, 2023)

Statistical Relational AI: Getoor and Taskar (2019)

- Classical combination of logical and probabilistic approaches to reasoning

Neurosymbolic AI

- Required for more advanced AI Hochreiter (2022)

- Add the paradigm of neural computing to logical reasoning
- Potential benefits from Statistical Relational AI Marra et al. (2024)
- Tensor based approaches Cohen et al. (2020)
- Badreddine et al. (2022) representation of logic using tensor networks and automated differentiation to optimize.

Knowledge Graphs Hogan et al. (2021)

- The advent of large Knowledge Graphs enables reasoning methods.
- The bottleneck is the efficiency of methods to identify formulas from data.
- Knowledge Graphs are stored in a sparse format, i.e. only true atoms instead of all + truth label.

1.3.2 Tensors Representations

Tensor Representation is natural, if systems are described in factored representation.

Tensor Network formats

- HT Format ?
- CP Format

Tensor Representation of Logics

- Tensor Networks have been applied in the automatization of logic reasoning Sakama et al. (2017); Sato (2017) apply Matrix multiplication in reasoning.
- Nickel et al. (2016) review over relational machine learning and latent features via matrix embeddings.

Tensor Representation of Knowledge Graphs

- Effective representation of queries
- Usage of tensor networks in embeddings Yang et al. (2015) and using complex extensions Trouillon and Nickel (2017); Trouillon (2017)

Tensor Representation of Graphical Models

- Duality of Graphical Models and Tensor Networks: Robeva and Seigal (2019)
- Expressivity studies Glasser et al. (2019)

Tensor Networks as Regressors

- Dynamical Systems learning Gelß et al. (2019); Goeßmann et al. (2020)

1.4 Outline

Part I: Review of classical approaches using tensor calculus

Part II: Neuro-symbolic learning

Part III: Numerical and Stochastic investigation of the calculus with tensors

2 Notation and Basic Concepts

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

2.1 Categorical Variables and Representations

We will in this work investigate systems, which are described by a set of properties, each called a categorical variable. 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 $[m_k]$.*

2.2 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 $[m_k]$. We call maps*

$$T[X_0, \dots, X_{d-1}] : \prod_{k \in [d]} [m_k] \rightarrow \mathbb{R}$$

tensor of order d and leg dimensions m_0, \dots, m_{d-1} . Evaluations of these maps at indices x_0, \dots, x_{d-1} are denoted by

$$T[X_0 = x_0, \dots, X_{d-1} = x_{d-1}] = T[X_0, \dots, X_{d-1}](x_0, \dots, x_{d-1}) .$$

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

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

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

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 $T[X_0, \dots, X_{d-1}]$ by $T[X_{[d]}]$. 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 is defined as the map*

$$\mathbb{I}[X_{[d]}] : \prod_{k \in [d]} [m_k] \rightarrow \{1\} \subset \mathbb{R}$$

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

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

2.3 One-hot encodings

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 4 (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 $e_x[X]$ by*

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

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

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

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

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

$$e_x[X] = [0 \quad \cdots \quad 0 \quad 1 \quad 0 \quad \cdots \quad 0], \quad (2)$$

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 5 (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*

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

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

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

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

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

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.*

2.4 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.

2.4.1 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 6. *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 1a). 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.

Both depiction schemes are simplifications of factor graphs, by highlighting the categorical variables in the depiction in Figure 1b) and the tensors in the depiction in Figure 1c). We in this work will prefer the simplification of the tensor network community, depicted in Figure 1b).

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

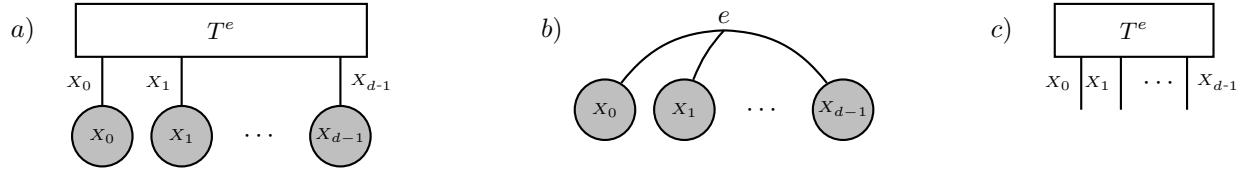


Figure 1: Depiction of Tensors a) As a factor in a factor graph, depicted by a block, and connected to categorical variables assigned to nodes. b) Highlighting only the variable dependencies by a hyperedge connecting the variables X_k to each axis $k \in [d]$. c) Highlighting the tensor by a blockwise 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 2. 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 7. 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

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

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

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

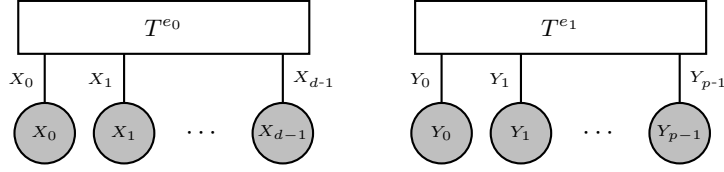
the Tensor Network of the decorated hypergraph \mathcal{G} .

2.4.2 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 2: Example of a tensor network. a) Hypergraph with edges $e_0 = \{X_0, X_1, X_2\}$, $e_1 = \{X_1, X_2\}$ and $e_2 = \{X_2, X_3\}$ decorated by tensor cores. b) Dual tensor network, depicting a contraction with leaving all variables open.



Definition 8 (Tensor Product). *Let there be two tensor*

$$T^{e_0} [X_{[d]}] : \prod_{k \in [d]} [m_k] \rightarrow \mathbb{R} \quad \text{and} \quad T^{e_1} [Y_{[p]}] : \prod_{l \in [r]} [m_l] \rightarrow \mathbb{R}$$

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

$$\langle T^{e_0} [X_{[d]}], T^{e_1} [Y_{[p]}] \rangle [X_{[d]}, Y_{[p]}] : \left(\prod_{k \in [d]} [m_k] \right) \times \left(\prod_{l \in [r]} [m_l] \right) \rightarrow \mathbb{R}$$

defined for $x_0, \dots, x_{d-1} \in \prod_{k \in [d]} [m_k]$ and $y_0, \dots, y_{p-1} \in \prod_{l \in [r]} [m_l]$ as

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

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

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

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

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

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

2.4.3 Generic Contractions

Contractions of Tensor Networks $\mathcal{T}^{\mathcal{G}}$ 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 9. *Let $\mathcal{T}^{\mathcal{G}}$ 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*

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

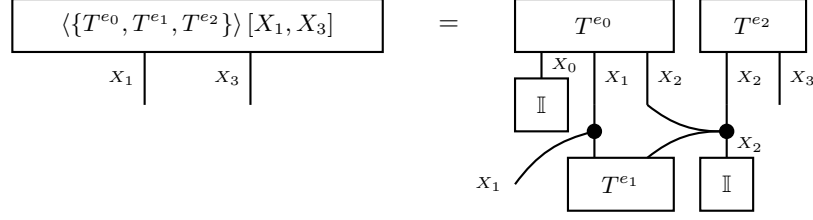


Figure 3: Example of a tensor network contraction of all but the variables X_1, X_3 . Contraction of variables can always be depicted by closing the open legs with trivial tensors \mathbb{I} performing index sums.

defined coordinatewise by the sum

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

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

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 Definition 9, 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 products by:

$$\overset{X_0}{\text{---}} \boxed{M} \overset{X_1}{\text{---}} \boxed{V} = \overset{X_0}{\text{---}} \boxed{\langle M[X_0, X_1], V[X_1] \rangle [X_0]}$$

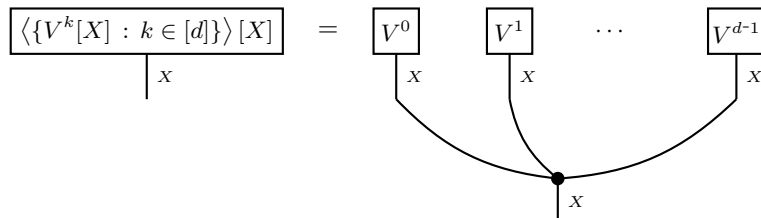
Example 3. Hadamard Products of Vectors A node appearing in arbitrary many hyperedges denotes a Hadamard product of the axis of the respective decorating tensors. To give an example, let $V^k[X] \in \mathbb{R}^m$ be vectors for $k \in [d]$. Their hadamard product is the vector

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

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



2.4.4 Decompositions

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

Definition 10. A Tensor Network Decomposition of a tensor $T[X_{\mathcal{V}}]$ is a Tensor Network $\mathcal{T}^{\mathcal{G}}$ such that

$$T[X_{\mathcal{V}}] = \langle \mathcal{T}^{\mathcal{G}} \rangle [X_{\mathcal{V}}] .$$

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

2.5 Properties of Tensors

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

Definition 11. We call a tensor $T[X_{[d]}]$ boolean, when $\text{im}(T) \subset [2]$, i.e. all coordinates are either 0 or 1.

Directionality represents constraints on the structure of tensors: Summing over outgoing trivializes the tensor.

Definition 12. A Tensor

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

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

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

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 normation.

Definition 13. A tensor $T[X_{\mathcal{V}}]$ 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

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

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

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

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 86 that normations 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:



2.6 Encoding schemes for functions

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.

2.6.1 Relational 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 15, see Definition 76.

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

$$f : \prod_{k \in [d]} [m_k] \rightarrow \prod_{l \in [r]} [m_l]$$

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 f is a tensor

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

defined by

$$\rho^f [X_0, \dots, X_{d-1}, Y_0, \dots, Y_{p-1}] = \sum_{x_0, \dots, x_{d-1} \in \prod_{k \in [d]} [m_k]} e_{x_0, \dots, x_{d-1}} [X_0, \dots, X_{d-1}] \otimes e_{f(x_0, \dots, x_{d-1})} [Y_0, \dots, Y_{p-1}] .$$

2.6.2 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*

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

the selection encoding of f is a tensor

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

defined by the basis decomposition

$$\gamma^f [X_{[d]}, L_{[n]}] = \sum_{x_0, \dots, x_{d-1} \in \prod_{k \in [d]} [m_k]} e_{x_0, \dots, x_{d-1}} [X_{[d]}] \otimes f(x_0, \dots, x_{d-1}) [L_{[n]}] .$$

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

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

Decomposition and Inference of Factored Representations

The computational automation of reasoning is rooted both in the probabilistic and the logical reasoning tradition. Both draw on the same ontological commitment that systems have a factored structure, that is their states are described by assignments to a set of variables. Based on this commitment both approaches bear a natural tensor representation of their states and a formalism of the respective reasoning algorithms based on multilinear methods.

3 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 normations. 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.

3.1 Tensor Representation of Distributions

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 in Chapter ??.

Definition 16 (Probability Tensor). *Let there be a factored system defined by a categorical variable X_k for each $k \in [d]$ taking values in $[m_k]$. A probability distribution over the states of \mathcal{F} is a tensor*

$$\mathbb{P}[X_0, \dots, X_{d-1}] : \bigtimes_{k \in [d]} [m_k] \rightarrow [0, 1] \subset \mathbb{R}$$

such that

$$\sum_{x_0, \dots, x_{d-1} \in \bigtimes_{k \in [d]} [m_k]} \mathbb{P}[X_0 = x_0, \dots, X_{d-1} = x_{d-1}] = 1.$$

We notice that there are two conditions for a tensor to be probability tensor. First, the tensor needs to have non-negative coordinates and second, the coordinates need to sum to 1.

The probability tensor to the distribution is an object

$$\mathbb{P}[X_0, \dots, X_{d-1}] \in \bigotimes_{k \in [d]} \mathbb{R}^{m_k}$$

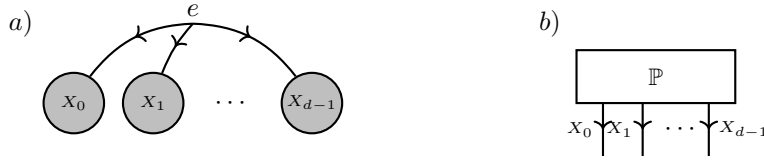
which is the sum over the one-hot encodings (see Lemma 20)

$$\mathbb{P}[X_0, \dots, X_{d-1}] = \sum_{x_0, \dots, x_{d-1} \in \bigtimes_{k \in [d]} [m_k]} \mathbb{P}[X_0 = x_0, \dots, X_{d-1} = x_{d-1}] \cdot e_{x_0, \dots, x_{d-1}}[X_0, \dots, X_{d-1}].$$

The normation condition of probability tensors can be expressed by the contraction equation $1 = \langle \mathbb{P} \rangle [\emptyset]$ since

$$1 = \sum_{x_0, \dots, x_{d-1}} \mathbb{P}[X_0 = x_0, \dots, X_{d-1} = x_{d-1}] = \sum_{x_0, \dots, x_{d-1}} \langle \mathbb{P}, e_{x_0, \dots, x_{d-1}} \rangle [\emptyset] = \langle \mathbb{P} \rangle [\emptyset].$$

Probability tensors are depicted as



3.1.1 Base measures

From a measure theoretic perspective, probabilities are measurable functions called probability densities, which integrals are 1. 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} as a probability density with respect to ν , any probabilistic interpretation will be through the contraction $\langle \mathbb{P}, \nu \rangle [X_{[d]}]$ and the normation condition reads as

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

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 ν , 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 ν , if for all $x_{[d]}$ with $\nu[X_{[d]} = x_{[d]}] = 0$ we have $\mathbb{P}[X_{[d]} = x_{[d]}] = 0$.

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

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

and can therefore safely ignore the base measures.

Starting with Chapter 5 we will further investigate boolean tensors and relate them with propositional formulas. In Chapter 6 we will connect the representation and positivity with respect to boolean base measures with the formalism of entailment.

We now investigate, which base measures ν can be chosen for a probability distribution \mathbb{P} , such that \mathbb{P} is representable by ν . Here we want to find a ν , 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 6 orders based on entailment and show the minimality in Theorem 39. 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} is positive with respect to a boolean base measure ν , if the distribution is representable by ν (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 $\mathbb{P}[X_0 = x_0, \dots, X_{d-1} = x_{d-1}] > 0$.

3.2 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 defined for each x_0 as the tensor

$$\mathbb{P}[X_0] : [m_0] \rightarrow \mathbb{R}$$

defined for $x_0 \in [m_0]$ by

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

Definition 19 generalizes to marginalizations of sets of variables, since we can always group a set of categorical variables and understand them as a single one.

Theorem 1. For any distribution $\mathbb{P}[X_0, X_1]$ the marginal distribution of the variable X is the contraction

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

Further, any marginal distribution is a probability distribution.

Proof. We have $\mathbb{P}[X_0] = \langle \{\mathbb{P}\} \rangle [X_0]$ by definition. To show that $\mathbb{P}[X_0]$ is a probability distribution, we need to show that $\langle \mathbb{P}[X_0] \rangle [\emptyset] = 1$. But this follows from the normation of \mathbb{P} and the commutativity of contractions (see Theorem 91 in Chapter 16) as

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

□

We depict the sum over the possible values of X_1 by contraction of the probability tensor with the trivial tensors \mathbb{I} as

$$\boxed{\mathbb{P}[X_0]}_{x_0 \downarrow} = \boxed{\mathbb{P}[X_0, X_1]}_{x_0 \downarrow, x_1 \downarrow} \boxed{\mathbb{I}}$$

Let us notice, that marginal distributions are probability tensors for themselves, which we again denote by a directed leg.

3.3 Conditional Probabilities

Normations 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} is normable on $\{X_1\}$. Then the distribution of X_0 conditioned on X_1 is defined by*

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

We show in the next theorem, that conditional distributions are calculated by normations.

Theorem 2. *The tensor $\mathbb{P}[X_0 | X_1]$ is the normation of $\mathbb{P}[X_0, X_1]$ on X_1 (see Definition 13), that is*

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

Further, for any $x_1 \in [m_1]$ the tensor $\mathbb{P}[X_0 | X_1 = x_1]$ is a probability tensor.

Proof. The first claim follows from a comparison of Definition 20 and 13. The second claim follows from the first and Theorem 86. Alternatively, the second claim can be showed using the diagrammatic notation as

$$\sum_{x_{X_0}} \mathbb{P}[X = x_{X_0} | Y = x_{X_1}] = \frac{\text{Diagram 1}}{\text{Diagram 2}} = 1.$$

Diagram 1: A box labeled $\mathbb{P}[X|Y]$ with an incoming wire labeled X from a box labeled \mathbb{I} and an outgoing wire labeled Y to a box labeled e_{x_Y} .

Diagram 2: A box labeled $\mathbb{P}[X, Y]$ with an incoming wire labeled X from a box labeled \mathbb{I} and an outgoing wire labeled Y to a box labeled e_{x_Y} .

□

Theorem 1 and 2 show that the formalism of contractions and normations is applied in basic operations of probabilistic reasoning.

We can 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 ingoing and the other outgoing.*

Proof. " \Rightarrow ": By Theorem 2 a conditional probability tensor $\mathbb{P}[X_0 | X_1]$ is the normation of a tensor and by Theorem 86 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 normations.

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

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

is a probability tensor, since

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

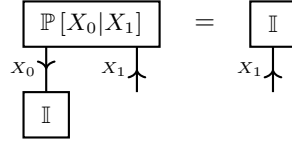
The conditional probability $\mathbb{P}[X_{\mathcal{V}^{\text{out}}} | X_{\mathcal{V}^{\text{in}}}]$ coincides with T , since

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

where in the last equation we used that the denominator is by definition trivial since T is normed.

□

Since conditional probabilities are directed tensors we therefore depict them by



Theorem 3 specifies a broad class of tensors to represent conditional probabilities. In combination with Theorem 95, which states that relational encodings are directed, we get that any relational encoding of a function is a conditional probability tensor.

3.4 Bayes Theorem and the Chain Rule

So far, we have connected concepts of probability theory such as marginal and conditional probabilities with contractions and normations 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. Directly from the more generic contraction equation Theorem 88, since by assumption of positivity of $\mathbb{P}[X_1]$, the tensor network \mathbb{P} is normable with respect to X_1 . \square

Probability distributions can be decomposed into conditional probabilities, as we demonstrate in the next theorem.

Theorem 5 (Chain Rule). *For any joint probability distribution \mathbb{P} of the variables $\mathbb{P}[X_0, \dots, X_{d-1}]$ we have*

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

where for $k = 0$ we denote by $\mathbb{P}[X_0|X_0, \dots, X_{-1}]$ the marginal distribution $\mathbb{P}[X_0]$.

Proof. This follows from Theorem 89. \square

3.5 Independent Variables

Independence leads to severe simplifications of conditional probabilities and is thus the key assumption to gain sparse decompositions. We will demonstrate this here applying the chain rule.

Definition 21 (Independence). *Given a joint distribution of variables X_0 and X_1 , we say that X_0 is independent from X_1 if for any values x_0, x_1 we have*

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

We give a criterion on independence based on a contraction equation of the probability distribution in the next theorem.

Theorem 6. *Given a probability distribution \mathbb{P} , X_0 is independent from X_1 , if and only if*

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

Proof. By Theorem 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

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

is equivalent to

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

We reorder the summations and arrive at

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

which is by Lemma 20 equal to the claim

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

\square

Independent variables result in decompositions of \mathbb{P} in a tensor product of marginal probability tensors. Having pairwise independent variables reduces the degrees of freedom from exponentially many in the number of atoms to linear.

In the tensor network decomposition we depict this by

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

Independence is a very strong assumption, which is often too restrictive. Conditional independence instead is a less demanding assumption, when certain conditional distribution variables are independent. This leads to tensor network decompositions with a more realistic assumption.

Definition 22 (Conditional Independence). *Given a joint distribution of variables X_0, X_1 and X_2 , we say X_0 is independent from X_1 conditioned on X_2 if for any incides x_0, x_1 and x_2*

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

Conditional independence is a relation between conditional probabilities and is therefore equivalent to a normation equation stated next.

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 from X_1 if and only if the contraction 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. Directly by Theorem 2 used on the conditional probabilities in Definition 22. □

We can exploit conditional independence to find tensor network decompositions of probability tensors, as we show in the next theorem.

Corollary 1. *If and only if X_0 is independent from X_1 conditioned on X_2 the probability distribution \mathbb{P} satisfies*

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

Proof. Follows from Theorem 7 and Theorem 4. □

Corollary 2. *Whenever X_0 is independent of X_1 given X_2 , we have*

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

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

Figure 4: Diagrammatic visualization of the contraction equation in Corollary 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.

Theorem 8 (Markov Chain). *Let there be a set of variables X_t where $t \in [T]$. When X_t is independent of $X_{0:t-2}$ conditioned on X_{t-1} (the Markov Property), then*

$$\mathbb{P} = \langle \mathbb{P}[X_t | X_{t-1}] : t \in [T] \rangle [X_0, \dots, X_{T-1}]$$

We depict this decomposition in Figure 5.

Proof. By the chain rule (Theorem 5) we have

$$\mathbb{P}[X_0, \dots, X_{T-1}] = \langle \{\mathbb{P}[X_t | X_{0:t}] : t \in [T]\} \rangle [X_{[T]}]$$

Using the conditional independence of X_t and $X_{0:t-2}$ conditioned on X_{t-1} we further have by Corollary 2

$$\mathbb{P}[X_t | X_{0:t} = x_{0:t}] = \mathbb{P}[X_t | X_{t-1} = x_{t-1}] .$$

Composing both equalities shows the claim. \square

Here we denoted by $X_{0:t}$ the tuple X_0, \dots, X_t .

Remark 3. Let us notice that the dimensionality dropped drastically through applying the independence assumption. The tensor space in the naive representation of any probability distribution has

$$\prod_{t \in [T]} m_t$$

coordinates, while the Markov Chain is represented by

$$\sum_{t \in [T]} m_t \cdot m_{t-1} .$$

Replacing exponential scaling with the number of variables to linear scaling is the advantage of tensor network decompositions.

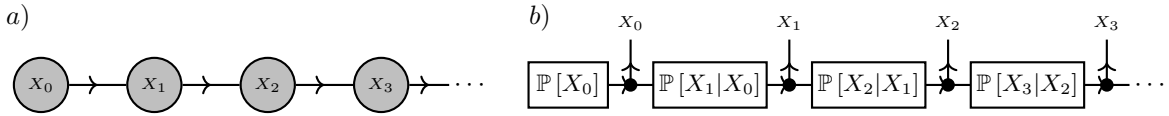


Figure 5: Depiction of a Markov Chain. a) Dependency Graph (of the corresponding chain Graphical Model). b) Dual Tensor Network representing the conditional probability factors.

3.6 Graphical Models

We have already depicted conditional dependency assumptions made for Markov Chains in Figure 5 and discussed the implied decomposition of the dual tensor networks. Graphical models provide a more general framework for conditional dependency assumptions and provide a generic approach to exploit independences in finding tensor network decompositions of \mathbb{P} .

Following the tensor network formalism we in this section introduce graphical models based on hypergraphs. Whether the hypergraph is directed or not distinguished between Bayesian Networks and Markov Networks.

3.6.1 Markov Networks

While typically Markov Networks are defined on graphs, we define them here on hypergraphs to establish a direct connection to tensor networks defined on the same hypergraph. Along that line, Markov Networks are tensor networks with non-negative tensors (see Definition 7), which are interpreted as probability distributions after normation.

Definition 23 (Markov Network). Let $\mathcal{T}^{\mathcal{G}}$ be a tensor network of non-negative tensors on a hypergraph \mathcal{G} . Then the Markov Network to $\mathcal{T}^{\mathcal{G}}$ is the probability distribution of X_v defined by the tensor

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

We call the denominator

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

the partition function of the Markov Network.

Often, we are only interested in the distribution of a subset of variables, which are called the observable variables, and call the other variables hidden variables. The marginalization of a Markov Network to \mathcal{T}^G on the variables $X_{\tilde{\mathcal{V}}}$ is

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

This can be derived from Theorem 91, 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

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

Definition 24 (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 an alternative graph-based definition of Markov Networks Koller and Friedman (2009).

Definition 25 (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 Theorem 10.

Theorem 9. *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. Let there be a hypergraph \mathcal{G} , a Markov Network \mathcal{T}^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

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

By Theorem 7, it now follows that X_A is independent of X_B conditioned on X_C . 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, which we will proof later in Section 14.3. \square

From the proof of Theorem 9 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)).

3.6.2 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 Definition 6 by introducing the property of acyclicity for hypergraphs.

Definition 26. *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 27 (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 Theorem 87 that their contraction defines a probability distribution.

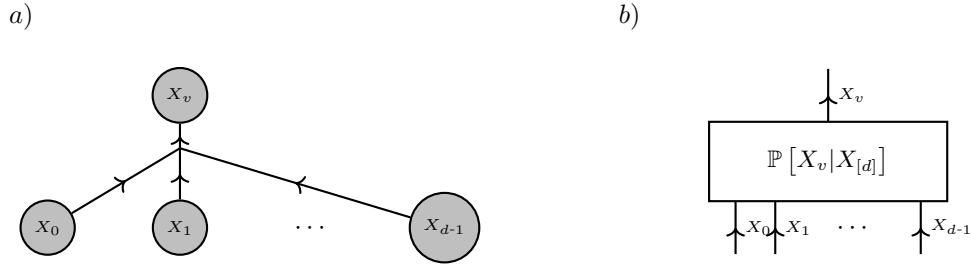


Figure 6: Example of a Factor of a Bayesian Network to the node X_v with parents X_0, \dots, X_{d-1} , as subgraph a) and dual tensor core b).

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 10 (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 Corollary 2 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

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

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 Theorem 7, it now follows that X_v is independent of $\text{NonDes}(v)$ conditioned on $\text{Pa}(v)$. \square

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

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

In particular, $\mathcal{T}^{\mathcal{G}}$ 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 $T^{(\emptyset, \{v\})}$ is then normed with no incoming variables and we thus have

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

and

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

$|\mathcal{V}| - 1 \rightarrow |\mathcal{V}|$: 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 $\mathcal{T}^{\tilde{\mathcal{G}}}$ the tensor network on the hypergraph $\tilde{\mathcal{G}} = \{\mathcal{V}/\{v\}, \mathcal{E}/\{(\text{Pa}(v), \{v\})\}\}$ with decorations inherited from $\mathcal{T}^{\mathcal{G}}$. With Theorem 91, the directionality of $T^{(\text{Pa}(v), \{v\})}$ and the induction assumption on $\mathcal{T}^{\tilde{\mathcal{G}}}$ we have

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

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

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

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

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

\square

Theorem 11 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 normation 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.

3.6.4 Example: Hidden Markov Models

We here extend the example of Markov Chains from Theorem 8 to a limited observation of the variables by observations. Let there be the variables X_t (states) and E_t (observations) with a discrete and finite time $t \in [T]$.

The conditional assumptions are

- X_{t+1} is independent of $X_{0:t-1}$ and $E_{0:t}$ conditioned on X_t
- E_t is independent of all other variables conditioned on X_t

Then the probability tensor has the decomposition

$$\mathbb{P}[X_{0:T}, E_{0:T}] = \prod_{t \in [T]} (\mathbb{P}[X_t | X_{0:t-1}, E_{0:t-1}] \cdot \mathbb{P}[E_t | X_{0:t}, E_{0:t-1}]) \quad (6)$$

$$= \mathbb{P}[X_0] \cdot \mathbb{P}[E_0 | X_0] \cdot \prod_{t \in [T], t > 0} (\mathbb{P}[X_t | X_{t-1}] \cdot \mathbb{P}[E_t | X_t]) \quad (7)$$

Here we used the Chain Rule decomposition of Theorem 5 in the first equation and the conditional independence assumptions in the second.

We notice, that this is a Bayesian Network on a directed acyclic hypergraph \mathcal{G} consistent in nodes to each state and each observation and directed hyperedges

- $(\{X_t\}, \{X_{t+1}\})$ for $t \in [T - 1]$
- $(\{X_t\}, \{E_t\})$ for $t \in [T]$

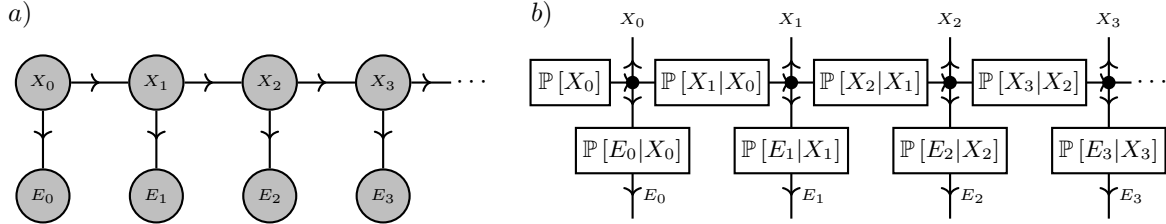


Figure 7: Depiction of a Hidden Markov Model. a) Dependency Graph (of the corresponding chain Graphical Model). b) Dual Tensor Network representing the conditional probability factors.

3.7 Exponential Families

Exponential families are collections of probability distributions, where each coordinate is determined by a base measure and a set ϕ of features as

$$\mathbb{P}[X_{[d]} = x_{[d]}] \propto \nu(x) \cdot \exp \left[\sum_{l \in [p]} \phi_l [X_{[d]} = x_{[d]}] \cdot \theta [L = l] \right].$$

We use the selection encoding to represent the weighted summation over the statistics, that is the tensor

$$\gamma^\phi [X_{[d]}, L] : \bigtimes_{k \in [d]} [m_k] \times [p] \rightarrow \mathbb{R}$$

with

$$\gamma^\phi [X_{[d]} = x_{[d]}, L = l] = \phi_l [X_{[d]} = x_{[d]}].$$

We then understand θ as a vector to the categorical variable L and use Theorem 85 to get

$$\sum_{l \in [p]} \theta [L = l] \cdot \phi_l [X_{[d]}] = \langle \gamma^\phi [X_{[d]}, L], \theta [L] \rangle [X_{[d]}].$$

Definition 28. Given a sufficient statistics

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

and a boolean base measure

$$\nu : \bigtimes_{k \in [d]} [m_k] \rightarrow [2]$$

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

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

is called the exponential family to ϕ . We further define for each member with parameters θ the associated energy tensor

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

and the cumulant function

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

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

As before, we restrict for boolean base measures, which satisfy $\langle \nu \rangle [\emptyset] \neq 0$. We notice, that by positivity of the exponential function, any distribution in an exponential family $\Gamma^{\phi, \nu}$ is positive with respect to ν (see Definition 18). In Chapter 8 we will investigate distributions, where the base measures and the sufficient statistics share a common computation framework.

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

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

Proof. By definition we have

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

□

Definition 29 (Minimal). We say that a statistic ϕ is minimal with respect to a boolean base measure ν , if there is no pair of a nonvanishing vector $V[L]$ and a scalar $\lambda \in \mathbb{R}$ with

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

3.7.1 Tensor Network Representation

We can use the relational encoding formalism to represent members of exponential families by a single contraction, as we show next. The central insight here is a relational encoding of the sufficient statistics, which enables representation by tensor network decomposition, when the sufficient statistic is decomposable.

Theorem 12 (Generic Representation of Exponential Families). Given any base measure ν and a sufficient statistic ϕ we enumerate for each coordinate $l \in [p]$ the image $\text{im}(\phi_l)$ by a variable X_{ϕ_l} taking values in $[\text{im}(\phi_l)]$, given an interpretation map

$$I^l : [\text{im}(\phi_l)] \rightarrow \text{im}(\phi_l) .$$

For any parameter vector $\theta[L] : [p] \rightarrow \mathbb{R}$ we build the activation cores

$$W^l [X_{\phi_l} = x_{\phi_l}] = \exp [\theta[L = l] \cdot I^l(x_{\phi_l})]$$

and have

$$\mathbb{P}^{(\phi, \theta, \nu)} = \langle \{\nu, \rho^\phi\} \cup \{W^l : l \in [p]\} \rangle [X_{[d]} | \emptyset] .$$

Proof. We use an extended image of ϕ by

$$\text{im}(\phi) = \bigtimes_{l \in [p]} \text{im}(\phi_l) .$$

Theorem 98 implies

$$\exp[\langle \gamma^\phi, \theta \rangle [X_{[d]}]] = \langle \rho^\phi, \exp[\langle \cdot, \theta \rangle]_{\text{im}(\phi)} \rangle [X_{[d]}] .$$

The claim follows from the equation

$$\exp[\langle \cdot, \theta \rangle]_{\text{im}(\phi)} [X_{\phi_0}, \dots, X_{\phi_{p-1}}] = \bigotimes_{l \in [p]} \exp[\cdot \theta [L = l]]_{\text{im}(\phi_l)} [X_{\phi_l}] = \bigotimes_{l \in [p]} W^l [X_{\phi_l}] .$$

□

We notice, that the relational encoding is the contraction of the relational encoding of its coordinate maps as

$$\rho^\phi [X_{[d]}, X_{\phi_{[p]}}] = \langle \rho^{\phi_0}, \dots, \rho^{\phi_{p-1}} \rangle [X_{[d]}, X_{\phi_{[p]}}] .$$

We will show this property in Theorem 99. One strategy to create ρ^ϕ is thus the creation of the encoding of all its coordinate maps. When the coordinate maps are sharing common components, a sparser representation can be derived through encodings of the components shared among the coordinate map encodings.

A tensor network representation of an exponential family is thus a Markov Network consistent of two types of cores. Computation cores are relational encodings of statistics ρ^{ϕ_l} . Our intuition is that they compute the hidden variable X_{ϕ_l} , based on Basis Calculus (see Chapter 15). Activation cores W^l exploit the computed variable and provide, when contracted with the relational encoding, a factor

$$\langle \rho^{\phi_l}, W^l [X_l] \rangle [X_{[d]}]$$

to the Markov Network reduced to the visible coordinates $X_{[d]}$. The activation cores are trivial, i.e. $W^l [X_{\phi_l}] = \mathbb{I} [X_{\phi_l}]$, when $\theta [L = l] = 0$. In that case

$$\langle \rho^{\phi_l}, W^l [X_l] \rangle [X_{[d]}] = \mathbb{I} [X_{[d]}]$$

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

Remark 4 (Comparison of relation and selection encodings). *Relation encodings are in general of higher dimensions than selection encodings. In can thus be intractible to instantiate the probability distribution as a tensor networks, while the energy tensor can still be efficiently represented based on selection encodings. In this case, energy-based reasoning algorithms are tractible while more direct methods are intractible.*

3.7.2 Mean Parameters

Mean parameters are an alternative way to represent members of exponential families.

Definition 30. *Let there be an exponential family defined by a statistic ϕ and a boolean base measure. We call the tensor*

$$\mu = \langle \mathbb{P}, \gamma^\phi \rangle [L]$$

the mean parameter tensor to a distribution \mathbb{P} of an exponential family. The set

$$\mathcal{M}_{\phi, \nu} = \{ \langle \mathbb{P}, \gamma^\phi, \nu \rangle [L] : 0 \prec \mathbb{P}, \langle \mathbb{P}, \nu \rangle [X_{[d]}] = \nu [X_{[d]}] \} ,$$

where Γ denotes the set of all probability distributions, is called the convex polytope of realizable mean parameters. The map

$$F^{(\phi, \nu)} : \mathbb{R}^p \rightarrow \mathbb{R}^p$$

with $F^{(\phi, \nu)}(\theta) = \langle \mathbb{P}^{(\phi, \theta, \nu)}, \gamma^\phi \rangle [L]$ is called the forward map of the exponential family and any map

$$B^{(\phi, \nu)} : \text{im} \left(F^{(\phi, \nu)} \right) \rightarrow \mathbb{R}^p$$

with $\mathbb{P}^{(\phi, B^{(\phi, \nu)}(F^{(\phi, \nu)}(\theta)), \nu)} = \mathbb{P}^{(\phi, \theta, \nu)}$ for any $\theta \in \mathbb{R}^p$ a backward map.

While introduced here as a property of a distribution, the mean parameters will be central to the discussion of probabilistic inference in Chapter 4. We now provide a simple characterization of the sets of mean parameters based on slices of the selection encoding of the statistic.

Theorem 13. For any statistic ϕ the polytope of mean parameters is the convex hull of the slices of γ^ϕ with fixed indices to $X_{[d]}$, that is

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

Proof. This follows from the fact, that the set of probability distributions is the convex hull of the one-hot encodings and the convex hull of mean parameters is a linear transform of that set. \square

The convex polytope $\mathcal{M}_{\phi,\nu}$ can further be characterized as an intersection of half-spaces, as we state next.

Theorem 14. For any statistic ϕ and base measure ν 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

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

Proof. This is a standard result of combinatorial optimization, see e.g. Ziegler (2013). \square

The determination of the the vectors $((a_i [L], b_i) : i \in [n])$ and is one reason for the intractability of probabilistic inference (see e.g. Wainwright and Jordan (2008)).

Definition 31. Given a mean parameter polytope $\mathcal{M}_{\phi,\nu}$ in the half space representation of Theorem 14, and any subset $\mathcal{I} \subset [n]$ we say that the set

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

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

Theorem 15. For any non-empty face $Q_{\phi,\nu}^{\mathcal{I}}$ to a subset $\mathcal{I} \subset [n]$ there is a vector $\theta [L]$, which we call a normal of the face, such that

$$Q_{\phi,\nu}^{\mathcal{I}} = \text{argmax}_{\mu \in \mathcal{M}_{\phi,\nu}} \langle \theta [L], \mu [L] \rangle [\emptyset].$$

For any collection of positive λ_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 $Q_{\phi,\nu}^{\mathcal{I}}$.

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 λ_i . We use that the face is non-empty, and thus there is a $\mu [L]$ with

$$\langle \mu [L], a_i [L] \rangle [\emptyset] = b_i$$

for all $i \in \mathcal{I}$. Since for any $\mu \in \mathcal{M}_{\phi,\nu}$

$$\langle \mu [L], a_i [L] \rangle [\emptyset] \leq b_i$$

it follows that

$$\max_{\mu \in \mathcal{M}_{\phi,\nu}} \langle \theta [L], \mu [L] \rangle [\emptyset] = \sum_{i \in \mathcal{I}} \lambda_i \cdot b_i.$$

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 $\mu \in Q_{\phi,\nu}^{\mathcal{I}}$. \square

In a slide abuse of notation, we denote in this case $Q_{\phi,\nu}^\theta = Q_{\phi,\nu}^{\mathcal{I}}$.

3.7.3 Examples

Example 4 (The minterm exponential family). *When taking as sufficient statistic the identity $\delta [X_{[d]}, L_{[d]}]$, we can represent any positive distribution \mathbb{P} as a member of the exponential family, namely when choosing the canonical parameter*

$$\theta = \ln [\mathbb{P}] .$$

The associated mean parameter is then, after relabeling the variables $X_{[d]}$ by $L_{[d]}$,

$$\mu = \mathbb{P}$$

and $\mathcal{M}_{\delta, \mathbb{I}}$ coincides with the set of probability tensors. For reasons to be explained in the Chapter 5 we refer to this family as the minterm exponential family.

Given a hypergraph with fixed node decoration, the different decorations of the hyperedges by tensors can be represented by an exponential family, as we show next.

Theorem 16 (Exponential Representation of Markov Networks). *For any hypergraph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ we define a sufficient statistics*

$$\phi = \bigotimes_{e \in \mathcal{E}} \phi_e$$

where

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

Given any Markov Network $\{T^e : e \in \mathcal{E}\}$ on \mathcal{G} with positive tensors T^e we define

$$\theta = \bigotimes_{e \in \mathcal{E}} \theta_e$$

where

$$\theta_e = \ln [T^e]$$

and \ln acts coordinatewise. Then, the Markov Network is in the member of the exponential family with trivial base measure, sufficient statistic ϕ and parameters θ .

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

$$\prod_{e \in \mathcal{E}} T^e [X_e = x_e] = \exp \left[\sum_{e \in \mathcal{E}} \theta_e [X_e = x_e] \right] = \exp \left[\sum_{e \in \mathcal{E}} \langle \theta_e [X_e], \phi_e(x_{\mathcal{V}}) \rangle [\emptyset] \right] . \quad (8)$$

Using that

$$\langle \phi, \theta \rangle [X_{\mathcal{V}}] = \sum_{e \in \mathcal{E}} \langle \phi_e, \theta_e \rangle [X_{\mathcal{V}}]$$

we get

$$\langle \{T^e : e \in \mathcal{E}\} \rangle [X_{\mathcal{V}}] = \exp [\langle \theta, \phi \rangle [X_{\mathcal{V}}]] . \quad (9)$$

This implies

$$\langle \{T^e : e \in \mathcal{E}\} \rangle [X_{\mathcal{V}} | \emptyset] = \langle \exp [\langle \theta, \phi \rangle [X_{\mathcal{V}}]] \rangle [X_{\mathcal{V}} | \emptyset] . \quad (10)$$

□

The mean parameter of the Markov Network exponential family is the cartesian product of the marginals $\mu_e [X_e]$ are often referred to as beliefs in the literature. **They are outer bounded by local consistency polytope, which leads to the motivation of message passing algorithms!**

3.8 Empirical Distributions

Let us now apply the formalism of probability distributions in tensor network representations to encode data.

Definition 32. Given a dataset $\left((x_0^j, \dots, x_{d-1}^j) : j \in [m] \right)$ 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

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

where by ρ^D we denote the relational encoding (see Definition 14) of the sample selector map, and the distributed variables $X_{[d]}$ are the head variables of the relational encoding.

The relational encoding of the sample selector map is the sum

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

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

$$\mathbb{P}^D [X_{[d]} = x_{[d]}] = \frac{1}{\langle \rho^D \rangle [\emptyset]} \left(\sum_{j \in [m]} e_{x_0^j, \dots, x_{d-1}^j} [X_{[d]} = x_{[d]}] \right) = \frac{|j \in [m] : (x_0^j, \dots, x_{d-1}^j) = (x_0, \dots, x_{d-1})|}{|j \in [m]|}$$

and is thus interpreted as the frequency of the corresponding world in the data.

The relational 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 Section 17.1.2 as basis CP decompositions. We now exploit this structure to find efficient tensor network decompositions (see Figure 8) based on matrices encoding its variables.

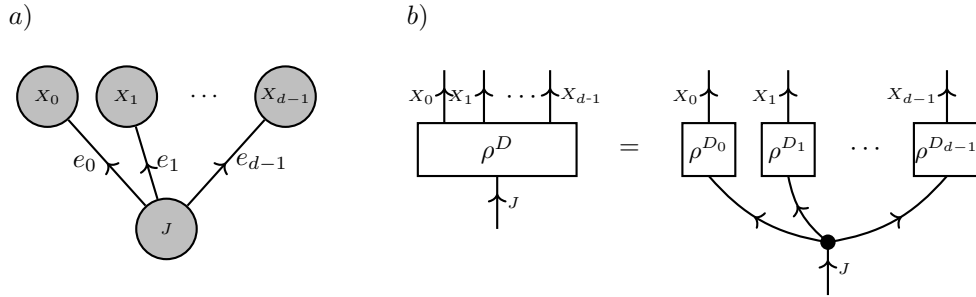


Figure 8: Representation of a dataset $\left((x_0^j, \dots, x_{d-1}^j) : j \in [m] \right)$. a) Interpretation as a data selection variable J selecting states for the variables $X_{[d]}$. b) Corresponding decomposition of the relational encoding ρ^D into a tensor network in the basis CP Format (see Section 17.1.2), where $T^{e_k} = \rho^{D_k}$.

Theorem 17. 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

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

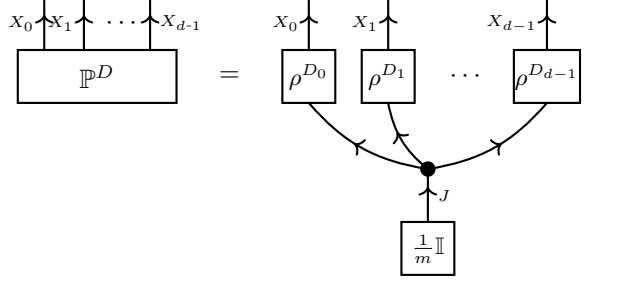
We then have

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

and

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

In a contraction diagram this reads



Proof. The first claim is a special case of Theorem 99, to be shown in Chapter 15. To show the second claim we notice

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

With the first claim it now follows that

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

□

The cores ρ^{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 17 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 ρ^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 8, the variable J served as single parent for each categorical variable X_k).

3.9 Discussion and Outlook

Remark 5 (Alternative definitions of graphical models). *In the literature, tensor networks are often called dual to the hypergraphs defining graphical models (see e.g. Robeva and Seigal (2019)). The duality becomes clear, when one interpretes the tensors as cores and their common variables as edges. We in this work avoid this ambiguity by directly defining tensor networks as decoration of hyperedges by tensors.*

Often, the tensors decorating hyperedges are called factors and their logarithm features Koller and Friedman (2009).

Further, we directly use hypergraphs instead of the more canonical association of factors with cliques of a graph. This avoids the discussion of non-maximal cliques as decorated with trivial tensors. Such hypergraphs follow the same line of thought compared with factor graphs, which are bipartite graphs with nodes either corresponding with single variables or with a collection of them affected by a factor.

4 Probabilistic Reasoning

We have investigated means to store the knowledge about a system and now turn to the retrieval of information, a process called inference.

Contraction of the relational encoding of a function with a Markov Network gives the statistics over the values of the functions. When contracting the function directly, we get the expectation.

One can increase the efficiency of inference algorithms by using approximative contractions. Here, message passing schemes can be applied as to be introduced in Chapter 16.

4.1 Queries

In the previous chapter, we have derived efficient representation schemes of probability distributions based on tensor network decompositions. We have argued that one should avoid naive instantiation of these distributions based on an storage of each coordinates. In the task of reasoning, we want to retrieve information encoded in the probability distribution. To derive an efficient approach one therefore needs to avoid instantiating the distribution in a coordiantewise manner in an intermediate step. We thus formalize a basic reasoning scheme by contractions of the decomposed distributions with query tensors.

4.1.1 Querying by functions

We can formalize queries by retrieving expectations of functions given a distribution specified by probability tensors. We exploit basis calculus in defining categorical variables X_f to tensors f , which are enumerating the set $\text{im}(f)$. More details on this scheme are provided in Chapter 15, see Definition 76 therein.

Definition 33. *The marginal query of a probability distribution $\mathbb{P}[X_{[d]}]$ by a tensor*

$$f : \bigtimes_{k \in [d]} [m_k] \rightarrow \mathbb{R}$$

is the vector $\mathbb{P}[X_f] \in \mathbb{R}^{|\text{im}(f)|}$ defined as the contraction

$$\mathbb{P}[X_f] = \langle \mathbb{P}[X_{[d]}], \rho^f[X_{[d]}, X_f] \rangle [X_f] .$$

The expectation query of \mathbb{P} by f is

$$\mathbb{E}[f] = \langle f, \mathbb{P} \rangle [\emptyset] .$$

Given another tensor $g : \bigtimes_{k \in [d]} [m_k] \rightarrow \mathbb{R}$ the conditional query of the probability distribution $\mathbb{P}[X_{[d]}]$ by the tensor f conditioned on the tensor g is the matrix $\mathbb{P}[X_f|X_g] \in \mathbb{R}^{|\text{im}(f)|} \otimes \mathbb{R}^{|\text{im}(g)|}$ defined as the normation

$$\mathbb{P}[X_f|X_g] = \langle \{ \mathbb{P}[X_{[d]}], \rho^f[X_{[d]}, X_f], \rho^g[X_{[d]}, X_g] \} \rangle [X_f|X_g] .$$

Expectation queries are contractions of marginal queries with identities, that is

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

This will be shown in more detail in Chapter 15 in Corollary 12.

Conditional probabilities are queries, where the tensors f and g are identity mappings in the respective variable state spaces. Conversely, we can understand the conditional query $\mathbb{P}[f|g]$ as the conditional probability of f conditioned on g , of the underlying Markov Network with cores $\{\mathbb{P}, \rho^f, \rho^g\}$ and variables X_f, X_g besides the variables distributed by \mathbb{P} .

We further denote event queries by

$$\mathbb{E}[f = z] = \langle \mathbb{P}, \rho^f, e_z \rangle [\emptyset]$$

where by e_z be denote the one hot encoding of the state z with respect to some enumeration. Let us note that they are further contraction of the queries in Definition 33 since by Theorem 91

$$\begin{aligned} \mathbb{E}[f = z] &= \langle \langle \mathbb{P}, \rho^f \rangle [X_f], e_z \rangle [\emptyset] \\ &= \langle \mathbb{P}[f], e_z \rangle [\emptyset] . \end{aligned}$$

4.1.2 MAP Queries

Find the maximal variable of a tensor is a problem, which can be approached by sampling methods as we discuss here.

Definition 34. *Given a tensor T the MAP query is the problem*

$$\text{argmax}_{x_0, \dots, x_{d-1}} T[X_0 = x_0, \dots, X_{d-1} = x_{d-1}] . \quad (11)$$

By coordinate calculus, we notice that

$$T[X_0 = x_0, \dots, X_{d-1} = x_{d-1}] \langle T, e_{x_0, \dots, x_{d-1}} \rangle [\emptyset] . \quad (12)$$

Given the image Γ^{EL} of one-hot encodings, the MAP query problem is equivalent to

$$\max_{x_0, \dots, x_{d-1}} T[X_0 = x_0, \dots, X_{d-1} = x_{d-1}] = \max_{\theta \in \Gamma^{\text{EL}}} \langle T, \theta \rangle [\emptyset] . \quad (13)$$

We can thus understand MAP queries as a Tensor Network approximation problem, where the approximating tensor are the one-hot encodings of states.

Remark 6 (MAP queries on energy and probability tensors). *Since the exponential function is monotonic, MAP queries on the energy tensor of an exponential family with uniform base measure are equivalent to MAP queries of their energies.*

4.1.3 Answering queries by energy contractions

Let us now interpret a probability tensor at hand as a member of an exponential family (see Section 3.7), which is always possible when taking the naive exponential family.

Lemma 2. *For any probability distribution \mathbb{P} with $\mathbb{P} = \langle \exp [E [X_{[d]}]] \rangle [X_{[d]}|\emptyset]$, disjoint subsets $A, B \subset [d]$ with $A \cup B = [d]$ and any x_B we have*

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

Proof. Since no summation is commuted. □

Thus, it suffices to build the selection encoding of the statistics, and we can avoid the usage of the relational encoding.

We notice, that Lemma 2 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.

Lemma 3. *For any probability distribution \mathbb{P} with $\mathbb{P} = \langle \exp [E [X_{[d]}]] \rangle [X_{[d]}|\emptyset]$, disjoint subsets $A, B \subset [d]$ and any x_B we have*

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

Proof. By splitting the contraction into terms to $A \cup B$. □

4.2 Sampling based on queries

Let us here investigate how to draw samples from distributions \mathbb{P} , based on queries on \mathbb{P} .

Since there are $\prod_{v \in \mathcal{V}} m_v$ coordinates stored in \mathbb{P} , naive methods are often infeasible. One can instead exploit a representation of \mathbb{P} by a Markov network or the energy term in an exponential family for efficient algorithms and sample from local proxy distributions resulting from contractions and interpreted as marginal and conditional probabilities.

4.2.1 Exact Methods

Forward Sampling (see Algorithm 1) uses a chain decomposition (see Theorem 5) of a probability distribution to iteratively sample the variables.

Algorithm 1 Forward Sampling

for $k \in [d]$ **do**

 Draw $x_k \in [m_k]$ from the conditional query

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

end for

Forward Sampling is especially efficient, when sampling from a Bayesian Network respecting the topological order of its nodes. The reason for this lies in trivializations of all conditional distributions, which heads are not included in the evidence of previously sampled variables. More technically, we can show that

$$\mathbb{P} [X_k | X_{\tilde{k}} = x_{\tilde{k}} : \tilde{k} < k] = \mathbb{P} [X_k | X_{\text{Pa}(k)} = x_{\text{Pa}(k)}] ,$$

which is only involving a single core of a Bayesian network. **This can be shown using Corollary 13 to be derived in Chapter 15.**

4.2.2 Approximate Methods

When there are many variables to be sample, the computation of the conditional probability to all variables can be infeasible. One way to overcome this is Gibbs Sampling: Iteratively resemble single variables given the rest as evidence.

Sample each variable independent from the marginal distribution. Then, alternate through the variables and sample each variable from the conditional distribution taking the others as evidence.

Algorithm 2 Gibbs Sampling

```

for  $k \in [d]$  do
  Draw State for atom  $k$  from initialization distributions.
end for
while Stopping criterion is not met do
  for  $k \in [d]$  do
    Draw  $x_k \in [m_k]$  from the conditional query

```

$$\mathbb{P} \left[X_k | X_{\tilde{k}} = x_{\tilde{k}} : \tilde{k} \neq k \right]$$

```

  end for
end while

```

Gibbs can be implemented based on the energy tensor E of the probability tensor, as follows form the Lemma 2.

This is in contrast with forward sampling, where we need to sum over many coordinates of the exponentiated energy tensor, which amounts to the representation of the probability distribution as a tensor network using relational encodings.

4.2.3 Simulated Annealing

MAP queries are approximated by sampling from annealed distributions: Use T as the energy tensor, e.g. as parameter tensor to the naive exponential family.

Here by the naive exponential family! Simulated annealing manipulates the probability used to sample x_k in terms of an inverse temperature parameter β , by

$$\mathbb{P} \rightarrow \frac{\exp [\beta \cdot \ln [\mathbb{P}]]}{\langle \exp [\beta \cdot \ln [\mathbb{P}]] \rangle [\emptyset]} .$$

When the temperature is larger than 1, the probability of states with low probability increases while the probability of states with large probability decreases and for low temperatures the opposite. Simulated annealing, that is the decrease of the temperature to 0 during Gibbs sampling biases the algorithm towards states with large probability.

For any exponential family the transformation

$$E \rightarrow \beta \cdot E$$

can be performed by rescaling the canonical parameters as

$$\theta \rightarrow \beta \cdot \theta .$$

4.3 Maximum Likelihood Estimation

Let us now turn to inductive reasoning tasks, where a probabilistic model is trained on given data.

4.3.1 Likelihood and Loss

Given a datapoint $D(j)$ consisting of the images of the data selecting map D (see Definition 32), the likelihood given a Markov Logic Network is denoted as

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

When all $D(j)$ are drawn independently from $\mathbb{P} [X_{[d]}]$, we can factorize into

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

It is convenient to apply a logarithm on the objective, which does not influence the optimum when optimizing this quantity. This is especially useful, when investigating the convergence of the objective for $m \rightarrow \infty$ (see Chapter 10).

Definition 35. We define the loss of a distribution \mathbb{P} as

$$\mathcal{L}_D(\mathbb{P}) = \frac{1}{m} \ln [\mathbb{P} [\{D(j)\}_{j \in [m]}]]$$

We now state the Maximum Likelihood Problem in the form

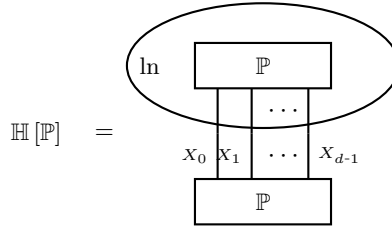
$$\operatorname{argmin}_{\mathbb{P} \in \Gamma} \mathcal{L}_D(\mathbb{P}) . \quad (\mathbb{P}_{\Gamma, \mathbb{P}^D}^{\mathcal{L}_D})$$

4.3.2 Entropic Interpretation

Definition 36 (Shannon entropy). *The information content or the Shannon entropy of a distribution is defined as*

$$\mathbb{H}[\mathbb{P}] := \mathbb{E}_{X_{[d]} \sim \mathbb{P}} [-\ln [\mathbb{P} [X_{[d]}]]] = \langle \mathbb{P}, -\ln [\mathbb{P}] \rangle [\emptyset] .$$

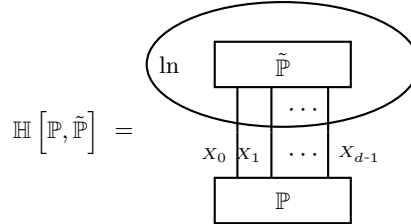
We depict this in a tensor network diagram with an ellipsis denoting a coordinatewise transform (see Chapter ??) with a natural logarithm \ln as:



Definition 37 (Cross entropy). *The cross entropy between two distributions is defined as*

$$\mathbb{H}[\mathbb{P}, \tilde{\mathbb{P}}] := \mathbb{E}_{X_{[d]} \sim \mathbb{P}} [-\ln [\tilde{\mathbb{P}}[X_{[d]}]]] = \langle \mathbb{P}, -\ln [\tilde{\mathbb{P}}] \rangle [\emptyset] .$$

We depict this in a tensor network diagram with an ellipsis denoting a coordinatewise transform (here the \ln) as :



We here use $\ln[0] = -\infty$ and $0 \cdot \ln[0] = 0$. Then we have $\mathbb{H}[\mathbb{P}, \tilde{\mathbb{P}}] = \infty$ if and only if there is a $x_{[d]}$ such that $\mathbb{P}[X_{[d]} = x_{[d]}] > 0$ and $\tilde{\mathbb{P}}[X_{[d]} = x_{[d]}] = 0$.

The Gibbs inequality states that

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

The difference between both sides is called the Kullback Leibler Divergence and a useful metric in reasoning, since it vanishes for $\mathbb{P} = \tilde{\mathbb{P}}$.

Definition 38 (Kullback Leibler Divergence). *The KL divergence between two distributions is defined as*

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

We are now ready to provide an entropic interpretation of the loss introduced in Definition 35.

Theorem 18. *Given a data selecting map D and a distribution \mathbb{P} we have*

$$\mathcal{L}_D(\mathbb{P}) = \mathbb{H}[\mathbb{P}^D, \mathbb{P}] . \quad (14)$$

Proof. We have

$$\begin{aligned}\mathcal{L}_D(\mathbb{P}) &= \frac{1}{m} \ln [\mathbb{P}[\{D(j)\}_{j \in [m]}]] = \frac{1}{m} \sum_{j \in [m]} \ln [\mathbb{P}[X_{[d]} = D(j)]] = \frac{1}{m} \sum_{j \in [m]} \langle \{\ln [\mathbb{P}], e_{D(j)}\} \rangle [\emptyset] \\ &= \langle \mathbb{P}^D, \ln [\mathbb{P}] \rangle [\emptyset] .\end{aligned}$$

Comparing with the negative log likelihood we notice that that loss coincides with the cross-entropy between the empirical distribution \mathbb{P}^D and \mathbb{P} , i.e.

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

□

We can therefore rewrite Problem $\mathcal{P}_{\Gamma, \mathbb{P}^D}^{\mathcal{L}_D}$ as minimization of cross-entropies and of Kullback Leibler divergences as

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

Most general, the Maximum Likelihood Problem is the M-Projection of a distribution \mathbb{P}^* onto a set Γ of probability tensors is

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

where the Maximum Likelihood Estimation is the special case $\mathbb{P}^* = \mathbb{P}^D$.

Example 5 (Cross entropy with respect to exponential families). *If $\tilde{\mathbb{P}}$ from an exponential family with boolean base measure, have with the representation from Lemma 1*

$$\mathbb{H}[\mathbb{P}, \mathbb{P}^{(\phi, \theta, \nu)}] = \langle \mathbb{P}, \ln [\mathbb{P}^{(\phi, \theta, \nu)}] \rangle [\emptyset] = \langle \mathbb{P}, \gamma^\phi \rangle [\emptyset] - A^{(\phi, \nu)}(\theta) + \langle \mathbb{P}, \ln [\nu] \rangle [\emptyset] .$$

For the trivial base measure we can further exploit the existence of the energy tensor and have the representation

$$\mathbb{H}[\mathbb{P}, \mathbb{P}^{(\phi, \theta, \nu)}] = \langle \mathbb{P}, (E^{(\phi, \theta, \nu)}[X_{[d]}] - A^{(\phi, \nu)}(\theta) \cdot \mathbb{I}) \rangle [\emptyset] = \langle \mathbb{P}, E^{(\phi, \theta, \nu)}[X_{[d]}] \rangle [\emptyset] - A^{(\phi, \nu)}(\theta) .$$

4.4 Forward Mapping in Exponential Families

Mean parameter coordinates are expectation queries to ϕ_l , by

$$\mu[L = l] = \mathbb{E}[\phi_l] .$$

Forward mappings have a closed form representation by

$$F^{(\phi, \nu)}(\theta) = \langle \gamma^\phi, \langle \nu, \exp[\langle \gamma^\phi, \theta \rangle [\emptyset]] \rangle [X_{[d]} | \emptyset] \rangle [L] .$$

This contraction can, however, be infeasible, since it requires the instantiation of the probability tensor, 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.

4.4.1 Variational Formulation

Besides the direct computation of the mean parameter tensor we can give a variational characterization of the forward mapping. This is especially useful, when the contraction is intractable, for example because the tensor $\mathbb{P}^{(\phi, \theta, \nu)}$ is infeasible to create.

Theorem 19. *We have*

$$F^{(\phi, \nu)}(\theta) = \operatorname{argmax}_{\mu \in \mathcal{M}_{\phi, \nu}} \langle \mu, \theta \rangle [\emptyset] + \mathbb{H}[\mathbb{P}^\mu]$$

where by \mathbb{P}^μ we denote a probability distribution with respect to a base measure ν , which reproduces the mean parameter μ .

Proof. Theorem 3.4 in Wainwright and Jordan (2008). □

Let us now characterize the image of the forward map, which turns out to be the interior of the mean polytope, if the statistic is minimal (see Definition 29).

Theorem 20. *For any statistics ϕ , which is minimal with respect to a base measure ν , the image $\text{im} (F^{(\phi, \nu)})$ of the forward map is the interior of the convex polytope $\mathcal{M}_{\phi, \nu}$.*

Proof. Theorem 3.3 in Wainwright and Jordan (2008). □

For the practice usage of this theorem, we need a characterization of the interior of $\mathcal{M}_{\phi, \nu}$.

Theorem 21. *For any minimal statistics ϕ and boolean base measure ν we have for some $\mu [L]$ that $\mu [L] \in \mathcal{M}_{\phi, \nu}$ if and only if there is a positive distribution with respect to ν such that*

$$\mu [L] = \langle \mathbb{P}, \gamma^\phi \rangle [L] .$$

Proof. " \Rightarrow ": By Theorem 20 we find a canonical parameter $\theta [L]$ such that

$$\mu [L] = \langle \mathbb{P}^{(\phi, \theta, \nu)} [X_{[d]}], \gamma^\phi [X_{[d]}, L] \rangle [L] .$$

We notice, that $\mathbb{P}^{(\phi, \theta, \nu)}$ is positive with respect to ν , as is any member of an exponential family with base measure ν .

" \Leftarrow ": Since by assumption the statistics is minimal, the convex set $\mathcal{M}_{\phi, \nu}$ 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 $\mu \in \mathcal{M}_{\phi, \nu}^\circ$ if for any non-vanishing vector $V [L]$ 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

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

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

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

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

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

and thus

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

We can interpret this quantity as the variance of the random variable $\langle \gamma^\phi [X_{[d]}, L], V [L] \rangle [X_{[d]} = x_{[d]}]$, 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 ϕ with respect to ν it follows, that this variable is not constant and we therefore have

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

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

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

□

4.4.2 Boundary of convex polytopes

For mean parameters $\mu [L]$ outside the interior of $\mathcal{M}_{\phi, \nu}$ we know by Theorem 21, that any distribution with mean parameter $\mu [L]$ is not positive with respect to ν 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 22. *Let there be a minimal ϕ with respect to the base measure ν and $\mu [L] \notin \mathcal{M}_{\phi, \nu}^\circ$. Then there is a $\theta [L]$ with*

$$\mu [L] \in \operatorname{argmax}_{\mu \in \mathcal{M}_{\phi, \nu}} \langle \theta [L], \mu [L] \rangle [\emptyset]$$

and all distributions with 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

$$\mathcal{U} = \operatorname{argmax}_{x_{[d]}} \langle \theta, \phi(x_{[d]}) \rangle [\emptyset] .$$

Proof. When $\mu \notin \mathcal{M}_{\phi, \nu}^\circ$ we find a face such that $\mu \in Q_{\phi, \nu}^T$. The existence of $\theta [L]$ follows from Theorem 15, in which also a construction procedure is provided given a half-space representation (see Theorem 14).

Now, we have

$$\mu [L] \in \operatorname{argmax}_{\mu \in \mathcal{M}_{\phi, \nu}} \langle \theta [L], \mu [L] \rangle [\emptyset]$$

and thus

$$\mu [L] \in \operatorname{conv} \left(\gamma^\phi X_{[d]} = x_{[d]}, L : x_{[d]} \in \operatorname{argmax}_{x_{[d]} : \nu[X_{[d]}=x_{[d]}]=1} \langle \theta [L], \gamma^\phi X_{[d]} = x_{[d]}, L \rangle [\emptyset] \right)$$

Thus, any distribution reproducing meanparam is a convex combination of the one-hot encodings of the states in $\operatorname{argmax}_{x_{[d]}} \langle \theta [L], \gamma^\phi X_{[d]} = x_{[d]}, L \rangle [\emptyset]$, and therefore representable with respect to the base measure $\tilde{\nu}$. \square

Each face of $\mathcal{M}_{\phi, \nu}$ thus defines a refinement of a base measure, which is sufficient to reproduce the mean parameters on that face.

Definition 39. *The base measure to the face of \mathcal{M} with normal θ is*

$$\nu^{\phi, \theta} = \mathbb{I}_{\operatorname{argmax}_{x_{[d]}} \langle \theta, \phi(x_{[d]}) \rangle [\emptyset]} [X_{[d]}] .$$

Theorem 22 therefore states, that when a mean parameter is on a face of $\mathcal{M}_{\phi, \nu}$, then each distribution reproducing the mean parameter has a representation with respect to the refined base measure

$$\tilde{\nu} [X_{[d]}] = \langle \nu, \nu^{\phi, \theta} \rangle [X_{[d]}] .$$

We now utilize these findings and provide in Algorithm 4.4.2 a procedure to refine the base measure until the reduced mean parameter is in the open set of a reduced mean parameter polytope.

Theorem 23. *For arbitrary inputs ν, ϕ and $\mu \in \mathcal{M}_{\phi, \nu}$, Algorithm 4.4.2 terminates in finite time and outputs a triple of base measure $\tilde{\nu}$, statistic $\tilde{\phi}$ and mean parameter $\tilde{\mu}$ such that the following holds. Any probability tensor \mathbb{P} reproducing μ is representable with respect to $\tilde{\nu}$ and $\tilde{\mu} \in (\mathcal{M}_{\tilde{\phi}, \tilde{\nu}})^\circ$.*

Proof. Let us first show, that Algorithm 4.4.2 always terminates. The inner while loop of Algorithm 4.4.2 always terminates, since ϕ 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

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

and thus $Q_{\phi, \nu}^\theta = \mathcal{M}_{\phi, \nu}$, which is a contradiction with the assumption of a non-trivial face.

The second claim follows from an iterative application of Theorem 22 and the fact, that a probability distribution reproduces μ in a non-minimal representation, if and only if it reproduces the corresponding reduced μ with respect to the reduced statistics. \square

Algorithm 3 Base Measure Refinement**Input:** Base measure ν , statistic ϕ and mean parameter $\mu \in \mathcal{M}_{\phi, \nu}$ **Output:** Refined base measure $\tilde{\nu}$, remaining statistic $\tilde{\phi}$ and remaining mean parameter $\tilde{\mu}$ **while** $\mu \notin (\mathcal{M}_{\phi, \nu})^\circ$ **do** **while** ϕ not minimal with respect to ν (see Definition 29) **do** Find non-vanishing vector $V[L]$ and scalar $\lambda \in \mathbb{R}$ such that

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

 Choose a coordinate $l \in [p]$ with $V[L = l] \neq 0$ and drop it from ϕ and μ **end while** Find a non-trivial face (i.e. a non-empty face, which is a proper subset of $\mathcal{M}_{\phi, \nu}$) with normal θ , such that

$$\mu \in Q_{\phi, \nu}^\theta$$

Refine base measure

$$\nu \leftarrow \langle \nu, \nu^{\phi, \theta} \rangle [X_{[d]}]$$

end while**return** ν, ϕ, μ

Example 6 (Faces with normals parallel to one-hot encodings). *To get some intuition how to represent face base measures, let us consider face normals $\theta \in \{\lambda \cdot e_l [L] : l \in [p], \lambda \in \mathbb{R}/\{0\}\}$. We use relational encodings of the coordinates ϕ_l of the statistic ϕ , with head variables X_{ϕ_l} with dimension m_{ϕ_l} enumerating the image $\text{im}(\phi_l) \subset \mathbb{R}$ in an ascending order. If $\theta [L] = \lambda \cdot e_l [L]$ with $\lambda > 0$, then $\text{argmax}_{x_{[d]}} \langle \theta, \phi(x_{[d]}) \rangle [\emptyset]$ consists of states $x_{[d]}$ with minimal statistic $\phi_l [X_{[d]} = x_{[d]}]$, that is*

$$\nu^{\phi, \lambda \cdot e_l} [X_{[d]}] = \left\langle \rho^{\phi_l} [X_{[d]}, X_{\phi_l}], e_{m_{\phi_l}-1} [X_{\phi_l}] \right\rangle [X_{[d]}] .$$

If $\theta [L] = \lambda \cdot e_l [L]$ with $\lambda < 0$, then at the states with minimal statistic $\phi_l [X_{[d]} = x_{[d]}]$, that is

$$\nu^{\phi, \lambda \cdot e_l} [X_{[d]}] = \left\langle \rho^{\phi_l} [X_{[d]}, X_{\phi_l}], e_0 [X_{\phi_l}] \right\rangle [X_{[d]}] .$$

4.4.3 Mode Search by annealing

Finding the mode of a distribution is related to the forward mapping of $\beta \cdot \theta$: μ to a delta distribution (or in the convex hull of multiple maxima) in the limit.

This is because

$$\text{argmax}_{\mu \in \mathcal{M}_{\phi, \nu}} \langle \mu, \theta \rangle [\emptyset]$$

is taken at an extreme point in $\mathcal{M}_{\phi, \nu}$ (since linear objective over closed convex set), which is a delta distribution of a set and

$$\text{argmax}_{\mu \in \mathcal{M}_{\phi, \nu}} \langle \mu, \beta \cdot \theta \rangle [\emptyset] + \mathbb{H} [\mathbb{P}^\mu] = \text{argmax}_{\mu \in \mathcal{M}_{\phi, \nu}} \langle \mu, \theta \rangle [\emptyset] + \frac{1}{\beta} \cdot \mathbb{H} [\mathbb{P}^\mu]$$

thus the entropy term is neglectible for large β . A more precise argument is using a limit of the maxima and can be found in Theorem 8.1 in Wainwright and Jordan (2008)

4.4.4 Mean Field Method

We rewrite

$$\max_{\mu \in \mathcal{M}_{\phi, \nu}} \langle \mu, \theta \rangle [\emptyset] + \mathbb{H} [\mathbb{P}^\mu] = \max_{\mathbb{P}} \langle E, \mathbb{P} \rangle [\emptyset] + \mathbb{H} [\mathbb{P}]$$

where

$$E = \langle \gamma^\phi, \theta \rangle [X_{[d]}] .$$

We now restrict the distributions in the maximum. Typically we use the family of independent distributions, also called naive mean field method. The naive mean field is the approximation by distributions of independent random variables V^k , that is

$$\text{argmax}_{V^k : k \in [d]} \langle \{E\} \cup \{V^k : k \in [d]\} \rangle [\emptyset] + \sum_{k \in [d]} \mathbb{H} [V^k] .$$

Theorem 24 (Update equations for the mean field approximation). *Keeping all legs but one constant, the problem*

$$\operatorname{argmax}_{V^k} \left\langle \{E\} \cup \{V^k : k \in [d]\} \right\rangle [\emptyset] + \sum_{k \in [d]} \mathbb{H}[V^k]$$

is solved at

$$V^k[X_k] = \left\langle \exp \left[\left\langle \{E[X_{[d]}]\} \cup \{V^{\tilde{k}}[X_{\tilde{k}}] : \tilde{k} \neq k\} \right\rangle [X_{[d]}] \right] \right\rangle [X_k | \emptyset] .$$

Proof. We have

$$\frac{\partial \mathbb{H}[V^k]}{\partial V^k} = -\ln[V^k[X_k]] + \mathbb{I}[X_k]$$

and by multilinearity of tensor contractions

$$\frac{\partial \left\langle \{E\} \cup \{V^{\tilde{k}} : \tilde{k} \in [d]\} \right\rangle [\emptyset]}{\partial V^k} = \left\langle \{E\} \cup \{V^{\tilde{k}} : \tilde{k} \in [d], \tilde{k} \neq k\} \right\rangle [X_k] .$$

Combining both, the condition

$$0 = \frac{\partial \left(\left\langle \{E\} \cup \{V^{\tilde{k}} : \tilde{k} \in [d]\} \right\rangle [\emptyset] + \sum_{k \in [d]} \mathbb{H}[V^k] \right)}{\partial V^k}$$

is equal to

$$\ln[V^k[X_k]] = \mathbb{I}[X_k] + \left\langle \{E\} \cup \{V^{\tilde{k}} : \tilde{k} \in [d], \tilde{k} \neq k\} \right\rangle [X_k] .$$

Together with the condition $\langle V^k \rangle [=] 1$ this is satisfied at

$$V^k[X_k] = \left\langle \exp \left[\left\langle \{E\} \cup \{V^{\tilde{k}} : \tilde{k} \neq k\} \right\rangle [X_k] \right] \right\rangle [X_k | \emptyset] .$$

□

Algorithm 4 is the alternation of legwise updates until a stopping criterion is met.

Algorithm 4 Naive Mean Field Approximation

for $k \in [d]$ **do**

$$V^k[X_k] \leftarrow \langle \mathbb{I} \rangle [X_k | \emptyset]$$

end for

while Stopping criterion is not met **do**

for $k \in [d]$ **do**

$$V^k[X_k] \leftarrow \left\langle \exp \left[\left\langle \{E[X_{[d]}]\} \cup \{V^{\tilde{k}}[X_{\tilde{k}}] : \tilde{k} \neq k\} \right\rangle [X_k] \right] \right\rangle [X_k | \emptyset]$$

end for

end while

4.4.5 Structured Variational Approximation

More generically, we restrict the maximum over the mean parameters of efficiently contractable distributions and get a lower bound. In this section we use any Markov Network as the approximating family.

Let \mathcal{G} be any hypergraph, we define the problem

$$\operatorname{argmax}_{\mathbb{P} \in \Gamma^{\mathcal{G}, \mathbb{I}}} \langle E, \mathbb{P} \rangle [\emptyset] + \mathbb{H}[\mathbb{P}] \quad (\mathbb{P}_{\Gamma^{\mathcal{G}, \mathbb{I}}, \mathbb{P}})$$

We approximate the solution of this problem again by an alternating algorithm, which iteratively updates the cores of the approximating Markov Network.

Theorem 25 (Update equations for the structured variational approximation). *The Markov Network $\mathcal{T}^{\mathcal{G}}$ with hypercores $\{T^e : e \in \mathcal{E}\}$ is a stationary point for Problem $P_{\Gamma^{\mathcal{G}}, \mathbb{I}, \mathbb{P}}$, if for all $e \in \mathcal{E}$*

$$T^e[X_e] = \lambda \cdot \exp \left[\frac{\langle \{E\} \cup \{T^{\tilde{e}} : \tilde{e} \neq e\} \rangle [X_e]}{\langle \{T^{\tilde{e}} : \tilde{e} \neq e\} \rangle [X_e]} - \sum_{\tilde{e} \neq e} \frac{\langle \{\ln [T^{\tilde{e}}]\} \cup \{T^{\tilde{e}} : \tilde{e} \neq \hat{e}\} \rangle [X_e]}{\langle \{T^{\tilde{e}} : \tilde{e} \neq \hat{e}\} \rangle [X_e]} \right]$$

for any $\lambda > 0$ (e.g. by the norm). Here, the quotient denotes the coordinatewise quotient.

Proof. We proof the theorem by first order condition on the objective $O(\mathcal{T}^{\mathcal{G}}) = \langle E, \langle \mathcal{T}^{\mathcal{G}} \rangle [X_{\mathcal{V}} | \emptyset] \rangle [\emptyset] + \mathbb{H}[\langle \mathcal{T}^{\mathcal{G}} \rangle [X_{\mathcal{V}} | \emptyset]]$.

To proof the theorem, we use Lemma 22, which shows a characterization of the derivative of functions

We have

$$\langle E, \langle \mathcal{T}^{\mathcal{G}} \rangle [X_{[d]} | \emptyset] \rangle [\emptyset] = \frac{\langle \{E\} \cup \mathcal{T}^{\mathcal{G}} \rangle [\emptyset]}{\langle \mathcal{T}^{\mathcal{G}} \rangle [\emptyset]}.$$

Further we have

$$\mathbb{H}[\langle \mathcal{T}^{\mathcal{G}} \rangle [X_{[d]} | \emptyset]] = \left(\sum_{\tilde{e} \in \mathcal{E}} \langle -\ln [T^{\tilde{e}}], \langle \mathcal{T}^{\mathcal{G}} \rangle [X_{[d]} | \emptyset] \rangle [\emptyset] \right) + \ln [\langle \mathcal{T}^{\mathcal{G}} \rangle [\emptyset]]$$

We define the tensor

$$\tilde{T}[X_{\mathcal{V}}] = E[X_{\mathcal{V}}] - \sum_{\tilde{e} \neq e} \ln [T^{\tilde{e}}[X_{\tilde{e}}]] \otimes \mathbb{I}[X_{\mathcal{V}/\tilde{e}}]$$

and notice, that \tilde{T} does not depend on T^e .

The objective has then a representation as

$$O(\mathcal{T}^{\mathcal{G}}) = \langle \tilde{T}[X_{\mathcal{V}}], \langle \mathcal{T}^{\mathcal{G}} \rangle [X_{\mathcal{V}} | \emptyset] \rangle [\emptyset] - \langle \ln [T^e], \langle \mathcal{T}^{\mathcal{G}} \rangle [X_{\mathcal{V}} | \emptyset] \rangle [\emptyset] + \ln [\langle \mathcal{T}^{\mathcal{G}} \rangle [\emptyset]]$$

Let us now differentiate all terms. With Lemma 22 we now get

$$\begin{aligned} \frac{\partial}{\partial T^e[Y_e]} \langle \tilde{T}[X_{\mathcal{V}}], \langle \mathcal{T}^{\mathcal{G}} \rangle [X_{\mathcal{V}} | \emptyset] \rangle [\emptyset] &= \left\langle \tilde{T}[X_{\mathcal{V}}], \delta[Y_e, X_e], \frac{\langle \mathcal{T}^{\mathcal{G}} \rangle [X_e]}{T^e[X_e]}, \langle \mathcal{T}^{\mathcal{G}} \rangle [X_{\mathcal{V}/e} | X_e] \right\rangle [Y_e, X_{\mathcal{V}}] \\ &\quad - \left\langle \tilde{T}[X_{\mathcal{V}}], \langle \mathcal{T}^{\mathcal{G}} \rangle [X_{\mathcal{V}} | \emptyset] \right\rangle [\emptyset] \otimes \left\langle \frac{\langle \mathcal{T}^{\mathcal{G}} \rangle [Y_e]}{T^e[Y_e]} \right\rangle [Y_e]. \end{aligned}$$

Further we have

$$\begin{aligned} \frac{\partial}{\partial T^e[Y_e]} \langle \ln [T^e], \langle \mathcal{T}^{\mathcal{G}} \rangle [X_{\mathcal{V}} | \emptyset] \rangle [\emptyset] &= \left\langle \ln [T^e[X_e]], \delta[Y_e, X_e], \frac{\langle \mathcal{T}^{\mathcal{G}} \rangle [X_e]}{T^e[X_e]}, \langle \mathcal{T}^{\mathcal{G}} \rangle [X_{\mathcal{V}/e} | X_e] \right\rangle [Y_e, X_{\mathcal{V}}] \\ &\quad - \langle \ln [T^e[X_e]], \langle \mathcal{T}^{\mathcal{G}} \rangle [X_{\mathcal{V}} | \emptyset] \rangle [\emptyset] \otimes \left\langle \frac{\langle \mathcal{T}^{\mathcal{G}} \rangle [Y_e]}{T^e[Y_e]} \right\rangle [Y_e] \\ &\quad - \left\langle \frac{1}{T^e[X_e]}, \langle \mathcal{T}^{\mathcal{G}} \rangle [X_{\mathcal{V}} | \emptyset] \right\rangle [\emptyset] \end{aligned}$$

and (see Proof of 21)

$$\frac{\partial}{\partial T^e[Y_e]} \ln [\langle \mathcal{T}^{\mathcal{G}} \rangle [\emptyset]] = \frac{\frac{\partial}{\partial T^e[Y_e]} \langle \mathcal{T}^{\mathcal{G}} \rangle [\emptyset]}{\langle \mathcal{T}^{\mathcal{G}} \rangle [\emptyset]} = \frac{\langle \mathcal{T}^{\mathcal{G}} \rangle [Y_e]}{T^e[Y_e]}.$$

Together, the first order condition

$$0 = \frac{\partial}{\partial T^e[Y_e]} O(\mathcal{T}^{\mathcal{G}})$$

is equal to all y_e satisfying

$$\begin{aligned} 0 = & \frac{\langle \mathcal{T}^{\mathcal{G}} \rangle [Y_e = y_e]}{T^e [Y_e = y_e]} \left(\left\langle \tilde{T} [X_{\mathcal{V}/e}, X_e = y_e], \langle \mathcal{T}^{\mathcal{G}} \rangle [X_{\mathcal{V}/e} | X_e = y_e] \right\rangle [\emptyset] \right. \\ & - \left\langle \tilde{T} [X_{\mathcal{V}}], \langle \mathcal{T}^{\mathcal{G}} \rangle [X_{\mathcal{V}} | \emptyset] \right\rangle [\emptyset] \\ & - \langle \ln [T^e [X_e = y_e]], \langle \mathcal{T}^{\mathcal{G}} \rangle [X_{\mathcal{V}/e} | X_e = y_e] \rangle [\emptyset] \\ & \left. + \langle \ln [T^e [X_e]], \langle \mathcal{T}^{\mathcal{G}} \rangle [X_{\mathcal{V}} | \emptyset] \rangle [\emptyset] \right). \end{aligned}$$

We notice, that by normation

$$\langle \ln [T^e [X_e = y_e]], \langle \mathcal{T}^{\mathcal{G}} \rangle [X_{\mathcal{V}/e} | X_e = y_e] \rangle [\emptyset] = \ln [T^e [X_e = y_e]]$$

and that the scalar

$$\lambda_1 = \left\langle \tilde{T} [X_{\mathcal{V}}], \langle \mathcal{T}^{\mathcal{G}} \rangle [X_{\mathcal{V}} | \emptyset] \right\rangle [\emptyset] - \langle \ln [T^e [X_e]], \langle \mathcal{T}^{\mathcal{G}} \rangle [X_{\mathcal{V}} | \emptyset] \rangle [\emptyset]$$

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

$$\ln [T^e [X_e = y_e]] = \left\langle \tilde{T} [X_{\mathcal{V}/e}, X_e = y_e], \langle \mathcal{T}^{\mathcal{G}} \rangle [X_{\mathcal{V}/e} | X_e = y_e] \right\rangle [\emptyset] + \lambda_1.$$

The claim follows when applying the exponential on both sides and with the observation, that

$$\left\langle \tilde{T} [X_{\mathcal{V}/e}, X_e = y_e], \langle \mathcal{T}^{\mathcal{G}} \rangle [X_{\mathcal{V}/e} | X_e = y_e] \right\rangle [\emptyset] = \frac{\left\langle \{\tilde{T}\} \cup \{T^{\tilde{e}} : \tilde{e} \neq e\} \right\rangle [X_e = y_e]}{\langle \{T^{\tilde{e}} : \tilde{e} \neq e\} \rangle [X_e = y_e]}$$

and reparametrization of λ_1 to

$$\lambda = \exp [\lambda_1].$$

□

The mean field method corresponds with minimization of the KL Divergence to the efficiently contractable family, i.e. the I-projection onto the family.

Theorem 26. *For any hypergraph \mathcal{G} and energy tensor E we have*

$$\operatorname{argmax}_{\mathbb{P} \in \Gamma^{\mathcal{G}, \mathbb{I}}} \langle E, \mathbb{P} \rangle [\emptyset] + \mathbb{H} [\mathbb{P}] = \operatorname{argmax}_{\mathbb{P} \in \Gamma^{\mathcal{G}, \mathbb{I}}} \mathrm{D}_{\mathrm{KL}} \left[\mathbb{P}^{(\mathcal{G}, \theta)} || \langle \exp [E] \rangle [X_{[d]} | \emptyset] \right]$$

Problem $\mathbb{P}_{\Gamma^{\mathcal{G}, \mathbb{I}}, \mathbb{P}}$ is thus the I-projection onto the exponential family $\Gamma^{\mathcal{G}, \mathbb{I}}$.

Proof. By rearranging the objective to the KL divergence.

□

4.5 Backward Mapping in Exponential Families

The parameters optimizing the likelihood, will be shown to coincide with the backward mapping evaluated on the expectation of the sufficient statistics (see Theorem ??). This is in most generality true for the parameters of the M-projection of any distribution onto the exponential family. We therefore investigate methods to compute the backward mapping, in most generality by alternating algorithms and in the special case of Markov Logic Networks by closed form representations.

We have that θ is a solution of the backward problem at μ^* , if and only if

$$\left\langle \mathbb{P}^{(\phi, \theta, \nu)}, \gamma^{\phi} \right\rangle [L] = \mu^* [L].$$

This contraction equation is called moment matching, since the moment of the empirical distribution is matched by the moment of the fitting distribution.

We find one backward mapping as the dual problem to the forward mapping.

4.5.1 Variational Formulation

The backward mapping to $\mu_D [L] = \langle \mathbb{P}^D, \gamma^\phi \rangle [L]$ is Maximum Likelihood estimation and the solution of the maximum entropy problem.

Theorem 27. *Let there be a sufficient statistic ϕ . The map $B^{(\phi, \nu)} : \mathbb{R}^p \rightarrow \mathbb{R}^p$ defined as*

$$B^{(\phi, \nu)}(\mu) = \operatorname{argmax}_{\theta \in \mathbb{R}^p} \langle \mu, \theta \rangle [\emptyset] - A^{(\phi, \nu)}(\theta) .$$

is a backward mapping.

Proof. We show the claim can be shown by the first order condition on the objective. It holds that

$$\begin{aligned} \frac{\partial}{\partial \theta [L]} A^{(\phi, \nu)}(\theta) &= \frac{\partial}{\partial \theta [L]} \ln [\langle \exp [\langle \gamma^\phi, \theta \rangle [X_{[d]}]] \rangle [\emptyset]] \\ &= \frac{\partial}{\partial \theta [L]} \frac{\langle \gamma^\phi [L], \exp [\langle \gamma^\phi, \theta \rangle [X_{[d]}]] \rangle [\emptyset]}{\langle \exp [\langle \gamma^\phi, \theta \rangle [X_{[d]}]] \rangle [\emptyset]} \\ &= F^{(\phi, \nu)}(\theta) [L] \end{aligned}$$

and thus

$$\frac{\partial}{\partial \theta [L]} \left(\langle \mu, \theta \rangle [\emptyset] - A^{(\phi, \nu)}(\theta) \right) = \mu [L] - F^{(\phi, \nu)}(\theta) [L] .$$

The first order condition is therefore

$$\mu [L] = F^{(\phi, \nu)}(\theta) [L]$$

and any θ satisfies this condition exactly when $\theta = B^{(\phi, \nu)}(\mu)$ for a backward map. \square

4.5.2 Interpretation by Maximum Likelihood Estimation

Backward mapping coincides with the Maximum Likelihood Estimation Problem ($\mathcal{P}_{\Gamma, \mathbb{P}^D}^{\mathcal{L}_D}$), when we take Γ to the distributions in an exponential family $\Gamma^{\phi, \nu}$ for a sufficient statistic ϕ .

The loss is the cross entropy between a distribution with μ and the distribution $\mathbb{P}^{(\phi, \theta, \nu)}$.

Theorem 28. *Let there be any exponential family, a mean parameter vector $\mu^* \in \operatorname{im}(F^{(\phi, \nu)})$ and a backward map $B^{(\phi, \nu)}$. Then $\hat{\theta} = B^{(\phi, \nu)}(\mu^*)$ is the parameter of the M-projection (Problem $\mathcal{P}_{\Gamma, \mathbb{P}^*}$) of any \mathbb{P}^* with $\langle \gamma^\phi, \mathbb{P}^* \rangle [L] = \mu^* [L]$ on to $\Gamma^{\phi, \nu}$, that is*

$$\mathbb{P}^{(\phi, \hat{\theta}, \nu)} \in \operatorname{argmax}_{\mathbb{P} \in \Gamma^{\phi, \nu}} \mathbb{H} [\mathbb{P}^*, \mathbb{P}] .$$

In particular, if $\mu = \mu_D$ for a data map D , the backward map is a maximum likelihood estimator.

Proof. We exploit the variational characterization of the backward map by Theorem 27, 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 θ we have with Example 5

$$\mathbb{H} [\mathbb{P}^*, \mathbb{P}^{(\phi, \theta, \nu)}] = \langle \mathbb{P}^*, \gamma^\phi, \theta \rangle [\emptyset] - A^{(\phi, \nu)}(\theta) .$$

We use that by assumption $\langle \mathbb{P}^*, \gamma^\phi \rangle [L] = \mu^* [L]$ and thus

$$\mathbb{H} [\mathbb{P}^*, \mathbb{P}^{(\phi, \theta, \nu)}] = \langle \mu^*, \theta \rangle [\emptyset] - A^{(\phi, \nu)}(\theta) .$$

This shows, that the backward map coincides with the M-projection onto $\Gamma = \Gamma^{\phi, \nu}$.

Further, if $\mu = \mu_D$ for a data map D , we have that the corresponding empirical distribution \mathbb{P}^D satisfies $\langle \gamma^\phi, \mathbb{P}^D \rangle [L] = \mu [L]$. The backward map of μ is therefore the M-projection of \mathbb{P}^D , which is with Theorem 18 the maximum likelihood estimator. \square

4.5.3 Connection with Maximum Entropy

The Maximum entropy problem with respect to matching expected statistics $\mu^* \in \mathcal{M}_{\phi, \nu}$

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

where the optimization is over all the distributions Γ^ν , which are representable with respect to the base measure ν .

Theorem 29. *Let ϕ be a statistic and ν a base measure. For any $\mu^* \in (\mathcal{M}_{\phi, \nu})^\circ$ the solution of Problem $\mathbb{P}_{\phi, \nu, \mu^*}^{\mathbb{H}}$ is the distribution $\mathbb{P}^{(\tilde{\phi}, \hat{\theta}, \tilde{\nu})}$, where $\hat{\theta} = B^{\tilde{\phi}, \tilde{\nu}}(\tilde{\mu})$.*

Proof. Since $\mu^* \in (\mathcal{M}_{\phi, \nu})^\circ$, Theorem 21 implies the existence of $\hat{\theta}$ such that

$$\mu^*[L] = \langle \mathbb{P}^{(\phi, \hat{\theta}, \nu)}, \gamma^\phi \rangle [L] .$$

We now follow the argumentation of the proof of Theorem 20.2 in Koller and Friedman (2009). Let $\tilde{\mathbb{P}}$ further be an arbitrary distribution, possibly different from $\mathbb{P}^{(\phi, \hat{\theta}, \nu)}$, such that

$$\mu^*[L] = \langle \tilde{\mathbb{P}}, \gamma^\phi \rangle [L] .$$

We then have

$$\mathbb{H}[\mathbb{P}^{(\phi, \hat{\theta}, \nu)}] = \mathbb{H}[\tilde{\mathbb{P}}, \mathbb{P}^{(\phi, \hat{\theta}, \nu)}]$$

With the Gibbs inequality we have if $\tilde{\mathbb{P}} \neq \mathbb{P}^{(\phi, \hat{\theta}, \nu)}$

$$\mathbb{H}[\mathbb{P}^{(\phi, \hat{\theta}, \nu)}] - \mathbb{H}[\tilde{\mathbb{P}}] = \mathbb{H}[\tilde{\mathbb{P}}, \mathbb{P}^{(\phi, \hat{\theta}, \nu)}] - \mathbb{H}[\tilde{\mathbb{P}}] > 0 .$$

Therefore, if $\tilde{\mathbb{P}}$ does not coincide with $\mathbb{P}^{(\phi, \hat{\theta}, \nu)}$, it is not a solution of Problem $\mathbb{P}_{\phi, \nu, \mu^*}^{\mathbb{H}}$. \square

Let us highlight the fact, that in Problem $\mathbb{P}_{\phi, \nu, \mu^*}^{\mathbb{H}}$ 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. Theorem 29 states, that when the maximum entropy problem has a solution (i.e. $\mu^* \in \mathcal{M}_{\phi, \nu}$), then the solution is in the exponential family to the statistic ϕ .

When $\mu^* \notin (\mathcal{M}_{\phi, \nu})^\circ$, the mean parameter is by Theorem 21 not reproducible by a member of the exponential family $\Gamma^{\phi, \nu}$. Instead, in combination with the base measure refinement Algorithm 4.4.2, we show that the solution is in a refined exponential family.

Theorem 30. *Let ϕ be a statistic and ν a base measure. For any $\mu^* \in \mathcal{M}_{\phi, \nu}$, let $\tilde{\phi}, \tilde{\nu}$ and $\tilde{\mu}$ be the outputs of Algorithm 4.4.2 when passing ϕ, ν and μ^* as input. Then, the distribution $\mathbb{P}^{(\tilde{\phi}, \hat{\theta}, \tilde{\nu})}$, where $\hat{\theta} = B^{\tilde{\phi}, \tilde{\nu}}(\tilde{\mu})$, solves Problem $\mathbb{P}_{\phi, \nu, \mu^*}^{\mathbb{H}}$.*

Proof. Theorem 23 and the above Lemma. \square

Theorem 30 further implies, that the base measure $\tilde{\nu}$ identified by Algorithm 4.4.2 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 by Theorem 23.

4.5.4 Alternating Algorithms to Approximate the Backward Map

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.

Alternate through the coordinates of the statistics and adjust $\theta [L = l]$ to a minimum of the likelihood, i.e. where for any $l \in [p]$

$$0 = \frac{\partial}{\partial \theta [L = l]} \mathcal{L}_D \left(\mathbb{P}^{(\phi, \theta, \nu)} \right) .$$

This condition is equal to the collection of moment matching equations

$$\langle \mathbb{P}^{(\phi, \theta, \nu)}, \gamma^\phi \rangle [L = l] = \langle \mathbb{P}^D, \gamma^\phi \rangle [\emptyset] L = l.$$

Lemma 4. For any sufficient statistic ϕ a parameter vector θ and a $l \in [p]$ we define

$$T[X_{\phi_l}] = \langle \{\rho^\phi\} \cup \{W^{\tilde{l}} : \tilde{l} \in [p], \tilde{l} \neq l\} \rangle [X_{\phi_l}].$$

Then the moment matching condition for ϕ_l relative to θ and μ is satisfied for any $\theta [L = l]$ with

$$\langle W^l, \text{Id}_{\text{im}(\phi_l)}, T[L_\phi] \rangle [\emptyset] = \langle W^l, T[L_\phi] \rangle [\emptyset] \cdot \mu[L = l].$$

Proof. We have

$$\mathbb{P}^{(\phi, \theta, \nu)} = \frac{\langle W^l, T \rangle [X_{[d]}]}{\langle W^l, T \rangle [\emptyset]}$$

and

$$\langle \mathbb{P}^{(\phi, \theta, \nu)}, \phi_l \rangle [\emptyset] = \frac{\langle W^l, \text{Id}_{\text{im}(\phi_l)}, T \rangle [X_{[d]}]}{\langle W^l, T \rangle [\emptyset]}.$$

Here we used

$$\phi_l = \langle W^l, \text{Id}_{\text{im}(\phi_l)} \rangle [X_{[d]}]$$

and redundancies of copies of relational encodings. It follows that

$$\langle \mathbb{P}^{(\phi, \theta, \nu)}, \phi_l \rangle [\emptyset] = \langle \mathbb{P}^D, \phi_l \rangle [\emptyset]$$

is equal to

$$\langle W^l, \text{Id}_{\text{im}(\phi_l)}, T[X_{\phi_l}] \rangle [\emptyset] = \langle W^l, T[X_{\phi_l}] \rangle [\emptyset] \cdot \mu[L = l].$$

□

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 θ , see Algorithm 5. Since the log likelihood is concave, the algorithm converges to a global minimum.

Algorithm 5 Alternating Moment Matching

Set $\theta [L] = 0$

Compute $\mu_D [L] = \langle \mathbb{P}^D, \gamma^\phi \rangle [L]$

while Stopping criterion is not met **do**

for $l \in [p]$ **do**

 Compute

$$T^l [X_{\phi_l}] \leftarrow \langle \{\rho^\phi\} \cup \{W^{\tilde{l}} : \tilde{l} \in [p], \tilde{l} \neq l\} \rangle [X_{\phi_l}]$$

 Set $\theta [L = l]$ to a solution of

$$\langle W^l, \text{Id}_{\text{im}(\phi_l)}, T^l \rangle [\emptyset] \leftarrow \langle W^l, T^l \rangle [\emptyset] \cdot \mu_D [L = l].$$

end for

end while

In general, if $\text{im}(\phi_l)$ 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 Section 9.3.

The computation of T^l in Algorithm 5 can be intractable and be replaced by an approximative procedure based on message passing schemes.

4.6 Discussion

Further in Wainwright and Jordan (2008): Convex Duality. Forward mapping coincides with gradient, i.e. $\mu = \nabla A^{(\phi, \nu)}(\theta)$.

In Wainwright and Jordan (2008): The objective is the conjugate dual $(A^{(\phi, \nu)})^*$ of $A^{(\phi, \nu)}$, and backward mapping has an expression by the gradient, i.e. $\theta = \nabla (A^{(\phi, \nu)})^*(\mu)$.

5 Propositional Logics

Propositional logics describes systems with d binary categorical 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_0, \dots, x_{d-1} the indices X_k encode whether the corresponding variable is True.

We here choose a semantic centric approach to propositional logic, by defining formulas as binary tensors. Then we investigate the corresponding syntax of formulas as specification of a tensor network decomposition of the relational encoding of formulas.

5.1 Encoding of Booleans

Propositional logic amounts to reason about Boolean variables, which are categorical variables with 2 possible values. Before applying this insight in the representation of propositional formulas, we first investigate how Boolean calculus can be represented by multilinear operations.

5.1.1 Booleans as categorical variables

To represent Booleans by categorical variables X with two states we use the following standard group homomorphism

$$(\{\text{True}, \text{False}\}, \wedge) \quad \text{and} \quad (\{1, 0\}, \cdot)$$

by the map

$$[\cdot] : \{\text{True}, \text{False}\} \rightarrow \{1, 0\}$$

defined as

$$[\text{True}] = 1 \quad , \quad [\text{False}] = 0.$$

The multiplication is performed in the binary tensor contractions and can thus be interpreted as the \wedge connective performed on Boolean coordinates.

While the conjunction of is in this embedding performed by multiplications, operations like the negation

$$[\neg X] = 1 - [X]$$

are affine linear. Direct applications of these affine linear operations to perform logical calculus will be discussed later in Section ??.

However, in this chapter we will circumvent the problems arising with affine linearity by using the one-hot encoding of Booleans.

5.1.2 One-hot Encoding

Booleans are categorical variables with $m = 2$ states, where we interpret the states $[2] = \{0, 1\}$ by $\{\text{True}, \text{False}\}$. The one-hot encoding of Booleans

$$e : [2] \rightarrow \{e_0[X], e_1[X]\} \subset \mathbb{R}^2$$

is thus understood as an encoding of truth values.

As discussed before, the one-hot encoding is rich enough to represent any function of the state by a linear function on the encoding. The truth of formulas is a function of the truth of atomic formulas, and thus representable by linear functions of the one-hot encodings.

5.2 Semantics of Propositional Formulas

The epistemological commitments 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.

5.2.1 Formulas

Logics is especially strong in interpreting binary 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 40. A propositional formula $f [X_0, \dots, X_{d-1}]$ depending on d atoms X_k is a tensor

$$f [X_0, \dots, X_{d-1}] : \prod_{k \in [d]} [2] \rightarrow [2] \subset \mathbb{R}.$$

We call $x_0, \dots, x_{d-1} \in \prod_{k \in [d]} [2]$ a model of a propositional formula f , if

$$f [X_0 = x_0, \dots, X_{d-1} = x_{d-1}] = 1$$

. If there is a model to a propositional formula, say the formula is satisfiable.

The propositional formulas coincide therefore with the binary tensors (see Definition 11).

Since propositional formulas are binary valued tensors, the generic decomposition of Lemma 20 simplifies to

$$f [X_0, \dots, X_{d-1}] = \sum_{x_0, \dots, x_{d-1} \in \prod_{k \in [d]} [2]} f [X_0 = x_0, \dots, X_{d-1} = x_{d-1}] \cdot e_{x_0, \dots, x_{d-1}} [X_0, \dots, X_{d-1}] \quad (15)$$

$$= \sum_{x_0, \dots, x_{d-1} \in \prod_{k \in [d]} [2] : f [X_0 = x_0, \dots, X_{d-1} = x_{d-1}] = 1} e_{x_0, \dots, x_{d-1}} [X_0, \dots, X_{d-1}]. \quad (16)$$

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 15 (see Definition 73).

We depict this decomposition in the diagrammatic notation by

$$\begin{array}{c} \boxed{f} \\ \hline X_0 \quad X_1 \quad \dots \quad X_{d-1} \end{array} = \sum_{\substack{x_0, \dots, x_{d-1} \in \prod_{k \in [d]} [2] \\ f(x_0, \dots, x_{d-1}) = 1}} \begin{array}{c} \boxed{e_{x_0}} \\ \hline \downarrow X_0 \end{array} \dots \begin{array}{c} \boxed{e_{x_{d-1}}} \\ \hline \downarrow X_{d-1} \end{array}$$

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.

5.2.2 Relational 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 Definition 40). 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 X_f representing the formula satisfaction. We can then build the relational encoding (Definition 14) of that map to represent the formula (see Figure 9).

Definition 41 (Relation Encoding of Formulas). Given a factored system with d atoms X_0, \dots, X_{d-1} and a propositional formula f , we define the relational encoding of f (see Definition 14)

$$\rho^f [X_0, \dots, X_{d-1}, X_f] \in \left(\bigotimes_{k \in [d]} \mathbb{R}^2 \right) \otimes \mathbb{R}^2$$

by

$$\rho^f [X_0, \dots, X_{d-1}, X_f] = \sum_{x_0, \dots, x_{d-1} \in \prod_{k \in [d]} [2]} e_{x_0, \dots, x_{d-1}} [X_0, \dots, X_{d-1}] \otimes e_{f(x_0, \dots, x_{d-1})} [X_f]. \quad (17)$$

We can build relational 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 5. For any formula f we have

$$\rho^f [X_{[d]}, X_f] = f [X_{[d]}] \otimes e_1 [X_f] + \neg f [X_{[d]}] \otimes e_0 [X_f] .$$

In particular

$$f [X_{[d]}] = \langle \{\rho^f [X_{[d]}, X_f], e_1 [X_f]\} \rangle [X_{[d]}] .$$

Proof. We can decompose relational encodings of formulas into the sum (see Figure 9)

$$\rho^f = e_0 \otimes \left(\sum_{x_0, \dots, x_{d-1} : f(x_0, \dots, x_{d-1})=0} e_{x_0, \dots, x_{d-1}} \right) \quad (18)$$

$$+ e_1 \otimes \left(\sum_{x_0, \dots, x_{d-1} : f(x_0, \dots, x_{d-1})=1} e_{x_0, \dots, x_{d-1}} \right) \quad (19)$$

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 15, we notice that the relational 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 relational decomposition ρ^f has coordinates

$$\langle \{\rho^f, e_{x_0, \dots, x_{d-1}}\} \rangle [X_f] = \begin{cases} e_1 & \text{if the world where } X_k = x_k \text{ is a model of } f \\ e_0 & \text{else.} \end{cases} \quad (20)$$

The contractions of the relational encoding therefore calculate whether an assignment of atoms is a model of the formula, using basis calculus (see Theorem 96).

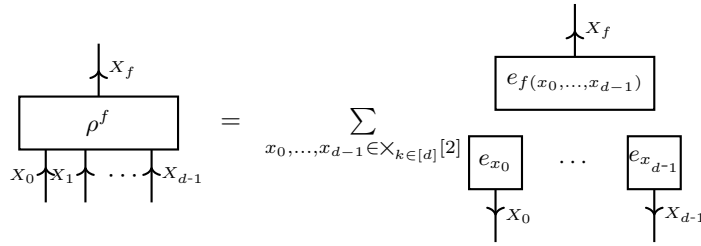


Figure 9: Relational encoding of a propositional formula. The encoding is a sum of the one hot encodings of all states of the factored system in a tensor product with basis vectors, which encode whether the state is a model of the formula. The tensor is directed, since any contraction with an encoded state results in the basis vector evaluating the formula, which we called basis calculus.

5.3 Syntax of Propositional Formulas

Relational encodings of propositional formulas are especially useful when representing function compositions by the representation of their components (see Theorem 97). 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 relational encodings into tensor networks of component encodings for binary tensors following propositional logic syntax.

5.3.1 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 42. The tensors $f_k [X_0, \dots, X_{d-1}]$ defined for x_0, \dots, x_{d-1} as

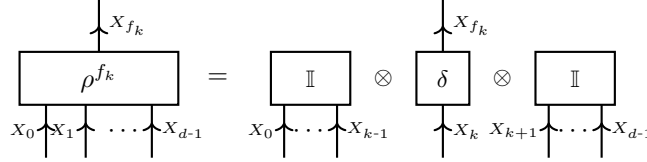
$$f_k [X_0 = x_0, \dots, X_{d-1} = x_{d-1}] = x_k$$

are called atomic formulas.

Theorem 31. The relational encoding of any atomic formula $f_k [X_0, \dots, X_{d-1}]$ has a tensor decomposition by

$$\rho^{X_k} [X_0, \dots, X_{d-1}, X_{f_k}] = \langle \{\delta [X_k, X_{f_k}]\} \rangle [X_0, \dots, X_{d-1}, X_{f_k}] .$$

The decomposition is depicted in a network diagram as



Proof. We have by definition

$$\begin{aligned} \rho^{X_k} [X_0, \dots, X_{d-1}, X_{f_k}] &= \sum_{x_0, \dots, x_{d-1} \in \times_{k \in [d]} [2]} e_{x_0, \dots, x_{d-1}} [X_0, \dots, X_{d-1}] \otimes e_{f_k [X_0=x_0, \dots, X_{d-1}=x_{d-1}]} [X_{f_k}] \\ &= (e_{0,0} [X_k, X_{f_k}] + e_{1,1} [X_k, X_{f_k}]) \otimes \mathbb{I} [X_l : l \neq k] \\ &= \langle \{\delta [X_k, X_{f_k}]\} \rangle [X_0, \dots, X_{d-1}, X_{f_k}] . \end{aligned}$$

□

Remark 7 (Representation of atomic formula tensors for connective action). *Need to represent this as $\langle \delta \otimes \mathbb{I}, e_1^k \rangle$, where the bracket indicates contraction along the k th axis. The core e_1^k can be replaced by further operations based on logical connectives. The axis k is changed from an axis associated with an atom truth to an axis associated with an formula truth.*

5.3.2 Syntactical combination of formulas

Formula tensors are elements of tensor spaces with $d + 1$ axis. The number of coordinates thus grows exponentially with the number of atoms, which is

$$\dim \left[\mathbb{R}^2 \otimes \bigotimes_{k \in [d]} \mathbb{R}^2 \right] = 2^{d+1} .$$

When the number of atoms in a theory is large, the naive representation of formula tensors will be thus intractable. In contrast, most 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 $\circ \in \{\neg, \wedge, \vee, \oplus, \Rightarrow, \Leftrightarrow\}$ can be represented by feasible tensor cores ρ° contracted along a tensor network.

Example 7. We use the following connectives:

- *negation* $\neg : [2] \rightarrow [2]$ by the vector

$$\neg[X_f] = \begin{bmatrix} 0 \\ 1 \end{bmatrix} \quad (21)$$

- *conjunctions* $\wedge : [2] \times [2] \rightarrow [2]$

$$\wedge[X_f, X_h] = \begin{bmatrix} 0 & 0 \\ 0 & 1 \end{bmatrix} \quad (22)$$

- *disjunctions* $\vee : [2] \times [2] \rightarrow [2]$

$$\vee[X_f, X_h] = \begin{bmatrix} 0 & 1 \\ 1 & 1 \end{bmatrix} \quad (23)$$

- *exact disjunction* $\oplus : [2] \times [2] \rightarrow [2]$

$$\oplus[X_f, X_h] = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix} \quad (24)$$

- *implications* $\Rightarrow : [2] \times [2] \rightarrow [2]$

$$\Rightarrow[X_f, X_h] = \begin{bmatrix} 1 & 1 \\ 0 & 1 \end{bmatrix} \quad (25)$$

- *biimplication* $\Leftrightarrow : [2] \times [2] \rightarrow [2]$

$$\Leftrightarrow[X_f, X_h] = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \quad (26)$$

Lemma 6. *Let there be formulas f and h depending on categorical variables $X_{[d]} = X_0, \dots, X_{d-1}$ and a map*

$$\circ : [2] \times [2] \rightarrow [2].$$

Then we have

$$\rho^{f \circ h}[X_{[d]}, X_{f \circ h}] = \langle \{ \rho^\circ[X_f, X_h, X_{f \circ h}], \rho^f[X_{[d]}, X_f], \rho^h[X_{[d]}, X_h] \} \rangle [X_{[d]}, X_{f \circ h}]$$

and for any map $\circ : [2] \rightarrow [2]$

$$\rho^{\circ f}[X_{[d]}, X_{\circ f}] = \langle \{ \rho^\circ[X_f, X_{\circ f}], \rho^f[X_{[d]}, X_f] \} \rangle [X_{[d]}, X_{\circ f}].$$

Proof. This follows from Theorem 97 to be shown in Chapter 15. \square

Theorem 32 (Composition of Formulas). *Let there be a set of binary variables $X_{\mathcal{V}}$ including atoms X_0, \dots, X_{d-1} and image variables to some formulas. For any formula $f[X_0, \dots, X_{d-1}]$, which has a syntactical composition into connectives $\{\circ_l[X_{\mathcal{V}_l}]\} : l \in [p]\}$ taking their inputs by variables $X_{\mathcal{V}_l} \subset X_{\mathcal{V}}$ and output by a variable X_{\circ_l} we have that*

$$\rho^f[X_0, \dots, X_{d-1}, X_f] = \langle \{ \rho^{\circ_l}[X_{\mathcal{V}_l}, X_{\circ_l} : l \in [p]] \} \rangle [X_0, \dots, X_{d-1}, X_f].$$

Proof. When a variable in $X_{\mathcal{V}}$ 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 Theorem 93 can be applied). The claim follows then from iterative application of Lemma 6. \square

Remark 8 (d -ary connectives such as \wedge and \vee). *Since the decomposition of relational encoding can be applied to generic function compositions (see Theorem 97), we can also allow for d -ary connectives*

$$\circ : \bigtimes_{k \in [d]} [2] \rightarrow [2]$$

in Theorem 97 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 relational encoding by any tree-structured composition of binary \wedge and \vee connectives.

In general, any d -ary logical connective is a map

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 the Hadamard product with the trivial tensor \mathbb{I} leaves any tensor invariant, and the contraction with the elementary matrix δ 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 11.

Remark 9 (HT Interpretation of Formula Tensor Networks). *The sketched decomposition of the formula tensor into a network is a hierarchical tree decomposition of the formula tensor, which we will describe in more detail in Section 17.5. At each decomposition of a formula into subformulas, two subspaces spanned by the respective atomic spaces are selected.*

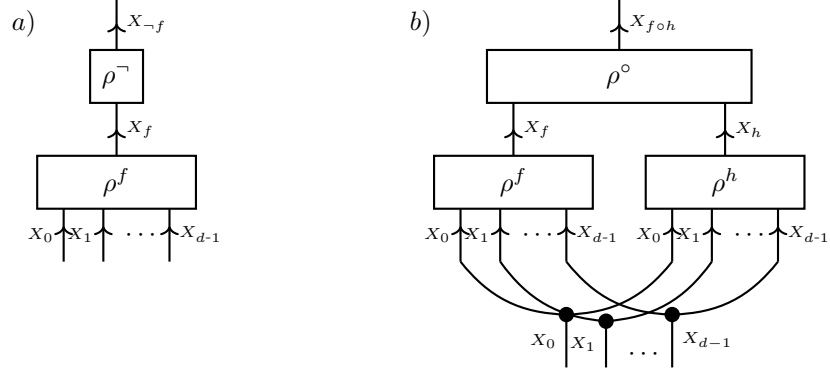


Figure 10: a) Relational encoding of a negated formula f as a tensor network of the encoded formula and the encoded connective \neg . b) Relational encoding of a composition of formulas f, h by a connective $\circ \in \{\wedge, \vee, \oplus, \Rightarrow, \Leftrightarrow\}$. The encoding is a contraction of encodings to f, h and \circ .

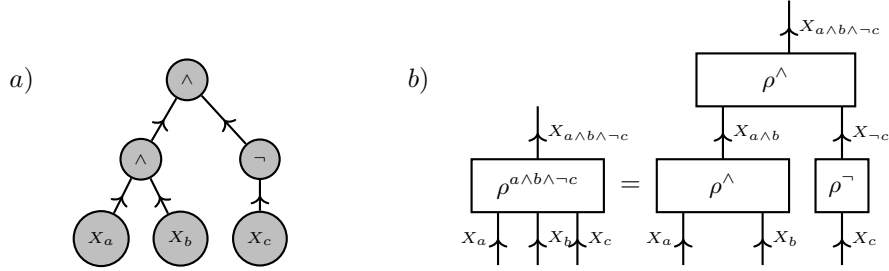


Figure 11: Decomposition of the formula tensor to $f = a \wedge b \wedge \neg c$ into unary (matrix) and binary (third order tensor) cores. a) Visualization of f as a graph. b) Dual Tensor Network decomposition of f . We can make use of the invariance of a Hadamard product with a constant tensor \mathbb{I} and thus not draw axis to atoms not affected by a formula.

5.3.3 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 Definition 40 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 43 (Terms and Clauses). *Given two disjoint subsets \mathcal{V}_0 and \mathcal{V}_1 of the $[d]$, the corresponding term is the formula defined on the indices $x_0, \dots, x_{d-1} \in \times_{k \in [d]} [2]$ by*

$$Z_{\mathcal{V}_0, \mathcal{V}_1}^{\wedge} = \left(\bigwedge_{k \in \mathcal{V}_0} \neg f_k \right) \wedge \left(\bigwedge_{k \in \mathcal{V}_1} f_k \right)$$

and the corresponding clause is the formula defined on the indices $x_0, \dots, x_{d-1} \in \times_{k \in [d]} [2]$ by

$$Z_{\mathcal{V}_0, \mathcal{V}_1}^{\vee} = \left(\bigvee_{k \in \mathcal{V}_0} f_k \right) \vee \left(\bigvee_{k \in \mathcal{V}_1} \neg f_k \right),$$

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 7. *Terms are contractions of one-hot encodings, that is for any disjoint subsets $\mathcal{V}_0, \mathcal{V}_1 \subset [d]$ we have*

$$Z_{\mathcal{V}_0, \mathcal{V}_1}^\wedge [X_{[d]}] = \langle e_{\{x_k=0:k \in \mathcal{V}_0\} \cup \{x_k=1:k \in \mathcal{V}_1\}} \rangle [X_{[d]}].$$

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

$$Z_{\mathcal{V}_0, \mathcal{V}_1}^\vee [X_{[d]}] = \mathbb{I} [X_{[d]}] - \langle e_{\{x_k=0:k \in \mathcal{V}_0\} \cup \{x_k=1:k \in \mathcal{V}_1\}} \rangle [X_{[d]}].$$

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 $x_0, \dots, x_{d-1} \in \times_{k \in [d]} [2]$ we have*

$$Z_{x_0, \dots, x_{d-1}}^\wedge = e_{x_0, \dots, x_{d-1}}$$

Maxterms are substraction of basis elements from the trivial tensor, that is for any $x_0, \dots, x_{d-1} \in \times_{k \in [d]} [2]$ we have

$$Z_{x_0, \dots, x_{d-1}}^\vee = \mathbb{I} [X_{[d]}] - e_{x_0, \dots, x_{d-1}}.$$

Proof. Follows from Lemma 7, since when $\mathcal{V}_0 \cup \mathcal{V}_1 = [d]$ the contractions 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 33. *For any binary tensor $T [X_{[d]}] \in \bigotimes_{k \in [d]} \mathbb{R}^2$ with two-dimensional axes we have*

$$T [X_{[d]}] = \left(\bigvee_{x_0, \dots, x_{d-1}: T[X_{[d]}=x_{[d]}]=1} Z_{\{k:x_k=0\}, \{k:x_k=1\}}^\wedge \right) [X_{[d]}]$$

and

$$T [X_{[d]}] = \left(\bigwedge_{x_0, \dots, x_{d-1}: T[X_{[d]}=x_{[d]}]=0} Z_{\{k:x_k=0\}, \{k:x_k=1\}}^\vee \right) [X_{[d]}].$$

Proof. To show the representation by minterms we use the decomposition

$$T [X_{[d]}] = \sum_{x_{[d]}: T[X_{[d]}=x_{[d]}]=1} e_{x_{[d]}} [X_{[d]}]$$

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

$$T [X_{[d]}] = \mathbb{I} [X_{[d]}] - \sum_{x_{[d]}: T[X_{[d]}=x_{[d]}]=0} e_{x_{[d]}} [X_{[d]}]$$

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 T and are thus identical to T . \square

The decompositions found in Theorem 33 are also called canonical normal forms to propositional formulas $T[X_{[d]}]$.

Remark 10 (Efficient Representation in Propositional Syntax). *The decomposition in Theorem 33 is a basis CP decomposition of the binary tensor and will further be investigated in Chapter 17. The formulas constructed in the proof of Theorem 33 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 Theorem 33. We will provide such constructions in Chapter 17, where we show that dropping the demand of directionality and investigating binary CP Decompositions will improve the sparsity of the polynomial formula representation.*

5.4 Discussion and Outlook

Further study of representing Knowledge Bases based on Tensor Networks of its formulas in Section 8.2 (see Theorem 51).

6 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 by showing how formulas can be interpreted by distributions and define logical entailment based on corresponding probabilistic queries. This enables us to define logical entailment based on the resulting conditional distributions.

Remark 11 (Interpretation of Contractions in Logical Reasoning). *The coordinates of contracted binary 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 worldcount.*

This is consistent with the probabilistic interpretation of contractions, when applying the frequentist interpretation of probability and defining normed worldcounts as probabilities.

6.1 Entailment in Propositional Logics

Definition 44 (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*

$$\forall_{x_{[d]} \in \times_{k \in [d]} [2]} (\mathcal{KB}[X_{[d]} = x_{[d]}] = 1) \rightarrow (f[X_{[d]} = x_{[d]}] = 1).$$

If $\mathcal{KB} \models \neg f$ holds, we say that \mathcal{KB} contradicts f .

Entailment can be understood by subset relations of the models of formulas. This perspective can be applied with subset encodings in Chapter 15.

6.1.1 Deciding Entailment by contractions

Theorem 34 (Contraction Criterion of Entailment). *We have $\mathcal{KB} \models f$ if and only if*

$$\langle \mathcal{KB}, \neg f \rangle [\emptyset] = 0.$$

Proof. If for a x_0, \dots, x_{d-1} we have $\mathcal{KB}[X_{[d]} = x_{[d]}] = 1$ but not $(f(x_0, \dots, x_{d-1}) = 1)$, the contraction would be at least 1. Conversely, if the contraction is at least one, we would find x_0, \dots, x_{d-1} 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

To decide whether a formula is entailed, or its negation is entailed (in which case one says that the formula is contradicted) by a single contraction, one can perform the contraction

$$T = \langle \mathcal{KB}[X_{[d]}], f[X_{[d]}, X_f] \rangle [X_f]$$

and use that

$$\langle \mathcal{KB}, \neg f \rangle [\emptyset] = T[X_f = 0]$$

and

$$\langle \mathcal{KB}, f \rangle [\emptyset] = T[X_f = 1] .$$

6.1.2 Contraction Knowledge Base

We now show how to implement a propositional Knowledge Base with the TELL and ASK operations based on Theorem 4.

Algorithm 6 Contraction Knowledge Base

```

ASK(formula  $f$ )
   $T[X_f] \leftarrow \langle \mathcal{KB}, \rho^f \rangle [X_f]$ 
  if  $T[X_f = 0] = 0$  then
    return Entailed
  end if
  if  $T[X_f = 1] = 0$  then
    return Contradicted
  end if
  return Contingent
TELL(formula  $f$ )
  if ASK( $f$ ) returns Contingent: then
     $\mathcal{KB} \leftarrow \mathcal{KB} \wedge f$ 
  end if

```

Comment: TELL checks whether the formula to be added is entailed, in which case it is redundant to add, and whether the formula to be added is contradicted, in which case the knowledge base would become unsatisfiable.

6.1.3 Sparse Representation of a Knowledge Base

Let us now investigate how to sparsely represent a Knowledge Base. Towards getting insights on this we first show that entailed formulas can be dropped from the Knowledge Base.

Theorem 35. *If and only if $\mathcal{KB} \models f$ we have*

$$\mathcal{KB}[X_{[d]}] = \langle \mathcal{KB}, f \rangle [X_{[d]}] .$$

Proof. For any world indexed by a coordinate x_0, \dots, x_{d-1} , $\mathcal{KB}[X_{[d]} = x_{[d]}]$ indicates whether the world is a model of \mathcal{KB} . We have entailment, when the models of $\mathcal{KB} \cup f$ coincide with those of \mathcal{KB} . \square

Remark 12 (Sparest Description of a Knowledge Base). *Given a set of worlds indexed by T , find the sparsest set of formulas \mathcal{KB} such that*

$$T = \mathcal{KB}$$

would be beneficial for small computational complexity. Since the formula tensors are invariant under entailment, we can drop entailed formulas.

6.2 Formulas as Random Variables

Aim here: Relate with the probabilistic reasoning concepts of marginal and conditional distributions.

Given a probability distribution \mathbb{P} of atoms we add a variable by building the Markov Network of \mathbb{P} and ρ^f to get a joint distribution of the atoms and a query formula f

There are two ways of interpreting formula tensors as conditional probabilities. The standard one, which we also used above, understands the atomic legs as conditions and calculates the truth of the formula. Another understands a formula as a condition.

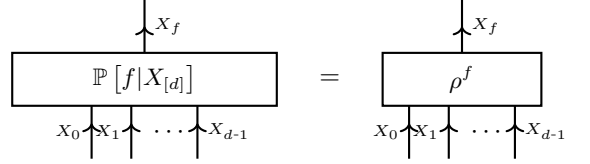
6.2.1 Conditioning on the atoms

Our main interpretation understands each tuple of indices x_0, \dots, x_{d-1} as conditions of a probability tensor. Given a truth assignment to the atomic variables X_k , that is a choice of indices x_k , determines the truth of the formula. We thus interpret the formula tensors as defining a conditional probability of f given the atoms X_k indexed by x_k .

Theorem 36. *The relational encoding of any propositional formula f coincides with the conditional probability of that formula conditioned on the identity on the atoms, that is*

$$\rho^f = \mathbb{P} [X_f | X_{[d]}] .$$

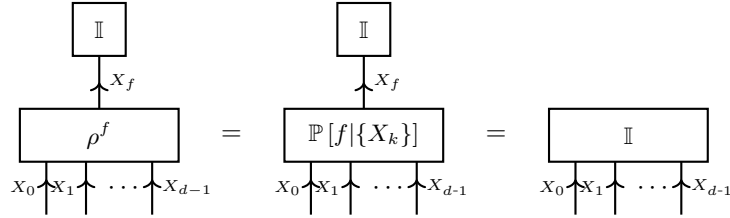
We depict this by



Proof. The distribution \mathbb{P} does not influence the conditional query, since the normation acts on any state. \square

The conditional query $\mathbb{P} [X_f | X_0, \dots, X_{d-1}]$ provides an interpretation of ρ^f as a conditional probability. This is also reflected in the fact that both $\mathbb{P} [f | X_{[d]}]$ and ρ^f are directed, since the first is a normation by Definition 33 and the second a encoding of a function.

This directly implies using Theorem ?? the trivialization of the formula tensor when contracting its head axis indexed by x_f with the trivial vector \mathbb{I} , depicted as

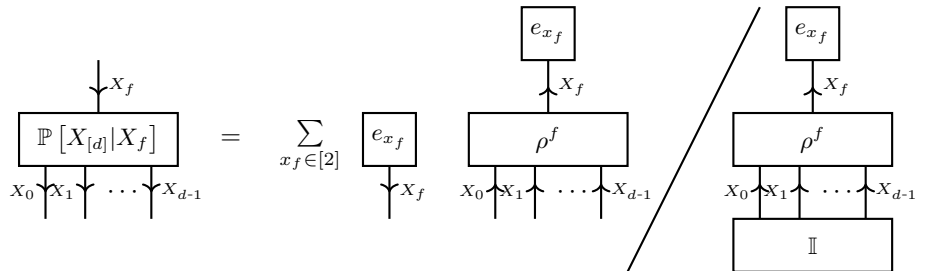


6.2.2 Conditioning on the formula

Let us now converse the order of conditioning from $\mathbb{P} [f | X_{[d]}]$ to $\mathbb{P} [X_{[d]} | f]$. In this way, we have propositional formulas defining probability distributions on the factored system of atoms.

Given a Markov Network \mathbb{P} with a single core ρ^f for a propositional formula f . By definition we have

$$\mathbb{P} [X_0, \dots, X_{d-1} | X_f] = \langle \rho^f \rangle [X_0, \dots, X_{d-1} | X_f] .$$



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

$$\mathbb{P} [X_{[d]} | f = 1] = \frac{1}{\langle f \rangle [\emptyset]} \sum_{x_0, \dots, x_{d-1} \in \times_{k \in [d]} [2] : f[X_{[d]} = x_{[d]}] = 1} e_{x_0, \dots, x_{d-1}} . \quad (27)$$

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 $f = 0$ we get a uniform distribution of the models of $\neg f$.

The probability distribution in Equation (27) is well defined except for the case that $\langle f \rangle [\emptyset] = 0$. In this case we have $f = 0$ 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 (27) correspond within the set of probability distributions with those having constant coordinates on their support. While the probability tensors with nonvanishing coordinates build a $2^d - 1$ -dimensional manifold, where the formulas parametrize 2^{2^d} probability tensors, most of which having vanishing coordinates.

6.2.3 Probability of a function given a Knowledge Base

We can 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 by Theorem 36

$$\begin{aligned} \mathbb{P}[X_f | X_{\mathcal{KB}}] &= \langle \mathbb{P}[X_f | X_{[d]}], \mathbb{P}[X_{[d]} | X_{\mathcal{KB}}] \rangle [X_f, X_{\mathcal{KB}}] \\ &= \langle \rho^f, \rho^{\mathcal{KB}} \rangle [X_f | X_{\mathcal{KB}}] \end{aligned}$$

Of special interest is the marginal probability of X_f given that $X_{\mathcal{KB}}$ is satisfied, that is

$$\begin{aligned} \mathbb{P}[X_f | X_{\mathcal{KB}} = 1] &= \langle \{\rho^f, \mathcal{KB}\} \rangle [X_f | \emptyset] \\ &= \frac{\langle \{\rho^f, \mathcal{KB}\} \rangle [X_f]}{\langle \{\mathcal{KB}\} \rangle [\emptyset]}. \end{aligned}$$

Remark 13 (Case of Unsatisfiable Knowledge Bases). *When the Knowledge Base is not satisfiable, one cannot normate it and the probability distribution is not dedfined.*

Theorem 37. *Given a satisfiable formula \mathcal{KB} , we have $\mathcal{KB} \models f$, if and only if*

$$\mathbb{P}[X_f = 0 | X_{\mathcal{KB}} = 1] = 0.$$

Proof. Since \mathcal{KB} is satisfiable, we have $\langle \mathcal{KB} \rangle [\emptyset] > 0$ and

$$\mathbb{P}[X_f = 0 | X_{\mathcal{KB}} = 1] = \frac{\langle \neg f, \mathcal{KB} \rangle [\emptyset]}{\langle \mathcal{KB} \rangle [\emptyset]}.$$

This term vanishes if and only if $\langle \neg f, \mathcal{KB} \rangle [\emptyset]$ vanish. Thus, the condition is equivalent to the condition in Theorem 34. \square

Given that \mathcal{KB} is satisfiable, we therefore have $\mathcal{KB} \models f$ if and only if

$$\mathbb{P}[X_f | X_{\mathcal{KB}} = 1] = e_1. \quad (28)$$

We depict this condition by the contraction diagram

$$\begin{array}{c} \uparrow^{X_f} \\ \boxed{\rho^f} \\ \uparrow^{X_0} \uparrow^{X_1} \cdots \uparrow^{X_{d-1}} \\ \boxed{\langle \mathcal{KB} \rangle [X_{[d]} | \emptyset]} \end{array} = \begin{array}{c} \uparrow^{X_f} \\ \boxed{e_1} \end{array}.$$

We can omit the normation by $\langle \mathcal{KB} \rangle [\emptyset]$ when deciding entailment, as we state next.

Corollary 4. *Given a satisfiable formula \mathcal{KB} , we have $\mathcal{KB} \models f$ (respectively $\mathcal{KB} \models \neg f$), if and only if*

$$\langle \mathcal{KB}, \rho^f \rangle [X_f = 0] = 0 \quad (\text{respectively } \langle \mathcal{KB}, \rho^f \rangle [X_f = 1] = 0).$$

Relating entailment to probability distributions motivates an extension of Definition44 of entailment to arbitrary probability distributions.

Definition 45. For any propositional formula f and a probability distribution \mathbb{P} we say that \mathbb{P} probabilistically entails f , denoted as $\mathbb{P} \models f$, if

$$\langle \mathbb{P}, \rho^f \rangle [X_f = 0] = 0.$$

If $\mathbb{P} \models \neg f$ we say that \mathbb{P} probabilistically contradicts f .

By Theorem 37 the definition of entailment reduces to propositional formulas by choosing $\mathbb{P} = \langle \mathcal{KB} \rangle [X_{[d]} | \emptyset]$

6.2.4 Knowledge Bases as Base Measures for Probability Distributions

Let us now relate the probabilistic entailment definition 45 with the logical entailment. Given a generic probability distribution \mathbb{P} we can build a Knowledge Base by the indicator function of the support as

$$\mathcal{KB}^{\mathbb{P}} = \chi \circ \mathbb{P}$$

where $\chi : \mathbb{R} \rightarrow \mathbb{R}$ is defined as $\chi(x) = 1$ if $x \neq 0$ and $\chi(x) = 0$ else.

Theorem 38. Any probability distribution \mathbb{P} probabilistically entails a formula f , if and only if the Knowledge Base $\mathcal{KB}^{\mathbb{P}}$ logically entails 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 $\mathcal{KB}^{\mathbb{P}} [X_{[d]} = x_{[d]}] = 1$ and

$$((\mathcal{KB}^{\mathbb{P}} [X_{[d]} = x_{[d]}] = 1) \rightarrow (f [X_{[d]} = x_{[d]}] = 1)).$$

is not satisfied. Together, $\mathbb{P} \models f$ does not holds if and only if

$$\forall X_{[d]} (\mathcal{KB}^{\mathbb{P}} [X_{[d]} = x_{[d]}] = 1) \rightarrow (f [X_{[d]} = x_{[d]}] = 1)$$

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 Definition 17) and positivity (see Definition 18) of distributions with respect to boolean base measures.

Theorem 39. A distribution \mathbb{P} of boolean variables is representable with respect to ν , if and only if $\chi \circ \mathbb{P} \models \nu$. A distribution \mathbb{P} of boolean variables is positive with respect to ν , if and only if $\nu = \chi \circ \mathbb{P}$.

Proof. To show the first claim, let \mathbb{P} be a distribution and ν be a base measure. With Definition 17, \mathbb{P} is representable with respect to ν , if and only if

$$\forall x_{[d]} \in \times_{k \in [d]} [2] (\nu [X_{[d]} = x_{[d]}] = 0) \rightarrow (\mathbb{P} [X_{[d]} = x_{[d]}] = 0).$$

This is equal to

$$\forall x_{[d]} \in \times_{k \in [d]} [2] (\chi \circ \mathbb{P} [X_{[d]} = x_{[d]}] = 1) \rightarrow (\nu [X_{[d]} = x_{[d]}] = 1)$$

and by definition Definition 44 equal to $\nu \models \chi \circ \mathbb{P}$.

To show the second claim, we show that when \mathbb{P} is in addition positive with respect to ν , then also $\nu \models \chi \circ \mathbb{P}$ and thus $\nu = \chi \circ \mathbb{P}$. Let \mathbb{P} be a distribution, which is representable with respect to ν . Then \mathbb{P} is positive with respect to ν , if and only if

$$\forall x_{[d]} \in \times_{k \in [d]} [2] (\nu [X_{[d]} = x_{[d]}] = 1) \rightarrow (\mathbb{P} [X_{[d]} = x_{[d]}] > 0)$$

This is equal to

$$\forall x_{[d]} \in \times_{k \in [d]} [2] (\nu [X_{[d]} = x_{[d]}] = 1) \rightarrow (\chi \circ \mathbb{P} [X_{[d]} = x_{[d]}] = 1)$$

and thus $\nu \models \chi \circ \mathbb{P}$. \square

6.2.5 Deciding entailment on Markov Networks

Theorem 40. Let $\mathcal{T}^G = \{T^e : e \in \mathcal{E}\}$ be a non-negative Tensor Network on a hypergraph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$, $\tilde{\mathcal{V}} \subset \mathcal{V}$ be a subset and

$$\mathbb{P}[X_{\tilde{\mathcal{V}}}] = \langle \{T[e] : e \in \mathcal{E}\} \rangle [X_{\tilde{\mathcal{V}}} | \emptyset]$$

and

$$\tilde{\mathbb{P}}[X_{\tilde{\mathcal{V}}}] = \langle \{\chi \circ T[e] : e \in \mathcal{E}\} \rangle [X_{\tilde{\mathcal{V}}} | \emptyset]$$

Then we have for any f that $\mathbb{P} \models f$ if and only if $\tilde{\mathbb{P}} \models f$.

Proof. We first show

$$\chi \circ \mathbb{P} = \chi \circ \tilde{\mathbb{P}}. \quad (29)$$

The claim follows then from Theorem 38. To show (29) let there be $X_{\tilde{\mathcal{V}}} = x_{\tilde{\mathcal{V}}}$ such that $\mathbb{P}[X_{\tilde{\mathcal{V}}} = x_{\tilde{\mathcal{V}}}] = 0$. Then for any $X_{\mathcal{V}} = x_{\mathcal{V}}$ extending $X_{\tilde{\mathcal{V}}} = x_{\tilde{\mathcal{V}}}$ we have $\langle \{T[e] : e \in \mathcal{E}\} \rangle [X_{\mathcal{V}} = x_{\mathcal{V}}] = 0$ and thus also $\langle \{\chi \circ T[e] : e \in \mathcal{E}\} \rangle [X_{\mathcal{V}} = x_{\mathcal{V}}] = 0$ and $\tilde{\mathbb{P}}[X_{\tilde{\mathcal{V}}} = x_{\tilde{\mathcal{V}}}] = 0$. One can similarly show, that when $\tilde{\mathbb{P}}[X_{\tilde{\mathcal{V}}} = x_{\tilde{\mathcal{V}}}] = 0$ then also $\mathbb{P}[X_{\tilde{\mathcal{V}}} = x_{\tilde{\mathcal{V}}}] = 0$. The support of the distributions \mathbb{P} and $\tilde{\mathbb{P}}$ is thus identical and (29) holds. \square

For any positive tensor T we have

$$\chi \circ T[X_e] = \mathbb{I}[X_e],$$

which does not influence the distribution and can be omitted from the Markov Network. By Theorem 40, 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 $\chi \circ T[X_e]$ 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.

6.3 Deciding Entailment by partial ordering

Classically entailment in propositional logics is defined by a model-theoretic approach. According to that approach, the entailment statement $\mathcal{KB} \models f$ holds, whenever any model of \mathcal{KB} is also a model of f . We will in the following show, that this is equal to our definition based on probabilistic queries.

Definition 46 (Partial ordering of tensors). We say that two tensors f and h in a tensor space $\bigotimes_{k \in [d]} \mathbb{R}^{m_k}$ are partially ordered, denoted by

$$f \prec h,$$

if for all $x_0, \dots, x_{d-1} \in \times_{k \in [d]} [m_k]$

$$f(x_0, \dots, x_{d-1}) \leq h(x_0, \dots, x_{d-1}).$$

We notice, that whenever $f \prec h$ holds, for any model x_0, \dots, x_{d-1} of f we have

$$1 = f(x_0, \dots, x_{d-1}) \leq h(x_0, \dots, x_{d-1})$$

and thus $h(x_0, \dots, x_{d-1}) = 1$. Therefore any model of f is also a model of h . We show in the next theorem, that this is equivalent to the entailment statement $f \models h$.

Theorem 41 (Partial Ordering Criterion). We have $\mathcal{KB} \models f$ if and only if $\mathcal{KB} \prec f$.

Proof. Directly by definition, since both \mathcal{KB} and f are Boolean and therefore for any $x_{[d]} \in \times_{k \in [d]} [2]$ we have that

$$\mathcal{KB}[X_{[d]} = x_{[d]}] \leq f[X_{[d]} = x_{[d]}]$$

is equivalent to

$$(\mathcal{KB}[X_{[d]} = x_{[d]}] = 1) \rightarrow (f[X_{[d]} = x_{[d]}] = 1).$$

Therefore, $\mathcal{KB} \prec f$ is equivalent to

$$\forall_{x_{[d]} \in \times_{k \in [d]} [2]} (\mathcal{KB}[X_{[d]} = x_{[d]}] = 1) \rightarrow (f[X_{[d]} = x_{[d]}] = 1),$$

which is equal to $\mathcal{KB} \models f$. \square

6.3.1 Monotonicity of Entailment

Vanishing local contractions provide sufficient but not necessary criterion to decide entailment, as we show in the next theorem.

Theorem 42 (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 $\mathbb{P}^{\tilde{\mathcal{G}}} \models f$.*

Proof. Based on the reduction to Knowledge Bases by Theorem 38 and the monotonicity of binary contractions as shown in Theorem 92. \square

Remark 14. *To make use of Theorem 42 we can exploit any entailment criterion. However, there is no claim about entailment being false, when the entailment Theorem 42 therefore just provides a sufficient but not necessary criterion of entailment with respect to $\mathbb{P}^{\mathcal{G}}$.*

6.4 Deciding Entailment by local contractions

Global entailment can become inefficient, when

- we are interested in batches of entailment checks. Here we can make use of dynamic programming (store partial contraction results in the Knowledge Cores).
- the network is large. Although efficient tensor network contraction often work, they might get infeasible when the tensor network has a large connectivity. For many

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.

Remark 15. *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.*

6.4.1 Knowledge Propagation

Let us now draw on these insights and store partial entailment results in Knowledge Cores, which is a use of the dynamic programming paradigm. We then iterate over local entailment checks, where we recursively add further entailment checks to be redone due to additional knowledge. We then call the local checks until convergence Entailment Propagation, since different stadia of knowledge are propagated through the network. We describe local Knowledge Propagation in a generic way in Algorithm 7.

Algorithm 7 Knowledge Propagation (KP)

```

Tensor Network  $\mathcal{T}^{\mathcal{G}}$ ,  $K^e = \mathbb{I}[X_e]$ 
while Stopping Criterion is not met do
  Choose  $e$ , subset  $M$  of  $\mathcal{T}^{\mathcal{G}}$  and of  $\{K^e : e \in \mathcal{E}\}$  containing  $K^e$ 
  Update
    
$$K^e \leftarrow \chi \circ \langle M \rangle [X_e]$$

end while

```

Each chosen subset M is understood as a local knowledge base, which is then applied for local entailment.

There are different ways of implementing Algorithm 7, by choosing an order of local knowledge bases M and a stopping criterion.

Theorem 43. *In Entailment Propagation Algorithm 7, K^e is monotonically decreasing with respect to the partial ordering and greater than \tilde{K}^e defined as*

$$\tilde{K}^e = \chi(\langle \mathcal{T}^{\mathcal{G}} \rangle [e]) .$$

Proof. We deduce the theorem from generic properties of the support of contractions, see Section 14.5. Monotonic decreasing follows from monotonicity of tensor contractions, see Theorem 92. By Theorem 93 we have during any state of the algorithm

$$\chi \circ \langle \mathcal{T}^G \rangle [X_V] = \chi \circ \langle \mathcal{T}^G \cup \{K^e : e \in \mathcal{E}\} \rangle [X_V] .$$

It follows that

$$\tilde{K}^e = \chi \circ \langle \mathcal{T}^G \cup \{K^e : e \in \mathcal{E}\} \rangle [X_e]$$

and by Theorem 92

$$\tilde{K}^e \prec K^e .$$

□

Corollary 5. *Whenever for a formula $f [X_V]$ and a K^e we have*

$$\langle K^e, \rho^f \rangle [X_f = 0] = 0$$

then the Markov Network \mathcal{T}^G probabilistically entails f .

Another way to use Algorithm 7 to decide entailment of formulas f is adding each ρ^f to \mathcal{T}^G and defining Knowledge Cores $K^f [X_f]$. Since then the Knowledge Core has only two dimensions, there are only four possible cores with the interpretation

- e_1 : the formula is known to be true
- e_0 : the formula is known to be false
- \mathbb{I} : the formula is not known
- 0 : the knowledge base is inconsistent

Part II

Neuro-Symbolic Learning

We now employ tensor networks to define architectures and algorithms for neuro-symbolic reasoning based on the logical and probabilistic foundations. Markov Logic Networks will be taken as generative models to be learned from data, using formula selecting tensor networks and likelihood optimization algorithms.

7 Formula Selecting Networks

In this chapter we will investigate efficient schemes to represent collections of formulas with similar structure in one tensor network.

Definition 47. *Given a set of p formulas $\{f_l : l \in [p]\}$, we define the formula selecting map as*

$$\mathcal{H} [X_{[d]}, L] : \bigtimes_{k \in [d]} [2] \times [p] \rightarrow [2]$$

defined for $l \in [p]$ by

$$\mathcal{H} [X_{[d]} = x_0, \dots, x_{d-1}, L = l] = f_l [X_{[d]} = x_0, \dots, x_{d-1}] .$$

We introduce a selection variable L and depict the formula selection in Figure 12.

A naive representation of the formula selecting map is as a sum

$$\mathcal{H} = \sum_{l \in [p]} f_l [X_{[d]}] \otimes e_l [L] .$$

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.

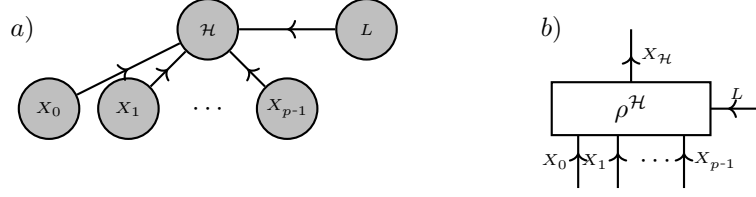


Figure 12: Representation of the Formula Selecting map as a a) Graphical Model with a selection variable \mathcal{H} . b) Dual Tensor Core with selection variable corresponding with an additional axis.

7.1 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.

7.1.1 Connective Selecting Tensors

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 48. Let $\{\circ_0, \dots, \circ_{p_C-1}\}$ be a set of connectives with d arguments. The associated connective selection map is

$$\mathcal{H}_C [X_{[d]}, L_C] : \bigtimes_{k \in [d]} [2] \times [p_C] \rightarrow [2]$$

defined for each $l_C \in [p_C]$ and $x_{[d]} \in \bigtimes_{k \in [d]} [2]$ by

$$\mathcal{H}_C [X_{[d]} = x_{[d]}, L_C = l_C] = \circ_{l_C} [X_{[d]} = x_{[d]}] .$$

We depict the relational encoding of connective selection maps in Figure 13.

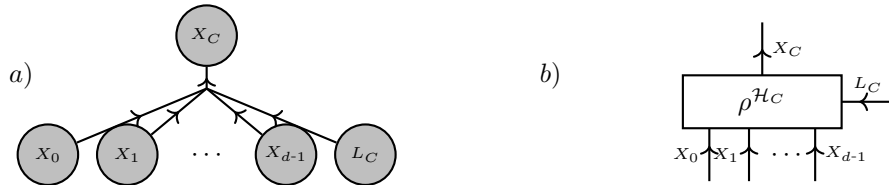


Figure 13: Connective Selector.

Remark 16 (HT Interpretation of Superposed Formula Tensor Networks). *Continuing Remark 9: Superposed Formula Tensors have a decomposition into a HT as sketched here, where we distinguish between formula selection subspaces (indices l_s) and atomic subspaces (indices x_k). At each formula selection we thus have a decomposition into three subspaces, two of atomic formulas and one for the formula selection.*

7.1.2 Variable Selecting Tensor Network

Definition 49. The selection of one out of p variables in a list $X_{[p]}$ is done by variable selecting maps

$$\mathcal{H}_V [X_{[p]}, L_V] : \left(\bigtimes_{l \in [p]} [2] \right) \times [p] \rightarrow [2] \quad (30)$$

are defined coordinatewise by

$$\mathcal{H}_V [X_0 = x_0, \dots, X_{p-1} = x_{p-1}, L_V = l_V] = x_l . \quad (31)$$

Variable selecting maps appear in the literature as multiplex gates (see Definition 5.3 in Koller and Friedman (2009)).

The relational encoding of the variable selection map has a decomposition

$$\rho^{\mathcal{H}_V} [X_{\mathcal{H}_V}, X_{[p_V]}] = \sum_{l_V \in [p_V]} \rho^{X_{l_V}} [X_{\mathcal{H}_V}, X_{l_V}] \otimes e_{l_V} [L_V] .$$

This structure is exploited in the next theorem to derive a tensor network decomposition of $\rho^{\mathcal{H}_V}$.

Theorem 44 (Decomposition of Variable Selecting Maps). *Given a list $X_{[p_V]}$ of variables, we define for each $l_V \in [p_V]$ the tensors*

$$T^{V_{l_V}} [X_{\mathcal{H}_V}, X_{l_V}, L_V] = \delta [X_{\mathcal{H}_V}, X_{l_V}] \otimes e_{l_V} [L_V] + \mathbb{I} [X_{\mathcal{H}_V}, X_{l_V}] \otimes (\mathbb{I} [L_V] - e_{l_V} [L_V]) .$$

Then we have (see Figure 14)

$$\rho^{\mathcal{H}_V} [X_{\mathcal{H}_V}, X_{[p]}, L_V] = \langle \{T^{V_{l_V}} [X_{\mathcal{H}_V}, X_{l_V}, L_V] : l_V \in [p_V]\} \rangle [X_{\mathcal{H}_V}, X_{[p]}, L_V] .$$

Proof. We show the equivalence of the tensors on an arbitrary coordinates. For $\tilde{l}_V \in [p_V]$, $X_{\mathcal{H}_V} \in [2]$ and $x_{[p_V]} \in \times_{k \in [p_V]} [2]$ we have

$$\begin{aligned} & \langle \{T^{V_{l_V}} [X_{\mathcal{H}_V}, X_{l_V}, L_V] : l_V \in [p_V]\} \rangle [X_{\mathcal{H}_V} = x_{\mathcal{H}_V}, X_{[p]} = x_{[p]}, L_V = \tilde{l}_V] \\ &= \prod_{l_V \in [p_V]} T^{V_{l_V}} [X_{\mathcal{H}_V} = x_{\mathcal{H}_V}, X_{l_V} = x_{l_V}, L_V = \tilde{l}_V] \\ &= T^{V_{\tilde{l}_V}} [X_{\mathcal{H}_V} = x_{\mathcal{H}_V}, X_{\tilde{l}_V} = x_{\tilde{l}_V}, L_V = \tilde{l}_V] \\ &= \begin{cases} 1 & \text{if } x_{\mathcal{H}_V} = x_{\tilde{l}_V} \\ 0 & \text{else} \end{cases} \\ &= \rho^{\mathcal{H}_V} [X_{\mathcal{H}_V} = x_{\mathcal{H}_V}, X_{[p]} = x_{[p]}, L_V = \tilde{l}_V] \end{aligned}$$

In the second equality, we used that the tensor $T^{V_{l_V}}$ have coordinates 1 whenever $\tilde{l}_V \neq l_V$. \square

The decomposition provided by Theorem 44 is in a CP format, as will be further discussed in Chapter 17. The introduced tensors $T^{V_{l_V}}$ are Boolean, but not directed and therefore encodings of relations but not functions (see Chapter 15).

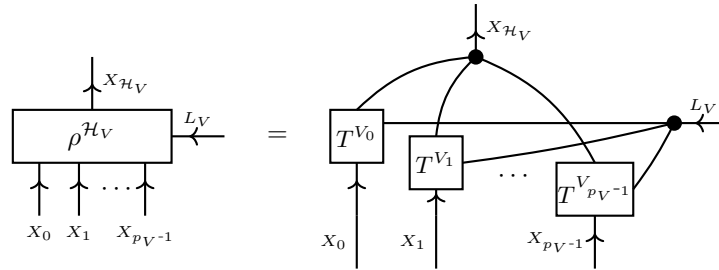


Figure 14: Decomposition of the relational encoding of a variable selecting tensor into a network of tensors defined in Theorem 44. The decomposition is in a CP-Format (see Chapter 17).

7.2 State Selecting Tensors

As an alternative, one can select a state of a categorical variable X .

Definition 50. *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 tensor*

$$\mathcal{H}_S [X_S, L_S] : [m_S] \times [p_S] \rightarrow [2]$$

is defined on $x_S \in [m_S]$ and $l_S \in [p_S]$ by

$$\mathcal{H}_S [X = x, L_S = l_S] = \begin{cases} 1 & \text{if } x = l_S \\ 0 & \text{else} \end{cases} .$$

State selecting tensors can also be realized by variable selecting tensors. In Section 8.2.4 we will describe methods to build atomic variables indicating the states of a categorical variable. This would, however, increase the number of variables in a tensor network and can thus lead to an exponential overhead of dimensions. State selecting tensors can therefore be seen as a mean to avoid such dimension increases.

Comment: State Selectors can be integrated in Variable Selection framework. In this perspective, Variable selection networks are the specific case to $X = 1$.

7.3 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 51. Let there be a formula $f_{l_0, \dots, l_{n-1}}$ for each index tuple in $l_0, \dots, l_{n-1} \in \times_{s \in [n]} [p_s]$, where $n, p_0, \dots, p_{n-1} \in \mathbb{N}$. The folded formula selector map (see Figure 15) is the map

$$\mathcal{H} [X_{[d]}, L_{[n]}] : \left(\times_{k \in [d]} [2] \right) \times \left(\times_{s \in [n]} [p_s] \right) \rightarrow [2]$$

with the coordinates at the indices $x_{[d]} \in \times_{k \in [d]} [2]$, $l_{[n]} \in \times_{s \in [n]} [p_s]$

$$\mathcal{H} [X_{[d]} = x_{[d]}, L_{[n]} = l_{[n]}] = f_{l_{[n]}} [X_{[d]} = x_{[d]}] .$$

We will find formula selector maps by composition variables selector maps (Definition 49) and connective selector maps (Definition 48). This is especially useful to provide efficient decompositions of relational encodings.

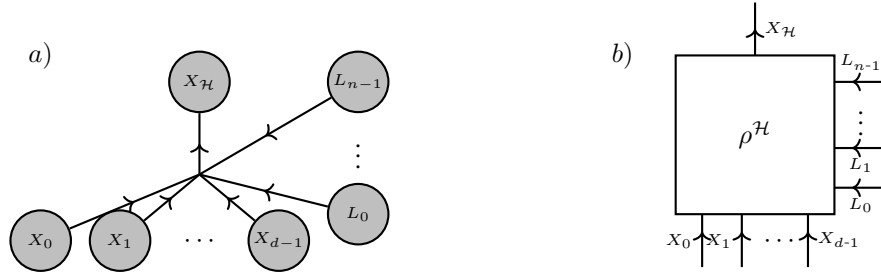


Figure 15: Relational encoding of the folded map \mathcal{H} .

7.3.1 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 with variable selection maps for each argument a formula selecting neuron.

Definition 52. Given an order $n \in \mathbb{N}$ let there be a connective selector L_C 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

$$\sigma [X_{[d]}, L_{[n]}] : \left(\times_{k \in [d]} [2] \right) \times [p_C] \times \left(\times_{s \in [n]} [p_s] \right) \rightarrow [2]$$

defined for $x_{[d]} \in \times_{k \in [d]} [2]$, $l_C \in [p_C]$ and $l_0 \dots l_{n-1} \in \times_{s \in [n]} [p_s]$ by

$$\sigma(x_0, \dots, x_{d-1}, l_C, l_0 \dots l_{n-1}) = \mathcal{H}_C(\mathcal{H}_{V,0}(x_0, \dots, x_{d-1}, l_0), \dots, \mathcal{H}_{V,n-1}(x_0, \dots, x_{d-1}, l_{n-1}), l_C) .$$

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 45. Decomposition of formula selecting neurons Let σ a logical neuron, defined for a connective selector L_C and variable selectors $\mathcal{H}_{V,0}, \dots, \mathcal{H}_{V,n-1}$. Then we have (see Figure 16 for the example of $n = 2$):

$$\begin{aligned} & \rho^\sigma [X_\sigma, X_{[d]}, L_C, L_{V,0}, \dots, L_{V,n-1}] \\ &= \langle \{ \rho^{\mathcal{H}_C} [X_\sigma, X_{V,0}, \dots, X_{V,n-1}] , \\ & \quad \rho^{\mathcal{H}_{V,0}} [X_{V,0}, X_{[d]}, L_{V,0}] , \dots , \rho^{\mathcal{H}_{V,n-1}} [X_{V,n-1}, X_{[d]}, L_{V,n-1}] \} \rangle [X_\sigma, X_{[d]}, L_C, L_{V,0}, \dots, L_{V,n-1}] . \end{aligned}$$

Proof. By composition Theorem 97. \square

Example of a formula selecting neuron: Given a skeleton expression and a set of candidates at each placeholder, we parameterize a set of formulas by the assignment of candidate atoms to each placeholder position. Let us denote the set of formulas, which are generated through choosing atoms from \mathcal{M}^s for the skeleton formula S by

$$\mathcal{F}_S := \{S(Z^1, \dots, Z^d) : Z^k \in \mathcal{M}^k\}$$

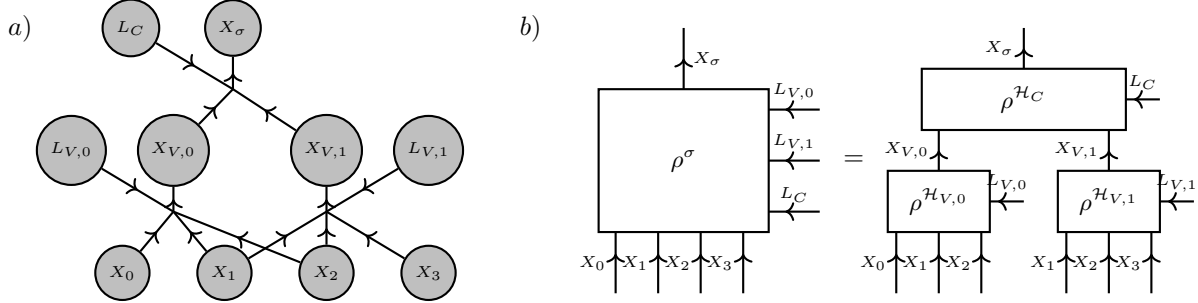


Figure 16: Example of a logical neuron σ of order $n = 2$. a) Selection and categorical variables and their interdependencies visualized in a hypergraph. b) Relational encoding of the logical neuron and tensor network decomposition into variable selecting and connective selecting tensors.

7.3.2 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 53. An architecture graph $\mathcal{G}^{\mathcal{A}} = (\mathcal{V}^{\mathcal{A}}, \mathcal{E}^{\mathcal{A}})$ 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}^{in} and nodes not appearing as outgoing nodes are the output neurons in the set \mathcal{A}^{out} (see Figure 17 for an example).

Given an architecture graph $\mathcal{G}^{\mathcal{A}} = (\mathcal{V}^{\mathcal{A}}, \mathcal{E}^{\mathcal{A}})$, a formula selecting neural network $\mathcal{H}_{\mathcal{A}}$ is a tensor network of logical neurons at each $\sigma \in \mathcal{V}^{\mathcal{A}}/\mathcal{A}^{\text{in}}$, such that each neuron depends on variables $X_{\text{Pa}(\sigma)}$ and on selection variables L_σ . The collection of all selection variable is notated by $L_{\mathcal{A}}$.

The activation tensor of each neuron $\sigma \in \mathcal{V}^{\mathcal{A}}/\mathcal{A}^{\text{in}}$ is

$$\sigma^{\mathcal{A}}[X_{\mathcal{A}^{\text{in}}}, L_{\mathcal{A}}] = \langle \{\rho^{\tilde{\sigma}} : \tilde{\sigma} \in \mathcal{V}^{\mathcal{A}}/\mathcal{A}^{\text{in}}\} \cup \{e_1[X_\sigma]\} \rangle [X_{\mathcal{A}^{\text{in}}}, L_{\mathcal{A}}] .$$

The activation tensor of the formula selecting neural network is the contraction

$$\mathcal{H}_{\mathcal{A}}[X_{\mathcal{A}^{\text{in}}}, L_{\mathcal{A}}] = \langle \{\rho^{\sigma^{\mathcal{A}}}[X_\sigma, X_{\text{Pa}(\sigma)}, L_{\mathcal{A}}] : \sigma \in \mathcal{V}^{\mathcal{A}}/\mathcal{A}^{\text{in}}\} \cup \{e_1[X_\sigma] : \sigma \in \mathcal{A}^{\text{out}}\} \rangle [X_{\mathcal{A}^{\text{in}}}, L_{\mathcal{A}}] .$$

The expressivity of a formula selecting neural network $\mathcal{H}_{\mathcal{A}}$ is the formula set

$$\mathcal{F}_{\mathcal{A}} = \left\{ \mathcal{A}[X_{\mathcal{A}^{\text{in}}}, L_{\mathcal{A}} = l_{\mathcal{A}}] : l_{\mathcal{A}} \in \prod_{s \in [n]} [p_s] \right\} .$$

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 $\sigma \in \mathcal{V}^{\mathcal{A}}/\mathcal{A}^{\text{in}}$ by an output variable X_σ and further add selection variables L_σ to the directed edges, that is to each directed hyperedge $(\{\sigma\}, \text{Pa}(\sigma)) \in \mathcal{E}^{\mathcal{A}}$ we construct a directed hyperedge $(\{X_\sigma\}, X_{\text{Pa}(\sigma)} \cup L_\sigma)$.

Theorem 46. Given fixed selection variables $L_{\mathcal{A}}$, the formula selecting neural network is the conjunction of output neurons, that is

$$\mathcal{H}_{\mathcal{A}}[X_{\mathcal{A}^{\text{in}}}, L_{\mathcal{A}}] = \bigwedge_{\sigma \in \mathcal{A}^{\text{out}}} \sigma[X_{\mathcal{A}^{\text{in}}}, L_{\mathcal{A}}] .$$

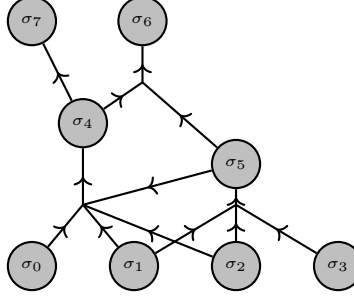


Figure 17: Example of an architecture graph \mathcal{G}^A with input neurons $\mathcal{A}^{\text{in}} = \{\sigma_0, \sigma_1, \sigma_2, \sigma_3\}$ and output neurons $\mathcal{A}^{\text{out}} = \{\sigma_6, \sigma_7\}$

Proof. By effective calculus (see Theorem 100), we have

$$\langle \rho^\wedge [X_\wedge, X_{[d]}], e_1 [X_\wedge] \rangle [X_{[d]}] = \bigotimes_{k \in [d]} e_1 [X_k]$$

and thus

$$\mathcal{H}_{\mathcal{A}} [X_{\mathcal{A}^{\text{in}}}, L_{\mathcal{A}}] = \langle \{\rho^\sigma : \sigma \in \mathcal{V}^{\mathcal{A}} / \mathcal{A}^{\text{in}}\} \cup \{\rho^\wedge [X_\wedge, X_\sigma : \sigma \in \mathcal{A}^{\text{out}}], e_1 [X_\wedge]\} \rangle [X_{\mathcal{A}^{\text{in}}}, L_{\mathcal{A}}] .$$

□

By the commutation of contractions, we can further use Theorem 45 to decompose each tensor ρ^σ into connective and variable selecting components to get a sparse representation of a formula selecting neural network $\mathcal{H}_{\mathcal{A}}$.

7.4 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.

7.4.1 Representation of selection encodings

In technical perspective: FSN provide efficient representation of $\gamma^{\mathcal{F}}$ -> Use for exponential families, structure learning.

Lemma 8. 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

$$\gamma^{\mathcal{F}} [X_{[d]}, L_{[n]}] = \mathcal{H} [X_{[d]}, L_{[n]}] .$$

Proof. For any indices $l_{[n]} \in \times_{s \in [n]} [p_s]$ and $x_{[d]} \in \times_{k \in [d]} [2]$ we have

$$\gamma^{\mathcal{F}} [X_{[d]} = x_{[d]}, L_{[n]} = l_{[n]}] = f_{l_0, \dots, l_{n-1}}(x_0, \dots, x_{d-1}) = \mathcal{H} [X_{[d]} = x_{[d]}, L_{[n]} = l_{[n]}] .$$

□

Technically, relational 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.

7.4.2 Efficient Representation of Formulas

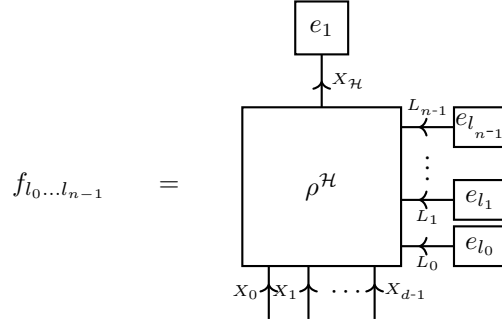
Weight contracted at the selection variables is elementary, then single formula retrieved.

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 extending the contractions by additional cores leaving less variable legs open. This is the case when contracting gradients of the parameter tensor networks in alternating least squares approaches. Other methods avoiding the bottleneck can be constructed by MCMC sampling, for example Gibbs Sampling. Here we only need to vary local components of the formula reflected in keeping only single variable legs open.

7.4.3 Batch contraction of parametrized formulas

Given a set \mathcal{F} of formulas, we build a formula selecting network parametrizing the formulas. The contraction

$$\langle \mathcal{T}^{\mathcal{G}}, \mathcal{H} \rangle [L_{[n]}]$$

is a tensor containing the contractions of the formulas $f_{l_{[n]}}$ with an arbitrary tensor network $\mathcal{T}^{\mathcal{G}}$ as

$$\langle \mathcal{T}^{\mathcal{G}}, f_{l_{[n]}} \rangle [\emptyset] = \langle \mathcal{T}^{\mathcal{G}}, \mathcal{H} \rangle [L_{[n]} = l_{[n]}] .$$

7.4.4 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 8 (Average satisfaction of formulas). *The average of the formula satisfactions in \mathcal{F} given a probability tensor \mathbb{P} is*

$$\frac{1}{\prod_{s \in [n]} p_s} \cdot \langle \mathbb{P}, \gamma^{\mathcal{F}} \rangle [\emptyset] .$$

Example 9 (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*

$$\langle \mathcal{KB}, \gamma^{\mathcal{F}} \rangle [\emptyset] = 0 .$$

We use Lemma 8 to get that $\gamma^{\mathcal{F}} = \mathcal{H}$. 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 $\langle \mathcal{KB}, \gamma^{\mathcal{F}} \rangle [\emptyset]$. This is further equal to

$$\mathcal{KB} \models \neg \left(\bigvee_{f \in \mathcal{F}} f \right).$$

7.5 Examples of formula selecting neural networks

7.5.1 Correlation

For example (see Figure 18) 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

$$T[L_{V,0}, L_{V,1}] = \langle \mathcal{H}_A \rangle [L_{V,0}, L_{V,1}]$$

to read of covariances as

$$\text{Cov}(X_k, X_l) = T[L_{V,0} = k, L_{V,1} = l] - T[L_{V,0} = k, L_{V,1} = k] \cdot T[L_{V,0} = l, L_{V,1} = l].$$

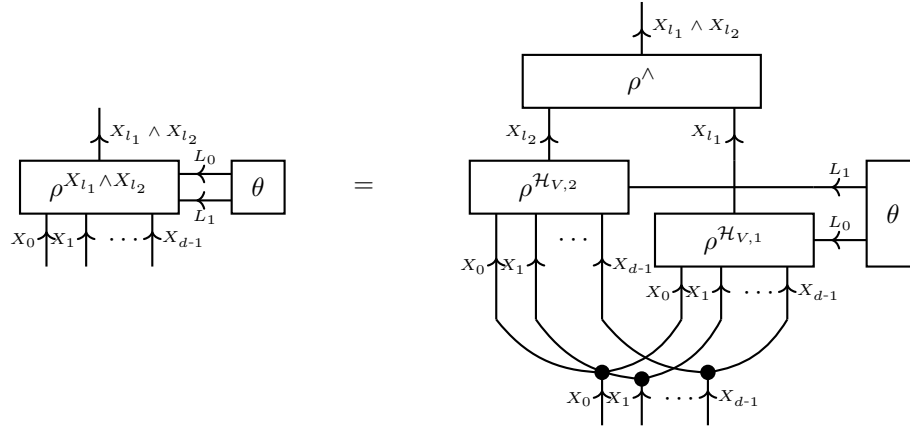


Figure 18: Superposition of the encoded formulas $\rho^{X_{l_1} \wedge X_{l_2}}$ with weight $\theta_{l_1 l_2}$

7.5.2 Conjunctive and Disjunctive Normal Forms

We can represent any propositional knowledge base by the following scheme: Literal selecting neurons by connective identity/negation (selecting positive/negative literal) and selecting one of the atoms. Single output neuron representing the disjunction combining the literal selecting neurons. 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).

The parameter core is in the basis CP format and each slice selects a clause of the knowledge base. When taking the slice values to infinity (e.g. by an annealing procedure), the represented member of the exponential family converges to the uniform distribution of the models of the knowledge base.

Useful to derive basis+ CP format based on CNF!

Remark 17 (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/Id connective selection neuron to each atom and the upper layer consists of a single dary conjunction.*

7.6 Extension to variables of larger dimension

Connective selecting tensors: Can encode arbitrary functions h_l of discrete variables, but need $X_{\mathcal{H}_C}$ to be an enumeration of the, 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

Example 10 (Discretization of a continuous neuron). *Let there be a neuron*

$$\sigma(w, y) : \mathbb{R} \times \mathbb{R}^d \rightarrow \mathbb{R}.$$

When $w \in \mathcal{U}^{weight} \subset \mathbb{R}^d$ and $x \in \mathcal{U}^x \subset \mathbb{R}^d$ have

$$|\sigma(\langle w, x \rangle)| \leq |\mathcal{U}^{weight}| \cdot |\mathcal{U}^x|.$$

To represent the discretization of the neuron, we use the subset encoding scheme of Definition 73. The variables O_{weight} indexing \mathcal{U}^{weight} will be understood as selection incoming variables and the variables O_{weight} indexing \mathcal{U}^{weight} as categorical incoming variables. We further define a variable O_σ indexing $\text{im}(\sigma|_{\mathcal{U}^{weight} \times \mathcal{U}^x})$ and have a tensor

$$\rho^\sigma [O_\sigma, O_{\mathcal{U}^x}, O_{weight}].$$

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 relational encodings for each, can be done.

8 Logic 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.

8.1 Markov Logic Networks

Markov Logic Networks exploit the efficiency and interpretability of logical calculus as well as the expressivity of graphical models.

8.1.1 Markov Logic Networks as Exponential Families

We introduce Markov Logic Networks in the formalism of exponential families (see Section 3.7).

Definition 54 (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 feature is contained in $[2]$, they are propositional formulas (see Def. 40).

Conversely, any binary feature ϕ_l of an exponential family defines a propositional formula (see Definition 40). Thus, any exponential family of distributions of $\prod_{k \in [d]} [2]$, such that $\text{im}(\phi_l) \subset \{0, 1\}$ 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 relational 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 18 (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 11.*

8.1.2 Tensor Network Representation

Based on the previous discussion on the representation of exponential families by tensor networks in Section 3.7 we now derive a representation for Markov Logic Networks.

Theorem 47 (Relational Encodings for Markov Logic Networks). *A Markov Logic Network to a set of formulas $\mathcal{F} = \{f_l : l \in [p]\}$ is represented as*

$$\mathbb{P}^{(\mathcal{F}, \theta)}[X_{[d]}] = \left\langle \{\rho^{f_l}[X_{f_l}, X_{[d]}] : l \in [p]\} \cup \{W^{f_l, \theta[L=l]}[X_{f_l}] : l \in [p]\} \right\rangle [X_{[d]} | \emptyset]$$

where we denote for each $l \in [p]$ an activation core

$$W^{f_l, \theta[L=l]}[X_{f_l}] = \left[\exp \left[\theta [L = l] \right] \right] [X_{f_l}].$$

Proof. The claim follows from Theorem 12 and the following contraction equations. We have with the grouped variable $X_{\mathcal{F}} = \{X_{f_l} : l \in [p]\}$

$$\rho^{\mathcal{F}}[X_{[d]}, X_{\mathcal{F}}] = \left\langle \{\rho^{f_l}[X_{f_l}, X_{[d]}] : l \in [p]\} \right\rangle [X_{[d]}, X_{\mathcal{F}}].$$

Since we have a Markov Logic Network we have $\text{im}(f_l) \subset [2]$ and thus

$$W^{f_l, \theta[L=l]}[X_{f_l} = x_{f_l}] = \begin{cases} 1 & \text{for } X_{f_l} = 0 \\ \exp[\theta[L=l]] & \text{for } X_{f_l} = 1 \end{cases}$$

Using these equations, the claim follows from Theorem 12. \square

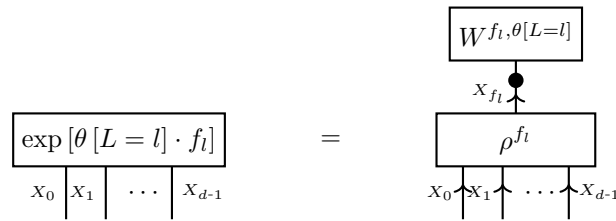


Figure 19: Factor of a Markov Logic Network to a formula f_l .

Since any member of an exponential family is a Markov Network with tensors to each coordinate of the statistic, also Markov Logic Networks are Markov Networks.

Corollary 6. *Given a set \mathcal{F} of formulas on atomic variables $X_{\mathcal{V}}$, we construct a $\mathcal{G} = (\mathcal{V}, \mathcal{E})$, where \mathcal{V} are decorated by the atoms and*

$$\mathcal{E} = \{\mathcal{V}^f : f \in \mathcal{F}\},$$

where by \mathcal{V}^f we denote the minimal set such that there exists a tensor $T[X_{\mathcal{V}^f}]$ with

$$f[X_{\mathcal{V}}] = T[X_{\mathcal{V}^f}] \otimes \mathbb{I}[X_{\mathcal{V}/\mathcal{V}^f}].$$

Any Markov Logic Network (\mathcal{F}, θ) is then a Markov Network given the graph $\mathcal{G}^{\mathcal{F}} \{\exp[\theta[L=l] \cdot f_l] : l \in [p]\}$.

Markov Logic Networks are Markov Networks with the factors given in a restricted form from the weighted truth of a formula. Each formula is seen as a factor of the graphical model.

There are two sparsity mechanisms drastically reducing the number of parameters (and loosing generality):

- Factors/Formulas contain only subsets of atoms (already in Corollary 6 exploited): The underlying assumptions of conditional independence loss generality.
- Structure in the factors: In MLN each factor corresponds with a formula evaluated on possible worlds. Again, any possible factor can be represented by a formula, but we concentrate on small formulas (see Theorem ??).

We can extend the set of variables, by including the hidden formulas, and get a Markov Network of the relational encodings of connectives and headcores. Here hidden variables are additional variables facilitating the decomposition, but not appearing in open variables of contractions when doing reasoning. One can then exploit redundancies and make sure that every subresult is computed just once, by dropping relational encodings with identical head functions.

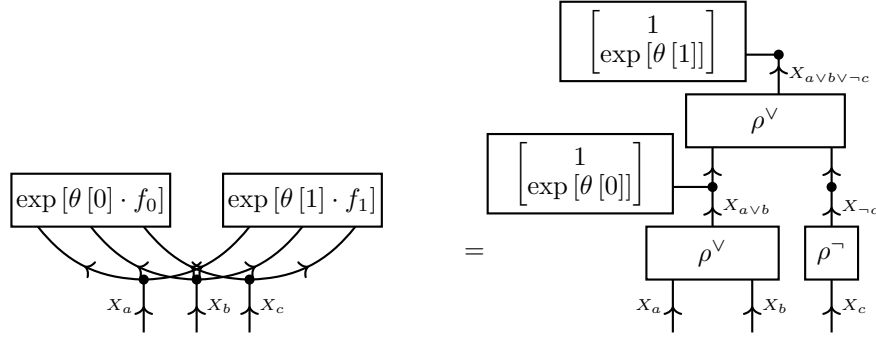


Figure 20: Example of a decomposed Markov Network representation of a Markov Logic Network with formulas $\{f_0 = a \vee b, f_1 = a \vee b \vee \neg c\}$. Since both formulas share the subformula $a \vee b$, their contracted factors have a representation by a connected tensor network.

8.1.3 Energy tensors

With the energy tensor

$$E^{(\mathcal{F}, \theta)} = \sum_{l \in [p]} \theta[L = l] \cdot f_l[X_{[d]}] = \langle \gamma^{\mathcal{F}}[X_{[d]}, L], \theta[L] \rangle [X_{[d]}] \quad (32)$$

the MLN is the distribution

$$\mathbb{P}^{(\mathcal{F}, \theta)}[X_{[d]}] = \langle \exp[E^{(\mathcal{F}, \theta)}] \rangle [X_{[d]} | \emptyset] . \quad (33)$$

In case of a common structure of the formulas in a Markov Logic Network, Formula selecting networks 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 θ in a tensor network format to improve the variance of estimations (see Chapter 10) 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 relational 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.

8.1.4 Expressivity

Based on Markov Logic Networks containing only maxterms and minterms (see Definition 43), we here provide an expressivity study. There are 2^d maxterms and 2^d minterms which are enough to represent any probability distribution as we show next.

Theorem 48. *Let there be a positive probability tensor*

$$\mathbb{P}[X_{[d]}] \in \bigotimes_{k \in [d]} \mathbb{R}^2.$$

Then the Markov Logic Network of minterms (see Definition 43)

$$\mathcal{F}_\wedge = \{Z_{x_0, \dots, x_{d-1}}^\wedge : x_0, \dots, x_{d-1} \in \bigtimes_{k \in [d]} [2]\}$$

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

$$E^{(\mathcal{F}, \theta)} = \ln \mathbb{P}[X_{[d]}]$$

and therefore

$$\mathbb{P}^{(\mathcal{F}, \theta)}[X_{[d]}] = \left\langle \exp \left[E^{(\mathcal{F}, \theta)} \right] \right\rangle [X_{[d]} | \emptyset] = \left\langle \exp \left[E^{(\mathcal{F}, \theta)} \right] \right\rangle [X_{[d]}] = \mathbb{P}[X_{[d]}].$$

In the case of minterms, we notice that for any $x_0, \dots, x_{d-1} \in \bigtimes_{k \in [d]} [2]$

$$Z_{x_0, \dots, x_{d-1}}^\wedge [X_{[d]}] = e_{x_0, \dots, x_{d-1}} [X_{[d]}]$$

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

$$Z_{x_0, \dots, x_{d-1}}^\vee [X_{[d]}] = \mathbb{I}[X_{[d]}] - e_{x_0, \dots, x_{d-1}} [X_{[d]}]$$

and thus that the maximal clauses coincide with the one-hot encodings of respective states. We thus have

$$\begin{aligned} & \sum_{x_0, \dots, x_{d-1} \in \bigtimes_{k \in [d]} [2]} (-\ln \mathbb{P}[X_0 = x_0, \dots, X_{d-1} = x_{d-1}]) \cdot Z_{x_0, \dots, x_{d-1}}^\vee [X_{[d]}] \\ &= \left(\sum_{\nu_0 \subset [d]} (-\ln \mathbb{P}[X_0 = x_0, \dots, X_{d-1} = x_{d-1}]) \cdot \mathbb{I}[X_{[d]}] \right) \\ &+ \left(\sum_{\nu_0 \subset [d]} (\ln \mathbb{P}[X_0 = x_0, \dots, X_{d-1} = x_{d-1}]) \cdot e_{x_0, \dots, x_{d-1}} [X_{[d]}] \right) \\ &= \ln \mathbb{P}[X_{[d]}] + \lambda \cdot \mathbb{I}[X_{[d]}], \end{aligned}$$

where λ is a constant. □

In general, this representation is redundant, since any offset of the weight by $\lambda \cdot \mathbb{I}$ results in the same distribution. However, the only θ are multiples of $\mathbb{I} [X_{[d]}]$.

Theorem 48 is the analogue in Markov Logic to Theorem 33, 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.

Remark 19 (Representation of Markov Networks). *If a probability distribution is representable as a Markov Network, we only need to activate clauses and terms, which variables are contained in factors. **Make a theorem out of that?***

8.1.5 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 49. *Let $\mathbb{P} [X_{[d]}]$ be a positive probability distribution, such that disjoint subsets of atoms 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

$$\theta^k = \ln \left[\frac{\mathbb{P}^k [X_k = 1]}{\mathbb{P}^k [X_k = 0]} \right]$$

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

$$\theta [L] = \sum_{k \in [d]} \theta [L = k] \cdot e_k [L]$$

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 49 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 Section 8.2).

All atomic formulas can be selected by a single variable selecting tensor, that is

$$E^{(\{X_k : k \in [d]\}, \theta)} [X_{[d]}] = \langle \mathcal{H}_V [X_{[d]}, L], \theta [L] \rangle [X_{[d]}] .$$

In case of negative coordinantes $\theta [L = k]$ 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 $-\theta [L = k]$.

8.1.6 Boltzmann machines

A Boltzmann machine is a member of an exponential family with the energy tensor (see e.g. Chapter 43 in MacKay (2003))

$$E^{W,b}[X_0 = x_0, \dots, X_{d-1} = x_{d-1}] = \sum_{k,l \in [d]} W[L_{V,0} = k, L_{V,1} = l] \cdot x_k \cdot x_l + \sum_{k \in [d]} b[L_{V,0} = k] \cdot x_k.$$

We notice, that this 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

$$\triangleleft[X_0 = x_0, X_1 = x_1] = \mathcal{H}_V[X_0 = x_0, X_1, L_V = 0].$$

The weight is

$$\theta = e_0[L_C] \otimes W + e_1[L_C] \otimes b[L_{V,0}] \otimes e_0[L_{V,0}]$$

And we have

$$E^{W,b}[X_{[d]}] = \langle \mathcal{H}_A[X_{[d]}, L_C, L_{V,0}, L_{V,1}], \theta[L_C, L_{V,0}, L_{V,1}] \rangle [X_{[d]}].$$

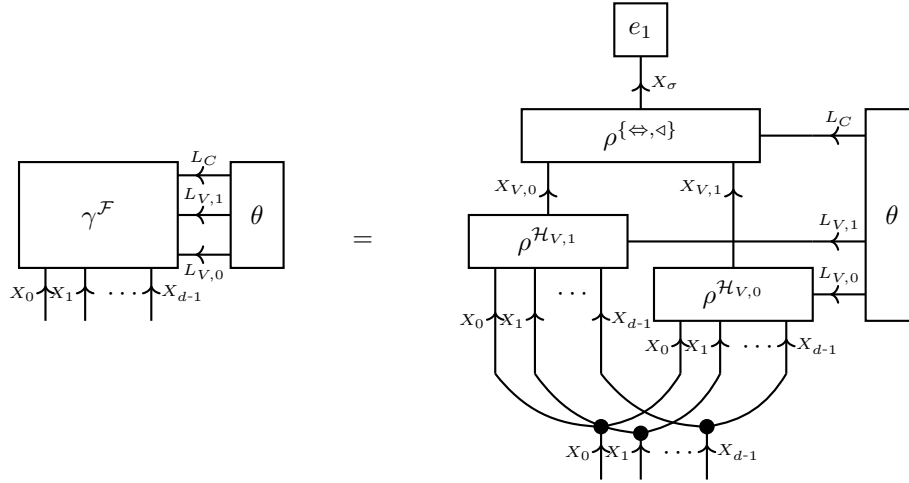


Figure 21: Tensor network representation of the energy of a Boltzmann machine

Often Boltzmann machines are formulated with hidden variables. To average those out, one needs to instantiate the probability distribution instead of the energy tensor and leave only visible variables open in a contraction.

Markov Logic Networks go beyond the Boltzmann machines already for binary formulas, by the flexibility to capture further dependencies beyond the correlation. We can use any binary logical connective and have an associated formula where we can put a weight on.

8.2 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.

8.2.1 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 9. *When taking the limit of large weights $\theta_f \rightarrow \infty$ we observe a coordinatewise (in the sense of a convergence of each coordinate of the tensor) convergence*

$$\langle \exp [\theta_f \cdot f] \rangle [X_{[d]} | \emptyset] \rightarrow \langle f \rangle [X_{[d]} | \emptyset] \quad (34)$$

Proof. We have

$$\mathcal{Z}((\mathcal{F}, \theta)) = \left(\prod_k m_k - \langle f \rangle [\emptyset] \right) + \langle f \rangle [\emptyset] \cdot \exp [\theta_f]$$

and therefore for any $x_0, \dots, x_{d-1} \in \times_{k \in [d]} [2]$ with $f(x_0, \dots, x_{d-1}) = 1$

$$\begin{aligned} \langle \exp [\theta_f \cdot f] \rangle [X_0 = x_0, \dots, X_{d-1} = x_{d-1} | \emptyset] &= \frac{\exp [\theta_f]}{(\prod_k m_k - \langle f \rangle [\emptyset]) + \langle f \rangle [\emptyset] \cdot \exp [\theta_f]} \\ &\rightarrow \frac{1}{\langle f \rangle [\emptyset]} = \langle f \rangle [X_0 = x_0, \dots, X_{d-1} = x_{d-1} | \emptyset]. \end{aligned}$$

For any $x_0, \dots, x_{d-1} \in \times_{k \in [d]} [2]$ with $f(x_0, \dots, x_{d-1}) = 0$ we have on the other side

$$\begin{aligned} \langle \exp [\theta_f \cdot f] \rangle [X_0 = x_0, \dots, X_{d-1} = x_{d-1} | \emptyset] &= \frac{1}{(\prod_k m_k - \langle f \rangle [\emptyset]) + \langle f \rangle [\emptyset] \cdot \exp [\theta_f]} \\ &\rightarrow 0 = \langle f \rangle [X_0 = x_0, \dots, X_{d-1} = x_{d-1} | \emptyset]. \end{aligned}$$

□

Theorem 50. *Let \mathcal{F} be a formulaset and θ 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 β the limit is a uniform distribution and coincides therefor with

$$\langle \mathcal{KB} \rangle [X_{[d]}].$$

□

Remark 20 (More generic situation of simulated annealing). *The process of taking $\beta \rightarrow \infty$ is known as simulated annealing, see Chapter 4. From the discussion there we have the more general statement, that the limiting distribution is the uniform distribution among the maxima of $\mathbb{P}^{(\mathcal{F}, \theta)}[X_{[d]}]$. If the formula \mathcal{KB} is not satisfiable the normation $\langle \mathcal{KB} \rangle [X_{[d]} | \emptyset]$ does not exist and the limit distribution has another syntactical representation, to be gained e.g. by minterm or maxterm representation (see Theorem 33).*

8.2.2 Tensor Network Representation

Hard Logic Network coincide with Knowledge Bases and are thus representable by contractions of formulas (which can be interpreted as an effective calculus scheme, see Section 15.4). We use \wedge symmetry to represent them as a contraction of the formulas building the Knowledge Base as conjunction.

Theorem 51 (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

$$\mathcal{KB} [X_{[d]}] = \langle \{\rho^f [X_f, X_{[d]}] : f \in \mathcal{F}\} \cup \{e_1 [X_f] : f \in \mathcal{F}\} \rangle [X_{[d]}] .$$

Proof. By the \wedge -symmetry, see effective calculus and

$$f [X_{[d]}] = \langle \{\rho^f [X_f, X_{[d]}], e_1 [X_f]\} \rangle [X_{[d]}]$$

□

Remark 21. \wedge symmetry does not generalize to Markov Logic Networks In Markov Logic, similar decompositions are not possible. For example, consider a MLN with a single formula $X_0 \wedge X_1$ and nonvanishing weight θ . 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.

8.2.3 Polynomial Representation

We now apply the representation symmetries to represent a propositional Knowledge Base in conjunctive normal form. A Knowledge Base in Conjunctive Normal Form is a conjunction of clauses, where clauses are disjunctions of literals being atoms (positive literals) or negated atoms (negative literals).

Formulas can be represented as sparse polynomials, which will be discussed in more detail in Chapter 17 (see Definition 81).

Lemma 10. *Any term is representable by a single monomial and any clause is representable by at most two monomials.*

Proof. Let \mathcal{V}_0 and \mathcal{V}_1 be disjoint subsets of \mathcal{V} , then we have

$$Z_{\mathcal{V}_0, \mathcal{V}_1}^\wedge = e_{\{x_k=0:k \in \mathcal{V}_0\} \cup \{x_k=1:k \in \mathcal{V}_1\}} [X_{\mathcal{V}_0 \cup \mathcal{V}_1}] \otimes \mathbb{I} [X_{\mathcal{V}/(\mathcal{V}_0 \cup \mathcal{V}_1)}]$$

and

$$Z_{\mathcal{V}_0, \mathcal{V}_1}^\vee = \mathbb{I} [X_{\mathcal{V}}] - e_{\{x_k=0:k \in \mathcal{V}_0\} \cup \{x_k=1:k \in \mathcal{V}_1\}} [X_{\mathcal{V}_0 \cup \mathcal{V}_1}] \otimes \mathbb{I} [X_{\mathcal{V}/(\mathcal{V}_0 \cup \mathcal{V}_1)}] .$$

We notice, that any tensors \mathbb{I} and $e_x \otimes \mathbb{I}$ 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. □

We apply Lemma 10 to show the following sparsity bound on the energy tensor of Markov Logic Networks.

Theorem 52. *Any formula f with a conjunctive normal form of n clauses satisfies*

$$\text{rank}^d(f) \leq 2^n .$$

For any set \mathcal{F} of formulas each with a conjunctive normal form of n_f clauses satisfies for any θ

$$\text{rank}^d \left(\sum_{f \in \mathcal{F}} \theta_f \cdot f \right) \leq \sum_{f \in \mathcal{F}} 2^{n_f} .$$

Proof. Let f have a CNF 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 Theorem 113 and Lemma 10 to get

$$\text{rank}^d(f) \leq \prod_{l \in [n]} \text{rank}^d(Z_{\mathcal{V}_0^l, \mathcal{V}_1^l}^\vee) \leq 2^n.$$

Given a collection of formulas \mathcal{F} , each with a CNF of n_f clauses we apply Theorem 112 and get

$$\text{rank}^d\left(\sum_{f \in \mathcal{F}} \theta_f \cdot f\right) \leq \sum_{f \in \mathcal{F}} \text{rank}^d(f) \leq \sum_{f \in \mathcal{F}} 2^{n_f}.$$

□

8.2.4 Categorical Constraints

We 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.

To capture the constraints resulting from this construction we introduce auxiliary parts. Such constraints can also be expressed by a formula but would result in an unnecessary large tensor network.

Definition 55 (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*

$$Z(x, x_A) = \begin{cases} 1 & \text{if } x_{[m]} = e_x \\ 0 & \text{else.} \end{cases}$$

We then call the variables X_0, \dots, X_{m-1} the atomization variables to the categorical variable X .

With Theorem 99 the tensor representation of ρ^Z decomposes in a basis CP format (see Figure 22b) of if its coordinate maps Z_x , where $x \in [m]$. For the cores

$$\rho^{Z_x} = e_x[X] \otimes e_1[X_x] + (\mathbb{I}[X] - e_x[X]) \otimes e_0[X_x] \quad (35)$$

we have by Theorem 99

$$\rho^Z[X, X_0, \dots, X_{m-1}] = \left\langle \{\rho^{Z(x)} : x \in [m]\} \right\rangle [X, X_0, \dots, X_{m-1}].$$

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 53. *For any tensor $T[X_{[d]}]$ and a categorical constraint defined by an ordered subset $X_A \subset X_{[d]}$, a variable $X \in X_{[d]}$ we have*

$$\langle \{T[X_{[d]}], Z\} \rangle [X_0 = x_0, \dots, X_{d-1} = x_{d-1}] = \begin{cases} T[X_0 = x_0, \dots, X_{d-1} = x_{d-1}] & \text{if } x_A = e_x \\ 0 & \text{else.} \end{cases}$$

Proof. For any $x_{[d]}$ we have

$$\langle \{T[X_{[d]}], Z\} \rangle [X_0 = x_0, \dots, X_{d-1} = x_{d-1}] = T[X_{[d]} = x_{[d]}] \cdot Z[X = x, X_A = x_A].$$

If $x_A = e_x$ we have $Z[X = x, X_A = x_A] = 1$ and thus

$$\langle \{T[X_{[d]}], Z\} \rangle [X_0 = x_0, \dots, X_{d-1} = x_{d-1}] = T[X_{[d]} = x_{[d]}].$$

If $x_A \neq e_x$ then $Z[X = x, X_A = x_A] = 0$ and

$$\langle \{T[X_{[d]}], Z\} \rangle [X_0 = x_0, \dots, X_{d-1} = x_{d-1}] = 0.$$

□

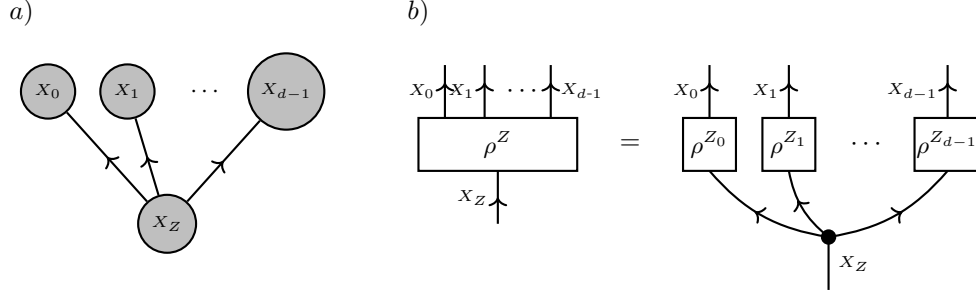


Figure 22: Representation of a categorical constraint in a CP Format tensor network. a) Representation of the dependency of the graphical model. b) Tensor Representation with further network decomposition. We average by contraction with the dashed tensor \mathbb{I} , if we do not specify the active atom.

Remark 22 (Combination of Constraints). *We can combine constraint cores by Hadamard products in the dual tensor network representation, as long as they can be satisfied together. An example, where this is not the case, are the categorical constraints to the three sets*

$$\{X_0, X_1, X_2, X_3\}, \{X_0, X_1\}, \{X_2, X_3\}.$$

Besides the categorical cores also the datacores have a similar bayesian network affecting the atoms by another hidden variable. Combining both is welldefined, only when all datapoints satisfy the categorical constraints (that is only one of the atoms in each constraint is active).

Example 11 (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 ?). 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 11 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.

Reasoning by Entailment propagation! Also, probabilistic choices possible when exact (!) contraction at a position not a basis vector; then can choose one possibility.

8.3 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).

$X_{0,0}$	$X_{0,1}$	$X_{0,2}$	$X_{0,3}$	$X_{0,4}$	$X_{0,5}$	$X_{0,6}$	$X_{0,7}$	$X_{0,8}$
$X_{1,0}$	$X_{1,1}$	$X_{1,2}$	$X_{1,3}$	$X_{1,4}$	$X_{1,5}$	$X_{1,6}$	$X_{1,7}$	$X_{1,8}$
$X_{2,0}$	$X_{2,1}$	$X_{2,2}$	$X_{2,3}$	$X_{2,4}$	$X_{2,5}$	$X_{2,6}$	$X_{2,7}$	$X_{2,8}$
$X_{3,0}$	$X_{3,1}$	$X_{3,2}$	$X_{3,3}$	$X_{3,4}$	$X_{3,5}$	$X_{3,6}$	$X_{3,7}$	$X_{3,8}$
$X_{4,0}$	$X_{4,1}$	$X_{4,2}$	$X_{4,3}$	$X_{4,4}$	$X_{4,5}$	$X_{4,6}$	$X_{4,7}$	$X_{4,8}$
$X_{5,0}$	$X_{5,1}$	$X_{5,2}$	$X_{5,3}$	$X_{5,4}$	$X_{5,5}$	$X_{5,6}$	$X_{5,7}$	$X_{5,8}$
$X_{6,0}$	$X_{6,1}$	$X_{6,2}$	$X_{6,3}$	$X_{6,4}$	$X_{6,5}$	$X_{6,6}$	$X_{6,7}$	$X_{6,8}$
$X_{7,0}$	$X_{7,1}$	$X_{7,2}$	$X_{7,3}$	$X_{7,4}$	$X_{7,5}$	$X_{7,6}$	$X_{7,7}$	$X_{7,8}$
$X_{8,0}$	$X_{8,1}$	$X_{8,2}$	$X_{8,3}$	$X_{8,4}$	$X_{8,5}$	$X_{8,6}$	$X_{8,7}$	$X_{8,8}$

Figure 23: 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.

Definition 56. Given a set of formulas \mathcal{F} with weights θ and set \mathcal{KB} of formulas, which conjunction is satisfiable, the hybrid logic network is the probability distribution

$$\mathbb{P}^{(\mathcal{F}, \theta, \nu^{\mathcal{KB}})}[X_{[d]}] = \langle \{f : f \in \mathcal{KB}\} \cup \{\exp[\theta_f \cdot f] : f \in \mathcal{F}\} \rangle [X_{[d]} | \emptyset] ,$$

which is the member of the exponential family with statistic by \mathcal{F} and the base measure

$$\nu^{\mathcal{KB}}[X_{[d]}] = \langle \{f : f \in \mathcal{KB}\} \rangle [X_{[d]}] .$$

Given a set of formulas \mathcal{F} , we define the set of hybrid logic networks realizable with \mathcal{F} as

$$\Lambda^{\mathcal{F}} = \bigcup_{\tilde{\mathcal{F}} \subset \mathcal{F}, \mu \in \{0,1\}^{|\mathcal{F}|}} \Gamma^{\mathcal{F}/\tilde{\mathcal{F}}, \nu^{\tilde{\mathcal{F}}}, \mu}$$

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 54. If any only if $\bigwedge_{f \in \mathcal{KB}} f$ is satisfiable, the tensor

$$\langle \{f : f \in \mathcal{KB}\} \cup \{\exp[\theta_f \cdot f] : f \in \mathcal{F}\} \rangle [X_{[d]}]$$

is normable.

Proof. We need to show that

$$\langle \{f : f \in \mathcal{KB}\} \cup \{\exp[\theta_f \cdot f] : f \in \mathcal{F}\} \rangle [\emptyset] > 0 . \quad (36)$$

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

$$\begin{aligned} \langle \{f : f \in \mathcal{KB}\} \cup \{\exp[\theta_f \cdot f] : f \in \mathcal{F}\} \rangle [X_{[d]} = x_{[d]}] &= \left(\prod_{f \in \mathcal{KB}} f[X_{[d]} = x_{[d]}] \right) \cdot \left(\prod_{f \in \mathcal{F}} \exp[\theta_f \cdot f][X_{[d]} = x_{[d]}] \right) \\ &= \left(\prod_{f \in \mathcal{F}} \exp[\theta_f \cdot f][X_{[d]} = x_{[d]}] \right) \\ &> 0. \end{aligned}$$

Condition (36) follows from this and the Hybrid Logic Network is well-defined. \square

8.3.1 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 55. *For any hybrid logic network we have*

$$\mathbb{P}^{(\mathcal{F}, \theta, \mathcal{KB})}[X_{[d]}] = \langle \{\rho^f[X_f, X_{[d]}] : f \in \mathcal{F} \cup \mathcal{KB}\} \cup \{e_1[X_f] : f \in \mathcal{KB}\} \cup \{W^f[X_f] : f \in \mathcal{F}\} \rangle [X_{[d]} | \emptyset] .$$

Proof. By Lemma 5. \square

While the statistics computing cores in the relational 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 47)

$$W^f[X_f] = \left[\exp \left[\theta \left[\frac{1}{L} \right] \right] \right] [X_{f_l}]$$

the hard formulas get activation cores by unit vectors

$$e_1[X_f] = \begin{bmatrix} 0 \\ 1 \end{bmatrix} [X_{f_l}] .$$

As shown in Section 8.2.1, the soft activation cores converge to these hard activation cores in the limit of large parameters, when imposing a local normation. 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 24.

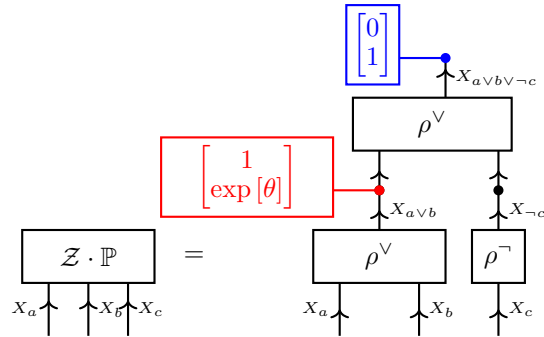


Figure 24: Diagram of a formula tensor with activated heads, containing **hard constraint cores** and **probabilistic weight cores**.

Remark 23. *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.

8.3.2 Reasoning Properties

Deciding probabilistic entailment (see Definition 45) with respect to Hybrid Logic Networks can be reduced to the hard logic parts of the network.

Theorem 56. *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 40 on the representation of Hybrid Logic Networks as Markov Networks in Theorem 55. \square

Formulas in \mathcal{F} , which are entailed or contradicted by \mathcal{KB} are redundant, as we show next.

Theorem 57. *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 θ , 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. \square

A similar statement holds for the hard formulas itself, as shown in Theorem 35. 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 6, which has been introduced in Chapter 6. By Theorem 56 only the hard logic parts of a Hybrid Logic Network are required in the ASK operation.

8.3.3 Expressivity

Hybrid Logic Networks extend the expressivity result of Theorem 48 to arbitrary probability tensors, dropping the positivity constraints for Markov Logic Networks.

Theorem 58. *Let $\mathbb{P} [X_{[d]}]$ a possibly not positive probability tensor we build a base measure*

$$\nu^{\mathcal{KB}} = \chi (\mathbb{P} [X_{[d]}])$$

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 θ , that is

$$\mathbb{P} [(\mathcal{F}_{\wedge}, \theta, \nu^{\mathcal{KB}})]$$

Proof. It suffices to show that

$$\langle \nu^{\mathcal{KB}}, \exp [\langle \gamma^{\mathcal{F}_{\wedge}} \theta \rangle [X_{[d]}]] \rangle [X_{[d]}] = \mathbb{P} [X_{[d]}].$$

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

$$\langle \nu^{\mathcal{KB}}, \exp [\langle \gamma^{\mathcal{F}_{\wedge}} \theta \rangle [X_{[d]}]] \rangle [X_{[d]} = x_{[d]}] = 0.$$

For indices $x_{[d]}$ with $\mathbb{P} [X_{[d]} = x_{[d]}] > 0$ we have $\nu^{\mathcal{KB}} [X_{[d]} = x_{[d]}] = 1$ and

$$\begin{aligned} \langle \nu^{\mathcal{KB}}, \exp [\langle \gamma^{\mathcal{F}_{\wedge}} \theta \rangle [X_{[d]}]] \rangle [X_{[d]} = x_{[d]}] &= \prod_{l_{[d]}} \exp [\theta [L_{[d]} = l_{[d]}] \cdot Z_{l_{[d}}}^{\wedge} [X_{[d]} = x_{[d]}]] \\ &= \exp [\theta [L_{[d]} = x_{[d]}]] \\ &= \mathbb{P} [X_{[d]} = x_{[d]}]. \end{aligned}$$

\square

8.4 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 48, 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 7 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 relational 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 ?). Solving such problems is then equivalent to sampling from the worlds in a logical interpretation, and can be approached by the methods of Chapter 6. Among these classed, we have only discussed the Sudoku game in Example 11. Extensions by Hybrid Logic Networks can be interpreted as implementations of preferences among possible solutions by probabilities.

9 Logic Network Inference

In this chapter we investigate the inference properties of Hybrid Logic Networks starting with characterizations of its mean parameter polytopes. We first investigate unconstrained parameter estimated for Markov Logic Networks and Hybrid Logic Networks, which are special cases of the backward maps introduced in Chapter 3. We then motivate structure learning based on sparsity constraints on the parameters on the minterm family and present heuristic strategies leading to efficient structure learning algorithms.

9.1 Mean parameters of Hybrid Logic Networks

The convex polytope of realizable mean parameters (see Definition 30) is for a statistic \mathcal{F} of propositional formulas and a base measure ν

$$\mathcal{M}_{\mathcal{H},\nu} = \left\{ \langle \mathbb{P}, \gamma^{\mathcal{F}} \rangle [L] : \mathbb{P} \in \Gamma \right\},$$

where by Γ we denote the set of all probability distributions. By Theorem 13 the convex polytope has a characterization as a convex hull

$$\mathcal{M}_{\mathcal{H},\nu} = \text{conv} \left(\gamma^{\mathcal{F}} [X_{[d]} = x_{[d]}, L] : x_{[d]} \in \prod_{k \in [d]} [2], \nu [X_{[d]} = x_{[d]}] = 1 \right). \quad (37)$$

The mean parameter polytopes are examples of 0/1-polytopes ?Gillmann (2007), from which a few obvious properties follow. Since those are convex subsets of the cube $[0, 1]^p$, which extreme points are all binary vectors, also each $\gamma^{\mathcal{F}} [X_{[d]} = x_{[d]}, L]$ (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 $\mathcal{M}_{\mathcal{H},\nu}$

To summarize some insights:

- If all mean parameters in $\{0, 1\}$ then an extreme points, exactly those are hard logic networks
- If some mean params in $\{0, 1\}$, then not in the interior. Back direction not correct: There are interior points where no coordinate in $\{0, 1\}$.
- If in the interior after dropping the hard coordinates: Then a hybrid logic network

9.1.1 Extreme points by hard logic networks

We exploit this characterization to show, that the extreme points of $\mathcal{M}_{\mathcal{H},\nu}$ are exactly those realizable by Hard Logic Networks.

Theorem 59. *Any parameter $\mu [L] \in \mathcal{M}_{\mathcal{H},\nu}$ is an extreme point of $\mathcal{M}_{\mathcal{H},\nu}$, if and only $\mu [L]$ 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, μ 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 μ be an extreme point of $\mathcal{M}_{\mathcal{H},\nu}$. Since by (37) $\mathcal{M}_{\mathcal{H},\nu}$ is a convex hull of vectors, there exists a $x_{[d]} \in \times_{k \in [d]} [2]$ such that

$$\mu [L] = \gamma^{\mathcal{F}} [X_{[d]} = x_{[d]}, L] .$$

By definition of $\gamma^{\mathcal{F}}$, $\mu [L]$ 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 $l \in [p]$

$$\gamma^{\mathcal{F}} [X_{[d]} = x_{[d]}, L = l] = f_l [X_{[d]} = x_{[d]}] = \mu [L = l]$$

and thus $\mu [L] = \gamma^{\mathcal{F}} [X_{[d]} = x_{[d]}, L]$. With the characterization (37) this establishes in particular $\mu [L] \in \mathcal{M}_{\mathcal{H},\nu}$. 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 $\mathcal{M}_{\mathcal{H},\nu} \subset [0, 1]^p$. \square

9.1.2 Mean parameters in the interior

By Theorem 20 the interior points are those realizable by a Hybrid Logic Network with statistics \mathcal{H} and base measure ν , as we state in the following Corollary.

Corollary 7. *If $\mu [L] \in \mathcal{M}_{\mathcal{H},\nu}^\circ$, or equivalently the statistic is minimal and $\mu [L]$ is reproduceable by a distribution positive with respect to ν , then there is $\theta [L]$ such that $\mathbb{P}^{\mathcal{F},\theta,\nu}$ reproduces $\mu [L]$.*

9.1.3 Mean parameters outside the interior

By Theorem 21 mean parameter vectors outside the interior of $\mathcal{M}_{\mathcal{H},\nu}$ are not realizable by distributions, which are positive with respect to the base measure ν . 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 $\mu [L] \notin \mathcal{M}_{\mathcal{H},\nu}^\circ$. In this case, the next theorem provides us with a procedure to refine the base measure in these cases.

Theorem 60 (Base measure refinement for mean coordinates in $\{0, 1\}$). *To any $\mu[L] \in \mathcal{M}_{\mathcal{H}, \nu}$ 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 ν and reproducing μ , that is

$$\langle \mathbb{P}, \gamma^{\mathcal{F}} \rangle [L] = \mu[L] .$$

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 Theorem 37 probabilistic entailment $\mathbb{P} \models f_l$ (see Definition 45). 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 Theorem 60 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 $\mu[L] \notin \mathcal{M}_{\mathcal{H}, \nu}^\circ$, see Example 12 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 Definition 39), which has been shown in more generality in Theorem 22. For hybrid logic network we characterize the face base measures in the next theorem.

Theorem 61 (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*

$$f_{\mathcal{F}, \theta} = \bigwedge_{l \in [p] : \theta[L=l] \neq 0} \neg^{(1-\psi(\theta[L=l]))} f_l$$

is satisfiable, then it is the base measure to the face with normal θ . Here $\psi(z)$ 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

$$\operatorname{argmax}_{x_{[d]}} \langle \theta[L], \gamma^{\mathcal{F}}[X_{[d]} = x_{[d]}, L] \rangle [\emptyset] .$$

For the first claim: We have since $\gamma^{\mathcal{F}}$ is boolean that

$$\|\theta[L]\|_1 \geq \max_{x_{[d]}} \langle \theta[L], \gamma^{\mathcal{F}}[X_{[d]} = x_{[d]}, L] \rangle [\emptyset] .$$

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 θ we have $\gamma^{\mathcal{F}}[X_{[d]} = x_{[d]}, L = l] = \psi(\theta[L = l])$. This condition is equal to $x_{[d]}$ being a model of $\neg^{(1-\psi(\theta[L=l]))} f_l$ for any l in the support of θ 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 $\neg^{(1-\psi(\theta[L=l]))} f_l$ being satisfied. In the worlds, for which $\max_{x[d]} \langle \theta[L], \gamma^{\mathcal{F}}[X[d] = x[d], L] \rangle [\emptyset]$ 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

$$\bigvee_{v[L] : \langle \theta, v \rangle [\emptyset] \in \max_{\mu} \langle \theta, \mu \rangle [\emptyset]} f_{\mathcal{F}, \langle \theta, v \rangle [L]} .$$

□

9.1.4 Expressivity of Hybrid Logic Networks

Let us recall, that the set $\Lambda^{\mathcal{F}}$ 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 $\Lambda^{\mathcal{F}}$ (see Definition 56), that is

$$\mathcal{M}_{\mathcal{H}, \nu} |_{\Lambda^{\mathcal{F}}} = \{ \langle \mathbb{P}, \gamma^{\mathcal{F}} \rangle [L] : \mathbb{P} \in \Lambda^{\mathcal{F}} \} .$$

While $\mathcal{M}_{\mathcal{H}, \nu} |_{\Lambda^{\mathcal{F}}} \subset \mathcal{M}_{\mathcal{H}, \nu}$ 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 $\Lambda^{\mathcal{F}}$. In the next example we provide a class of formulas, for which $\Lambda^{\mathcal{F}}$ does not have sufficient expressivity.

Example 12 (Insufficient expressivity of $\Lambda^{\mathcal{F}}$ 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 $\Lambda^{\mathcal{F}}$ has sufficient expressivity, let us first prove a generic criterion.

Theorem 62. *The set of mean parameters realizable by $\Lambda^{\mathcal{F}}$ coincides with the set of mean parameters realizable by any distribution, if and only if for any boolean vector $v[L]$ the formula $f_{\mathcal{F}, \theta}$ is satisfiable.*

Proof. " \Leftarrow ": If these formulas are satisfiable, the base measures to the facets coincides with those realizable by $\Lambda^{\mathcal{F}}$. " \Rightarrow ": We show that implication by refutation. If one of these formulas are not satisfiable, we find a facet, which mean parameters optimize. The uniform distribution on the worlds taking the maximum is not representable by any base measure. □

Example 13 (Sufficient expressivity in cases of atomic formulas). *The assumption of Theorem 62 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.*

When the assumptions of Theorem 62 are not satisfied, there are mean parameters, which can not be reproduced by a distribution in $\Lambda^{\mathcal{F}}$. In that case, we can flexibilize the distribution, to also represent the base measures used for refinement in Theorem 61. This can be done by adding activation cores with multiple variables, or further computation cores calculating the disjunctions of formulas.

To summarize, for any $\mu \in \mathcal{M}_{\mathcal{H}, \nu}$ we have one of the following (see also Figure 9.1.4):

- $\mu[L = l] \in \mathcal{M}_{\mathcal{H}, \nu}^{\circ}$: Then reproducible by a Hybrid Logic Network with base measure ν .
- $\mu[L = l] \notin \mathcal{M}_{\mathcal{H}, \nu}^{\circ}$: Then on a facet and reproducible by a Hybrid Logic Network with refined base measure $\tilde{\nu}$. Whether the base measure can be realized by the computation network of \mathcal{H} depends on the satisfiability of the face formula.

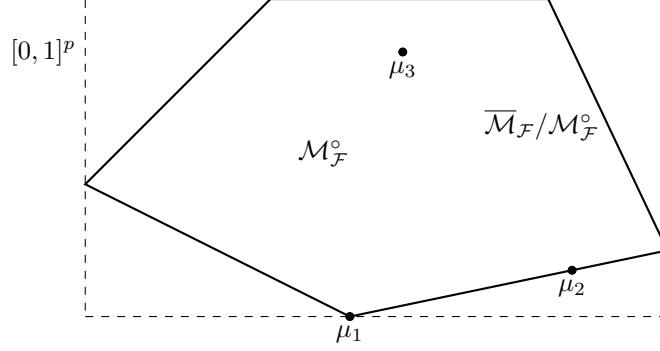


Figure 25: Sketch of the convex polytope $\mathcal{M}_{\mathcal{H},\nu}$ as a subset of the p -dimensional cube $[0, 1]^p$ (here as a 2-dimensional projection) with example mean parameters μ_1, μ_2, μ_3 . The boundary points $\mu_1, \mu_2 \in \overline{\mathcal{M}_{\mathcal{F}}} / \mathcal{M}_{\mathcal{F}}^\circ$ are examples of mean parameters, which can be realized by a Hard Logic Networks (respectively Hybrid Logic Network). Any extreme point $\mu_1 \in \mathcal{M}_{\mathcal{H},\nu} \cup \{0, 1\}^p$ is realizable by a Hard Logic Network, while a non-extreme boundary point $\mu_2 \in \mathcal{M}_{\mathcal{H},\nu} / \{0, 1\}^p$ is realizable by a Hybrid Logic Network. Any interior point $\mu_3 \in \mathcal{M}_{\mathcal{H},\nu}^\circ$ is realizable by a Markov Logic Network.

9.1.5 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.

9.2 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 4 for the generic exponential families.

9.2.1 Maximum Likelihood in Hybrid Logic Networks

The Maximum Likelihood Problem on Markov Logic Networks is the M-projection

$$\operatorname{argmax}_{\theta[L] \in \mathbb{R}^p} \quad \mathbb{H} \left[\mathbb{P}, \mathbb{P}^{(\phi, \theta, \nu)} \right]$$

in the case $\mathbb{P} = \mathbb{P}^D$ for a data map D .

The M-projection coincides, after dropping constant terms in case of non-trivial base measure, with the backward map

$$\operatorname{argmax}_{\theta[L] \in \mathbb{R}^p} \quad \langle \theta[L], \mu[L] \rangle [\emptyset] - A^{(\phi, \nu)}(\theta[L])$$

where

$$\mu[L] = \langle \gamma^{\mathcal{F}}, \mathbb{P} \rangle [L] \quad \text{and} \quad A^{(\phi, \nu)}(\theta[L]) = \langle \exp[\langle \gamma^{\mathcal{F}}[X_{[d]}, L], \theta[L] \rangle [X_{[d]}], \nu \rangle [\emptyset].$$

We now extend to Hybrid Logic Networks

$$\operatorname{argmax}_{\tilde{\mathbb{P}} \in \Lambda^{\mathcal{F}}} \quad \mathbb{H} \left[\mathbb{P}, \tilde{\mathbb{P}} \right]$$

Corollary 8. Let $\mu[L] = \langle \mathbb{P}, \gamma^{\mathcal{F}} \rangle [L]$ 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 reproducible 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 $\Gamma^{\mathcal{F}/\tilde{\mathcal{F}}, \nu^{\tilde{\mathcal{F}}, \mu}}$.

9.2.2 Maximum Entropy in Hybrid Logic Networks

The Maximum Entropy Problem for Markov Logic Networks is

$$\operatorname{argmax}_{\mathbb{P}} \mathbb{H}[\mathbb{P}] \quad \text{subject to} \quad \langle \mathbb{P}, \gamma^{\mathcal{F}} \rangle [L] = \mu [L] \quad (38)$$

Corollary 9 (of Theorem 30). *Let \mathbb{P}^D be a distribution such that there is a positive distribution \mathbb{P} with $\langle \mathbb{P}, \gamma^{\mathcal{F}} \rangle [L] = \langle \mathbb{P}^D, \gamma^{\mathcal{F}} \rangle [L]$. 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 θ 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 63. *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*

- $\mathcal{KB} = \{f_l : l \in [p], \mu [L = l] = 1\} \cup \{\neg f_l : l \in [p], \mu [L = l] = 0\}$
- $\mathcal{F} = \{f_l : l \in [p], \mu [L = l] \in (0, 1)\}$
- θ being the backward map evaluated at the vector μ 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

9.3 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 Lemma 4 in the following.

Lemma 11. *Let there be a base measure ν , a formula selecting map $\mathcal{F} = \{f_l : l \in [p]\}$ and weights θ , and choose $l \in [p]$ such that $f_l \notin \{\mathbb{I}[X_{[d]}], 0[X_{[d]}]\}$. The moment matching condition relative to θ , $l \in [p]$ and $\mu_D [L = l] \in (0, 1)$ is then satisfied, if*

$$\theta [L = l] = \ln \left[\frac{\mu_D [L = l]}{(1 - \mu_D [L = l])} \cdot \frac{T [X_{f_l} = 0]}{T [X_{f_l} = 1]} \right] \quad (39)$$

where by $T [X_{f_l}]$ we denote the contraction

$$T [X_{f_l}] = \left\langle \{\rho^{f_l} : l \in [p]\} \cup \{W^{\tilde{l}} : \tilde{l} \in [p], \tilde{l} \neq l\} \cup \{\nu\} \right\rangle [X_{f_l}].$$

Proof. Since $\operatorname{im}(f_l) \subset [2]$ we have

$$\operatorname{Id}|_{\operatorname{im}(f_l)} = e_1 [X_{f_l}]$$

and the moment matching condition is by Lemma 4 satisfied if

$$\langle W^l, e_1, T \rangle [\emptyset] = \langle W^l, T \rangle [\emptyset] \cdot \mu_D [L = l].$$

This is equal to

$$\exp [\theta [L = l]] \cdot T [X_{f_l} = 1] = (\exp [\theta [L = l]] \cdot T [X_{f_l} = 1] + T [X_{f_l} = 0]) \cdot \mu_D [L = l].$$

Rearranging the equations this is equal to

$$T [X_{f_l}] = \left\langle \{\rho^{f_l}\} \cup \{W^{\tilde{l}} : \tilde{l} \in [p], \tilde{l} \neq l\} \cup \{\nu\} \right\rangle [L].$$

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 $T [X_{f_l} = 0], T [X_{f_l} = 1] \neq 0$ 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 Section 8.2 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 $T [X_{f_l} = 1] = 0$, respectively $T [X_{f_l} = 1] = 1$ 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 will have been investigated in Section 8.2.

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 8, 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 $\mathcal{L}_D (\mathbb{P}^{(\phi, \theta, \nu)})$ as an objective depending on the vector θ .

Algorithm 8 Alternating Weight Optimization (AWO)

$\mathcal{KB} = \mathbb{I}, \tilde{\mathcal{V}} = \emptyset$

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

$$T [X_{f_l}] \leftarrow \langle \rho^{f_l} \rangle [X_{f_l}]$$

 Set

$$\theta [L = l] \leftarrow \ln \left[\frac{\mu_D [L = l]}{(1 - \mu_D [L = l])} \cdot \frac{T [X_{f_l} = 0]}{T [X_{f_l} = 1]} \right]$$

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

$$T [X_{f_l}] = \left\langle \{\rho^{f_l} : l \in [p]\} \cup \{W^{\tilde{l}} : \tilde{l} \in [p], \tilde{l} \neq l\} \cup \{\nu\} \right\rangle [X_{f_l}]$$

 Set

$$\theta [L = l] = \ln \left[\frac{\mu_D [L = l]}{(1 - \mu_D [L = l])} \cdot \frac{T [X_{f_l} = 0]}{T [X_{f_l} = 1]} \right]$$

end for

end while

In the initialization phase of Algorithm 8, 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 39 requires inference of a current model by answering the query to the formula \tilde{f} . This can be a bottleneck and circumvented by approximative inference, see e.g. CAMEL ?.

Remark 24 (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 θ^t we compute*

$$\mu^t = \left\langle \mathbb{P}^{(\phi, \theta^t)}, \phi \right\rangle [L]$$

and build the update

$$\theta^{t+1} = \theta^t + \ln [\mu^D] \mu^t.$$

Then, θ^{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 $\neg f$.

9.4 Forward and backward mappings in closed form

We recall from Chapter 4, 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.

9.4.1 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 Section 8.1.4 how to represent any distribution as a MLN of maxterms (respectively minterms), see Theorem 48.

We use the tuple enumeration of the maxterms and minterms by $\times_{k \in [d]} [2]$ introduced in Section 5.3.3. 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 64. *Given the Markov Logic Networks to the formula sets*

$$\mathcal{F}_\wedge := \{Z_{x_0, \dots, x_{d-1}}^\wedge : x_0, \dots, x_{d-1} \in \times_{k \in [d]} [2]\} \quad \text{and} \quad \mathcal{F}_\vee := \{Z_{x_0, \dots, x_{d-1}}^\vee : x_0, \dots, x_{d-1} \in \times_{k \in [d]} [2]\}$$

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 θ .

Possible choices of the backward mappings are

$$B^\wedge(\mu) = \ln [\mu] \quad \text{and} \quad B^{\mathcal{F}_\vee}(\mu) = -\ln [\mu].$$

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 μ 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

$$\begin{aligned} F^\vee(\theta) &= \langle \exp[\langle \{\mathcal{F}_\wedge, \theta\} \rangle [X_{[d]}]] \rangle [X_{[d]}|\emptyset] \\ &= \langle \{ \exp[\langle \{\mathbb{I}, \theta\} \rangle [X_{[d]}]] , \exp[-\langle \theta \rangle [X_{[d]}]] \} \rangle [X_{[d]}|\emptyset] \\ &= \langle \exp[-\theta] \rangle [X_{[d]}|\emptyset] \end{aligned}$$

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 normation. 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.

9.4.2 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 65. *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

$$\begin{aligned} F^{[d]}(\theta[L])[L = k] &= \langle X_k, \mathbb{P}^{(\mathcal{F}_{[d]}, \theta)}[X_{[d]}] \rangle [\emptyset] \\ &= \langle X_k, \langle \exp[\theta[L = k] \cdot X_k] \rangle [X_k|\emptyset] \rangle [\emptyset] \\ &= \frac{\exp[\theta[L = k]]}{1 + \exp[\theta[L = k]]}. \end{aligned}$$

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 μ in the image of the forward map

$$F^{[d]}(B^{[d]}(\mu)) = \mu$$

and is therefore a backward map. \square

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 17, 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 25 (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.*

9.5 Constrained parameter estimation in the minterm family

We approach structure learning as constrained parameter estimation in the naive exponential family (see Example 4), which coincides with the minterm family \mathcal{F}_\wedge . The minterm family is defined by the statistic $\phi = \delta [X_{[d]}, L_{[d]}]$ 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 Theorem 48 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 Theorem 58.

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

$$\operatorname{argmax}_{\theta \in \bigotimes_{k \in [d]} \mathbb{R}^{m_k}} \langle \theta, \mathbb{P}^D \rangle [\phi] - A^{(\phi, \nu)}(\theta)$$

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 Theorem 58 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

$$\operatorname{argmax}_{\theta \in \Theta} \langle \theta, \mathbb{P}^D \rangle [\phi] - A^{(\phi, \nu)}(\theta). \quad (\text{P}_{\Theta, \mathbb{P}^D})$$

Problem ?? has two important types of instantiation, which we discuss in the next sections.

9.5.1 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 Section 9.2. 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 26 (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.

9.5.2 Structure Learning

The problem of structure learning arises, when the set of parameters in Problem ?? is chosen as

$$\Theta^{\mathcal{H}} = \bigcup_{\mathcal{F} \in \mathcal{H}} \text{span}(\mathcal{F}) .$$

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 $f \in \mathcal{F}$.

9.6 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 7 and select from them a small subset.

9.6.1 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

$$\text{argmin}_{f \in \mathcal{H}} \text{argmin}_{\theta \in \mathbb{R}^{|\mathcal{F}|+1}} \mathbb{H} \left[\mathbb{P}^D, \mathbb{P}^{(\mathcal{F} \cup \{f\}, \theta, \nu)} \right] . \quad (\mathbb{P}_{D, \mathcal{F}, \mathcal{H}})$$

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)).

9.6.2 Gain Heuristic

In the gain heuristic, only the parameters of the new formula are optimized and the others left unchanged. This amounts to

$$\text{argmin}_{f \in \mathcal{H}} \left(\min_{\theta_{[|\mathcal{F}|]} \in \mathbb{R}} \mathbb{H} \left[\mathbb{P}^D, \mathbb{P}^{(\mathcal{F} \cup \{f\}, \theta, \nu)} \right] \right) . \quad (\mathbb{P}_{D, \mathcal{F}, \mathcal{H}}^{\text{gain}})$$

Here we denote by θ 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 $\theta_{[|\mathcal{F}|]}$.

Lemma 12. *The gain heuristic objective is an upper bound on the true greedy objective.*

Proof. Since

$$\text{argmin}_{f \in \mathcal{H}} \left(\text{argmin}_{\theta \in \mathbb{R}^{|\mathcal{F}|+1}} \mathbb{H} \left[\mathbb{P}^D, \mathbb{P}^{(\mathcal{F} \cup \{f\}, \theta, \nu)} \right] \right) \leq \text{argmin}_{f \in \mathcal{H}} \left(\text{argmin}_{\theta_{[|\mathcal{F}|]} \in \mathbb{R}} \mathbb{H} \left[\mathbb{P}^D, \mathbb{P}^{(\mathcal{F} \cup \{f\}, \theta, \nu)} \right] \right) .$$

□

Further, this is Problem $(\mathbb{P}_{\Theta, \mathbb{P}^D})$ in the case

$$\Theta = \ln \left[\tilde{\mathbb{P}} \right] + \cup_{f \in \mathcal{F}} \text{span}(f) .$$

Let us choose a formula $f \in \mathcal{F}$ and consider Problem $\mathbb{P}_{\Theta, \mathbb{P}^D}$ in the case

$$\Theta^f = \ln \left[\tilde{\mathbb{P}} \right] + \text{span}(f) .$$

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 4 and characterize the solution by the θ 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 66. Problem $(\mathbb{P}_{D,\mathcal{F},\mathcal{H}}^{\text{gain}})$ is solved at any

$$\hat{\theta} = \theta_{\hat{f}} \cdot \hat{f}$$

where the formula \hat{f} is in

$$\hat{f} \in \operatorname{argmax}_{f \in \mathcal{F}} \mathbb{D}_{\text{KL}} \left[\langle \mathbb{P}^D, f \rangle [\emptyset] \parallel \langle \tilde{\mathbb{P}}, f \rangle [\emptyset] \right]$$

and $\theta_{\hat{f}}$ is the weight of \hat{f} in the solution of Problem $\mathbb{P}_{\Theta, \mathbb{P}^D}$ with $\Gamma = \tilde{\mathbb{P}} + \operatorname{span}(f)$. Here we denote by $\mathbb{D}_{\text{KL}} [p_1 || p_2]$ the Kullback-Leibler divergence between Bernoulli distributions with parameters $p_1, p_2 \in [0, 1]$, that is

$$\mathbb{D}_{\text{KL}} [p_1 || p_2] = p_1 \cdot \ln \left[\frac{p_1}{p_2} \right] + (1 - p_1) \cdot \ln \left[\frac{(1 - p_1)}{(1 - p_2)} \right]$$

Proof. For any formula f , the inner minimum of Problem $(\mathbb{P}_{D,\mathcal{F},\mathcal{H}}^{\text{gain}})$ is by Lemma 11 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 55 and get

$$\begin{aligned} \langle \tilde{\mathbb{P}}, \exp [\theta_f \cdot f] \rangle [\emptyset] &= \langle \tilde{\mathbb{P}}, \rho^f [X_f], W^f [X_f] \rangle [\emptyset] \\ &= (1 - \tilde{\mu}) + \tilde{\mu} \cdot \exp [\theta_f] \\ &= (1 - \tilde{\mu}) + \frac{\mu_D \cdot (1 - \tilde{\mu})}{(1 - \mu_D)} \\ &= (1 - \tilde{\mu}) \cdot \frac{1}{(1 - \mu_D)} . \end{aligned}$$

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 μ_D and $\tilde{\mu}$. Since the gain in the likelihood loss when restricting to $\Theta = \operatorname{span}(f)$ is thus given by $\mathbb{D}_{\text{KL}} [\langle \mathbb{P}^D, f \rangle [\emptyset] || \langle \tilde{\mathbb{P}}, f \rangle [\emptyset]]$, we have that Problem ?? in the case $\Theta = \bigcup_{f \in \mathcal{F}} \operatorname{span}(f)$ is solved at $\hat{\theta} = \theta_{\hat{f}} \cdot \hat{f}$ where

$$\hat{f} = \mathbb{D}_{\text{KL}} [\langle \mathbb{P}^D, f \rangle [\emptyset] || \langle \tilde{\mathbb{P}}, f \rangle [\emptyset]] .$$

□

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.

9.6.3 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

$$\mathcal{L}_D \left(\mathbb{P}^{(\delta, \tilde{\theta})} \right) = \left\langle \mathbb{P}^D, \gamma^\delta, \tilde{\theta} \right\rangle [\emptyset] - \ln \left[\left\langle \exp \left[\left\langle \gamma^\delta, \tilde{\theta} \right\rangle [X_{[d]}] \right] \right\rangle [\emptyset] \right]$$

as

$$\begin{aligned} \nabla_{\tilde{\theta}[L]} \mathcal{L}_D \left(\mathbb{P}^{(\delta, \tilde{\theta})} \right) &= \left\langle \gamma^\delta, \mathbb{P}^D \right\rangle [L] - \left\langle \gamma^\delta, \mathbb{P}^{(\delta, \tilde{\theta})} \right\rangle [L] \\ &= \mathbb{P}^D - \mathbb{P}^{(\delta, \tilde{\theta})} . \end{aligned}$$

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

$$f(\theta) = \tilde{\theta} + \left\langle \theta, \gamma^{\mathcal{H}} \right\rangle [X_{[d]}] .$$

At $\theta = 0$ we have the gradient of the loss of the parametrized formula by

$$\begin{aligned} \nabla_{\theta|0} \mathcal{L}_D \left(\mathbb{P}^{(\delta, f(\theta), \nu)} \right) &= \left\langle \nabla_{f(\theta)}|_{\tilde{\theta}} \mathcal{L}_D \left(\mathbb{P}^{(\delta, f(\theta), \nu)} \right), \nabla_{\theta|0} f(\theta) \right\rangle [\emptyset] \\ &= \left\langle \mathbb{P}^D, \gamma^\phi \right\rangle [L] - \left\langle \mathbb{P}^{(\delta, \tilde{\theta}, \nu)}, \gamma^\phi \right\rangle [L] . \end{aligned}$$

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 μ_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.

Algorithm 9 Greedy Structure Learning

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_{Θ, \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 8.

$$\tilde{\mathbb{P}} \leftarrow \mathbb{P}^{\mathcal{F}, \theta}$$

end while**9.6.4 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 9, 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 θ 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 27 (Sparsification by Thresholding). *To maintain a small set of active formulas, one could combine greedy learning approaches with thresholding on the coordinates of θ . This is a standard procedure in Iterative Hard Thresholding algorithms of Compressed Sensing, but note that here we do not have a linear in θ objective.*

9.7 Proposal distribution

Let us now understand the likelihood gradient as the energy tensor of a probability distribution, which we call the proposal distribution.

Definition 57 (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 MAP query, for which we can apply the methods introduced in Chapter 4, such as Gibbs Sampling or Mean Field Approximations in combination with annealing.

9.7.1 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

$$\mathcal{M}_{\mathcal{H}^T} = \text{conv} \left(\gamma^{\mathcal{H}^T} L = l, X_{[d]} : l \in [p] \right) = \text{conv} \left(f[X_{[d]}] : f \in \mathcal{H} \right)$$

As it was the case for Markov Logic Networks, the mean parameter polytopes are instances of a 0/1-polytopes ?Gillmann (2007).

The extreme points are the formulas selectable by the formula selecting map \mathcal{H} .

9.8 Discussion

Remark 28 (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 θ 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}, θ) solving the maximum a posteriori problem*

$$\text{argmax}_{(\mathcal{F}, \theta)} \mathbb{P}^{(\mathcal{F}, \theta)}[\{D(j)\}_{j \in [m]}] \cdot \mathbb{P}[(\mathcal{F}, \theta)] . \quad (40)$$

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.

10 Probabilistic Success Guarantees

When drawing data independently from a random distribution, we in this chapter derive guarantees, that the learning algorithms in the Parameter Estimation and Structure Learning Chapter succeed.

10.1 Fluctuations

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.

10.1.1 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

$$e_{D(j)}[X_{[d]}] .$$

The expectation of this random tensor is

$$\mathbb{E}[e_{D(j)}] = \sum_{x_{[d]} \in \times_{k \in [d]} [m_k]} \mathbb{P}^*[X_{[d]} = x_{[d]}] e_{x_{[d]}}[X_{[d]}] = \mathbb{P}^*[X_{[d]}] .$$

The empirical distribution is then the average of independent random one-hot encodings, namely the random tensor

$$\mathbb{P}^D = \frac{1}{m} \sum_{j \in [m]} e_{D(j)}[X_{[d]}] .$$

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 Definition 60.

When the marginal of each datapoint is \mathbb{P}^* , the expectation of the empirical distribution is

$$\mathbb{E} [\mathbb{P}^D] = \frac{1}{m} \sum_{j \in [m]} \mathbb{E} [e_{D(j)}] = \mathbb{P}^* .$$

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.

10.1.2 Fluctuation tensors

Let us now investigate random tensors, which result from the forwarding of the fluctuation of the empirical distribution by sufficient statistics.

Definition 58. Given a statistic ϕ , $m \in \mathbb{N}$ and a dataset we define the fluctuation tensor as the random tensor

$$\eta^{\phi, \mathbb{P}^*, D} = \langle (\mathbb{P}^D - \mathbb{P}^*), \gamma^\phi \rangle [L]$$

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 naive 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 $\eta^{\phi, \mathbb{P}^*, D}$ is the centered mean parameter to the empirical distribution, that is

$$\mu_D - \mathbb{E} [\mu_D] = \langle \gamma^\phi, \mathbb{P}^D - \mathbb{P}^* \rangle [L] .$$

10.1.3 Widths of random tensors

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 59. Given a set $\Gamma \subset \bigotimes_{k \in [d]} \mathbb{R}^{m_k}$ and η a random tensor taking values in $\bigotimes_{k \in [d]} \mathbb{R}^{m_k}$ we define the width as the random variable

$$\omega_\Gamma(\eta) = \sup_{\theta \in \Gamma} |\langle \theta, \eta \rangle [\emptyset]| .$$

Bounds on the widths are also called uniform concentration bounds ? and generic probabilistic bounds will be provided in Chapter 19.

10.2 Fluctuations in Markov Logic Networks

For Markov Logic Networks we have statistics consistent of boolean features. In this case the marginal distributions of the coordinates of $\eta^{\phi, \mathbb{P}^*, D}$ are scaled and centered binomials, which we show now.

Lemma 13. Let ϕ be a statistic of boolean features ϕ_l for all $l \in [p]$, i.e. let $\text{im}(\phi_l) \subset \{0, 1\}$. Then, the marginal distribution of the coordinate $\eta^{\phi, \mathbb{P}^*, D}[L = l]$ is

$$\frac{1}{m} \left(B(p^{(l)}, m) - p^{(l)} \right) ,$$

where by $B(p^{(l)}, m)$ we denote the binomial distribution with mean parameter

$$p^{(l)} = \langle \phi_l, \mathbb{P}^* \rangle [\emptyset] .$$

Proof. We notice that when forwarding a random sample $D(j)$ of \mathbb{P}^* is the random tensor

$$e_{D(j)} [X_{[d]}]$$

and since $\text{im}(\phi_l) \subset \{0, 1\}$ the contraction

$$\langle \phi_l, e_{D(j)} [X_{[d]}] \rangle [\emptyset]$$

is a random variable taking values in $\{0, 1\}$. The variable therefore follows a Bernoulli distribution with mean parameter

$$p^{(l)} = \mathbb{E} [\langle \phi_l, e_{D(j)} [X_{[d]}] \rangle [\emptyset]] = \langle \phi_l, \mathbb{P}^* \rangle [\emptyset] .$$

□

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

$$\mu_D [L] = \langle \gamma^{\mathcal{H}}, \mathbb{P}^D \rangle [L] .$$

The expectation of this random tensor is

$$\mathbb{E} [\mu_D] = \langle \gamma^{\mathcal{H}}, \mathbb{E} [\mathbb{P}^D] \rangle [L] = \langle \gamma^{\mathcal{H}}, \mathbb{P}^* \rangle [L] = \mu^* ,$$

where we used that the expectation and contraction operation can be commuted due to the multilinearity of contractions.

10.2.1 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

$$\mathbb{E} [E^{\gamma^{\mathcal{H}^T}, \mathbb{P}^D - \tilde{\mathbb{P}}}] = \mathbb{E} [\langle \gamma^{\mathcal{H}^T}, \mathbb{P}^D - \tilde{\mathbb{P}} \rangle [L]] = \langle \gamma^{\mathcal{H}^T}, \mathbb{E} [\mathbb{P}^D - \tilde{\mathbb{P}}] \rangle [L] = \langle \gamma^{\mathcal{H}^T}, \mathbb{P}^* - \tilde{\mathbb{P}} \rangle [L] = \mathbb{E} [E^{\gamma^{\mathcal{H}^T}, \mathbb{P}^* - \tilde{\mathbb{P}}}] .$$

The fluctuation of this random tensor is

$$\mathbb{E} [E^{\gamma^{\mathcal{H}^T}, \mathbb{P}^D - \tilde{\mathbb{P}}}] - \mathbb{E} [E^{\gamma^{\mathcal{H}^T}, \mathbb{P}^* - \tilde{\mathbb{P}}}] = \mathbb{E} [E^{\gamma^{\mathcal{H}^T}, \mathbb{P}^D - \mathbb{P}^*}]$$

and coincides with $\eta^{\mathcal{H}, \mathbb{P}^*, D}$.

10.2.2 Naive Exponential Family

In case of the naive exponential family, which coincides with the minterm Markov Logic Network, we have $\phi = \delta [X_{[d]}, L]$ and the fluctuation tensor is

$$\eta^{\delta, \mathbb{P}^*, D} = \mathbb{P}^D - \mathbb{P}^* .$$

This fluctuation tensor is related to tensor encodings of multinomial distributions, which we now define as multinomial random tensors.

Definition 60. A multinomial random tensor is the sum of the one-hot encodings of independent random variables Z_j each distributed by \mathbb{P}

$$Z^{\mathbb{P}, m} = \sum_{j \in [m]} e_{Z_j} .$$

In the case of naive exponential families, the fluctuation tensor is a multinomial, as we show next. This characterization goes beyond the characterization of the marginal distributions as centered binomial variables, which holds for generic Markov Logic Networks.

Lemma 14. The fluctuation $\mathbb{P}^D - \mathbb{P}^*$ is a by $\frac{1}{m}$ rescaled centered multinomial random tensor with parameters \mathbb{P}^* and m .

Proof. By the above construction we have

$$\mathbb{P}^D - \mathbb{P}^* = \frac{1}{m} (e_{D(j)} [X_{[d]}] - \mathbb{E} [e_{D(j)} [X_{[d]}]]) .$$

□

10.3 Expected and Empirical Risk Minimization

We take a frequentist approach here and study the distributions of estimated parameters depending on random data.

10.3.1 Maximum Likelihood Estimation on random data

We here investigate the statistical errors of the maximum likelihood estimators. The empirical distribution is understood as an estimation of an underlying distribution. When sampling by independent copies of the underlying distribution, its mean is the underlying distribution. However, due to fluctuations around this mean the solution of estimation problems with respect to the empirical distribution does in general not coincide with the solution given the underlying distribution.

To be more precise, when generating data $\{D(j)\}_{j \in [m]}$ by independent copies of a distribution \mathbb{P}^* we define the minimization problems

$$\theta^D = \operatorname{argmin}_{\theta \in \Gamma} \mathbb{H} \left[\mathbb{P}^D, \mathbb{P}^{(\phi, \theta, \nu)} \right] \quad (P_{\{D(j)\}_{j \in [m]}, \Gamma})$$

and

$$\theta^* = \operatorname{argmin}_{\theta \in \Gamma} \mathbb{H} \left[\mathbb{P}^*, \mathbb{P}^{(\phi, \theta, \nu)} \right] \quad (P_{\mathbb{E}, \Gamma})$$

where by $\mathbb{P}^{(\phi, \theta, \nu)}$ we denote the element of an exponential family with sufficient statistics ϕ and parameter θ .

We have at each hypothesis $\theta \in \Gamma$

$$\mathbb{E} \left[\mathbb{H} \left[\mathbb{P}^D, \mathbb{P}^{(\phi, \theta, \nu)} \right] \right] = \mathbb{H} \left[\mathbb{P}^*, \mathbb{P}^{(\phi, \theta, \nu)} \right] \quad (41)$$

and thus, the objective in Problem $P_{\{D(j)\}_{j \in [m]}, \Gamma}$ converges in the limit $m \rightarrow \infty$ by the law of large numbers at each hypothesis $\theta \in \Gamma$ to the objective in Problem $P_{\mathbb{E}, \Gamma}$.

To use this insight in the derivation of bounds of the distance of θ^D and θ^* , we need to quantifying this convergence

- Non-asymptotically: Since we typically have access to limited amounts of data, that is finite m , we need to quantify the concentration in non-asymptotic cases.
- Uniform: Since the problems are optimized at extreme situations, the convergence of the objective has to happen uniformly at multiply $\theta \in \Gamma$.

10.3.2 Solution of the Expected Problem

Let us first investigate the solution of Problem $P_{\mathbb{E}, \Gamma}$, to get a reference for Problem $P_{\{D(j)\}_{j \in [m]}, \Gamma}$.

The solution does not change when manipulating the objective as

$$\operatorname{argmin}_{\theta \in \Gamma} \mathbb{H} \left[\mathbb{P}^*, \mathbb{P}^{(\phi, \theta, \nu)} \right] = \operatorname{argmin}_{\theta \in \Gamma} \mathbb{H} \left[\mathbb{P}^*, \mathbb{P}^{(\phi, \theta, \nu)} \right] - \mathbb{H} [\mathbb{P}^*, \mathbb{P}^*] \quad (42)$$

We notice that the objective of the right side problem is the Kullback-Leibler divergence

$$D_{\text{KL}} \left[\mathbb{P}^* || \mathbb{P}^{(\phi, \theta, \nu)} \right] \quad (43)$$

Thus, Problem $P_{\mathbb{E}, \Gamma}$ coincides with the approximation of \mathbb{P}^* by an element $\mathbb{P}^{(\phi, \theta, \nu)}$ with $\theta \in \Gamma$ in the Kullback-Leibler divergence.

We decompose:

$$D_{\text{KL}} \left[\mathbb{P}^{(\phi, \theta^*, \nu)} || \mathbb{P}^{(\phi, \theta, \nu)} \right] = \log \left[\frac{\mathcal{Z}(\theta)}{\mathcal{Z}(\theta^*)} \right] + \frac{1}{\mathcal{Z}(\theta^*)} \langle \theta^* - \theta, \exp [\theta^*] \rangle \quad (44)$$

Some insights can be drawn based on this decomposition:

- The first term vanishes, when both partition functions are the same. We can always adjust θ by constant offsets on all coordinates (of course without changing the distribution), such that the partition functions are equal. This is done by the map

$$\theta \rightarrow \theta + \lambda \cdot \mathbb{I}$$

where we choose $\lambda \in \mathbb{R}$ by

$$\lambda = \frac{\mathcal{Z}(\theta^*)}{\mathcal{Z}(\theta)}$$

to ensure a vanishing first term.

- In the second term the differences in coordinates are weighted by exponentiated solution θ^* . Where the probability mass is small errors in θ have a small influence on the Kullback-Leibler divergence.

Theorem 67. *Let us assume $\mathbb{P}^* = \mathbb{P}^{(\phi, \theta^*, \nu)}$ for some $\theta^* \in \Gamma$. Then one solution of Problem $P_{\mathbb{E}, \Gamma}$ coincides with θ^* up to constant offsets.*

Proof. This follows directly from the Gibbs inequality, since

$$D_{\text{KL}} \left[\mathbb{P}^{(\phi, \theta^*, \nu)} \parallel \mathbb{P}^{(\phi, \theta, \nu)} \right] \leq 0 \quad (45)$$

with equality if only if $\mathbb{P}^{(\phi, \theta^*, \nu)}$ and $\mathbb{P}^{(\phi, \theta, \nu)}$ are equal. The uniqueness follows from Theorem ?? □

10.3.3 Recovery Guarantee based on Widths

Theorem 68. *Let us assume $\theta^* \in \Gamma$ and that $\mathcal{Z}(\theta)$ is constant among $\theta \in \Gamma$. Then for any solution θ^D of the empirical problem we have*

$$D_{\text{KL}} \left[\mathbb{P}^{(\phi, \theta^*, \nu)} \parallel \mathbb{P}^{(\phi, \theta^D, \nu)} \right] \leq \omega_{\Gamma}(\eta). \quad (46)$$

Proof. First we notice

$$\operatorname{argmin}_{\theta \in \Gamma} \mathcal{L}_D \theta = \operatorname{argmin}_{\theta \in \Gamma} \mathcal{L}_D \theta - \mathcal{L}_D \theta^* \quad (47)$$

When $\theta^* \in \Gamma$ the minimum of the empirical loss with respect to θ^* is negative since

$$\mathcal{L}_D \theta^D - \mathcal{L}_D \theta^* \leq \mathcal{L}_D \theta^* - \mathcal{L}_D \theta^* \leq 0 \quad (48)$$

We separate expectations and fluctuations and get

$$D_{\text{KL}} \left[\mathbb{P}^{\theta^*} \parallel \mathbb{P}^{\theta^D} \right] \leq D_{\text{KL}} \left[\mathbb{P}^{\theta^*} \parallel \mathbb{P}^{\theta^D} \right] - (\mathcal{L}_D \theta^D - \mathcal{L}_D \theta^*) \leq \omega_{\Gamma}(). \quad (49)$$

□

The supremum of the differences between expected and empirical risks is the width of the fluctuation tensor, as we state next.

Lemma 15. *For any Γ and D we have*

$$\omega_{\Gamma}(\eta^{\delta, \mathbb{P}^*, D}) = \sup_{\theta \in \Gamma} \mathbb{H} \left[\mathbb{P}^D, \mathbb{P}^{(\phi, \theta, \nu)} \right] - \mathbb{H} \left[\mathbb{P}^*, \mathbb{P}^{(\phi, \theta, \nu)} \right]$$

Proof. Using the decomposition of cross entropy in the naive exponential family

$$\mathbb{H} \left[\mathbb{P}^D, \mathbb{P}^{(\phi, \theta, \nu)} \right] = \langle \mathbb{P}, \ln [\mathbb{P}^*] \rangle [-] A^{(\phi, \nu)}(\ln [\mathbb{P}^*]).$$

□

Corollary 10. *At the solution θ^D of Problem $P_{\{D(j)\}_{j \in [m]}, \Gamma}$ we have*

$$\mathbb{H} \left[\mathbb{P}^D, \mathbb{P}^{(\phi, \theta^D, \nu)} \right] - \mathbb{H} \left[\mathbb{P}^*, \mathbb{P}^{(\phi, \theta^D, \nu)} \right] \leq \omega_{\Gamma}(\eta^{\delta, \mathbb{P}^*, D}).$$

10.4 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.

Definition 61. The max gap of a tensor T $[X_{[d]}]$ is the quantity

$$\delta(T) = \min_{x_{[d]} \neq x_{[d]}^{max}} T[X_{[d]} = x_{[d]}^{max}] - T[X_{[d]} = x_{[d]}]$$

where

$$x_{[d]}^{max} \in \operatorname{argmax}_{x_{[d]}} T[X_{[d]} = x_{[d]}] .$$

If multiple maxima exist, the

Theorem 69. Whenever the energy tensor of the expected proposal distribution has a gap of δ , then for every $\epsilon > 0$ the mode of the expected proposal distribution coincides with the empirical proposal distribution with probability at least $1 - \exp[-\frac{1}{\epsilon^2}]$, provided that

$$m > C \frac{\left(\sum_{k \in [d]} \ln[m_k]\right)}{\delta^2}$$

where C is a universal constant.

To proof the theorem we first use a deterministic recovery guarantee involving the width of the fluctuation tensor and then apply the width bound of Theorem 120.

Lemma 16. Whenever

$$\delta(\langle \mathbb{P}^*, \mathcal{H} \rangle [L]) > 2 \cdot \omega_{\{e_{x_{[d]}} : x_{[d]} \in \times_{k \in [d]} [m_k]\}}(\eta^{\mathcal{H}, \mathbb{P}^*, D}),$$

then the mode of the proposal distribution to the empirical distribution coincides with the mode of the proposal distribution to the generating distribution.

Proof. If different, then the expected objective at the solutions of the empirical and expected is at least δ . But at the empirical, the difference is at most twice the width different, so that is a contradiction to the assumption. \square

Proof of Theorem 69. Given the assumed bound, the sub-gaussian norm of the width is upper bounded by $C_2 \cdot \delta$, thus for any $\epsilon > 0$ we have

$$\omega_{\{e_{x_{[d]}} : x_{[d]} \in \times_{k \in [d]} [m_k]\}}(\eta^{\mathcal{H}, \mathbb{P}^*, D}) < 2\delta$$

with probability at least $1 - \exp[-\frac{1}{\epsilon^2}]$. The claim thus follows with Lemma 16. \square

Example 14 (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

$$\delta(E^{(\mathcal{F}, \mathbb{P}^{\{f\}, \theta_f}) - \langle \mathbb{I} \rangle [X_{[d]} | \emptyset]}) = \frac{1}{2^d - 1 + \exp[-\theta_f]}$$

If $\theta_f > 0$ 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

$$\delta(E^{(\mathcal{F}, \mathbb{P}^{\{f\}, \theta_f}) - \langle \mathbb{I} \rangle [X_{[d]} | \emptyset]}) = \frac{1}{1 + (2^d - 1) \cdot \exp[-\theta_f]}$$

For large θ_f and d , the gap tends to 1.

10.5 Guarantees for Parameter Estimation

This is mean parameter fluctuation interpretation of the random tensor.

Lemma 17. For any \mathcal{H} and D drawn from \mathbb{P}^* we have

$$\|\mu_D - \mu^*\|_2 = \omega_S(\eta^{\mathcal{H}, \mathbb{P}^*, D}),$$

where $\mu_D = \langle \mathcal{H}, \mathbb{P}^D \rangle [L]$ and $\mu^* = \langle \mathcal{H}, \mathbb{P}^* \rangle [L]$.

We can thus apply the sphere bounds.

Theorem 70. *For any $\epsilon \in (0, 1)$ we have the following with probability at least $1 - \epsilon$. Let $\hat{\theta}$ and $\tau > 0$, then*

$$\left| \mathbb{H} \left[\mathbb{P}^*, \mathbb{P}^{(\mathcal{F}, \theta^D, \nu)} \right] - \mathbb{H} \left[\mathbb{P}^D, \mathbb{P}^{(\mathcal{F}, \theta^D, \nu)} \right] \right| \leq \tau \cdot \|\theta^D\|_2$$

provided that

$$m \geq \frac{\langle \mu^* \rangle [\emptyset] - \langle (\mu^*)^2 \rangle [\emptyset]}{\epsilon \tau^2}.$$

Proof. We have by Cauchy Schwartz

$$|\langle \mu_D - \mu^*, \theta^D \rangle [\emptyset]| \leq \|\mu_D - \mu^*\|_2 \cdot \|\theta^D\|_2$$

and with Lemma 17

$$|\langle \mu_D - \mu^*, \theta^D \rangle [\emptyset]| \leq \omega_S(\eta^{\mathcal{H}, \mathbb{P}^*, D}) \cdot \|\theta^D\|_2.$$

We show in Part III that in Theorem 121 that

$$\omega_S(\eta^{\mathcal{H}, \mathbb{P}^*, D}) \leq \tau$$

when m satisfies the assumed lower bound, from which the claim follows. \square

11 Statistical Models of First Order Logic (FOL)

We now extend the tensor representation to structured representations, whereas the above was on factored representation of systems. The models/events in this situation are precise relations between objects.

11.1 World Tensors

We introduce here object enumerating variables O are valued in $[r]$, with $|A| = r$ and an enumeration of A by a_o .

Formulas in first order logic can contain variables, which are placeholder for specific individuals. Given a model and an assignment of objects to the arguments of a formula, the truth of the formula can be interpreted. This truth interpretation defines thus for any model a tensor, which we call the grounding tensor.

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 A .

Definition 62 (FOL World). *Given a set of objects $A = \{a_o : o \in [r]\}$ enumerated by $[r]$ and a finite set $\{g_0, \dots, g_{d-1}\}$ of n -ary predicates a world is a Boolean tensor*

$$W[L, O_{[n]}] : [d] \times \left(\bigtimes_{l \in [n]} [r] \right) \rightarrow [2]. \quad (50)$$

We interpret the world tensor as encoding in the coordinate $W[L = k, O_{[n]} = o_{[n]}]$, whether the substitution of the variables in the k th predicate by the individuals $a_{o_0}, \dots, a_{o_{n-1}}$ is satisfied on a corresponding world, that is

$$g_k(a_{o_0}, \dots, a_{o_{n-1}}) = W[L = k, O_{[n]} = o_{[n]}].$$

When the assumptions 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 $e_0 [O]$, where we add an object a_0 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.

11.1.1 Special Case of Propositional Logics

Propositional Logics is the case of $n = 0$. In that case, worlds are

$$W : [d] \rightarrow [2]$$

and represented by the tuples $x_k = W(k)$ as before.

Compared with propositional formulas, the grounding tensor does not take as input a specific world, but is defined on a given world. We show in this chapter, how both tensor interpretations can be transformed, i.e. by extracting samples from a FOL world W interpreted as an empirical distribution over PL worlds, and by generating FOL worlds by a set of samples generated from a PL distribution.

11.1.2 Enumeration of worlds

We here enumerate worlds coinciding in their domains. This is typical in database semantics, where only those worlds are considered which domains have a one-to-one map to the constant symbols appearing in a knowledge base. As above we use term variables O to represent objects in the domain A of a world.

Propositional worlds have been enumerated by indices of d Booleans, that is for a world $W : [d] \rightarrow [2]$ we take the index

$$x_0, \dots, x_{d-1} \quad \text{where} \quad x_k = W(k).$$

When we want to enumerate the first order logic worlds to a fixed set of objects A , we flatten the tensor W and get indices

$$\{x_{k,o_0,\dots,o_{n-1}} : k \in [d], o_0, \dots, o_{n-1} \in [r]\} \quad \text{where} \quad x_{k,o_0,\dots,o_{n-1}} = W(k, o_0, \dots, o_{n-1}).$$

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.

11.1.3 Random World Tensors

Random worlds with fixed domains A are probability distributions among possible tensors.

All possible world tensors given a set of objects A are enumerated by a Boolean to each coordinate of the tensor, indicating whether the corresponding assertion holds in the world. Random worlds are then probability distributions

$$\mathbb{P} : \bigtimes_{k \in [d], o_0, \dots, o_{n-1} \in [r]} [2] \rightarrow \mathbb{R}.$$

We notice, that by definition these probability distributions are distributions of $d \cdot r^n$ Booleans with $2^{(d \cdot r^n)}$ many states. We will later in this chapter investigate methods to handle such high-dimensional distributions in the formalism of exponential families.

11.2 Formulas in first order logics

Given a FOL World W , arbitrary formulas in first order logics are interpreted in terms of the satisfactions of their groundings.

Definition 63 (Grounding of FOL Formulas). *Given a specific world W , with domain A enumerated by $[r]$, the grounding of a formula q with n_q variables is the tensor*

$$q|_W [O_q] : \bigtimes_{l \in [n_q]} [r] \rightarrow [2].$$

We interpret the coordinates as encoding, whether the substitution of the variables in the formula is satisfied given a world W , that is

$$q|_W [o_0, \dots, o_{n_q-1}] = \begin{cases} 1 & \text{if } q(a_{o_0}, \dots, a_{o_{n_q-1}}) = 1 \text{ given the world } W \\ 0 & \text{else} \end{cases}.$$

Slicing the grounding tensor of a FOL Formula amounts to substitution of the respective variable by the constant at the enumeration index.

The above definition of formulas can be interpreted as a map from worlds to grounding tensors

$$q : W \rightarrow q|_W .$$

When interpreting this map as a basis encoding, formulas are tensors in the tensor space

$$\left(\bigotimes_{k \in [d], o_0, \dots, o_{n-1} \in [r]} \mathbb{R}^2 \right) \otimes \left(\bigotimes_{k \in [d], o_0, \dots, o_{n_q} \in [r]} \mathbb{R}^2 \right) .$$

Atomic formulas are stored in the slices to the first axis of W and we have

$$g_k|_W = \langle \{W[L, O_{[n]}], e_k[L]\} \rangle [O_{[n]}] . \quad (51)$$

11.2.1 Syntactical Decomposition of quantifier-free formulas

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 ρ -formalism. While the atomic formulas where delta tensors copying states, they are more involved here.

Theorem 71. *For any connective \circ and formulas q_1 and q_2 we have*

$$(q_1 \circ q_2)|_W [O_{[n]}] = \left\langle \{ \rho^{q_1|_W} [O_{[n]}, X_{q_1}], \rho^{q_2|_W} [O_{[n]}, X_{q_2}], \rho^\circ [X_{q_1 \circ q_2}, X_{q_1}, X_{q_2}], e_1 [X_{q_1 \circ q_2}] \} \right\rangle [O_{[n]}] . \quad (52)$$

Proof. By the semantic interpretation of the groundings, which has to be sound. \square

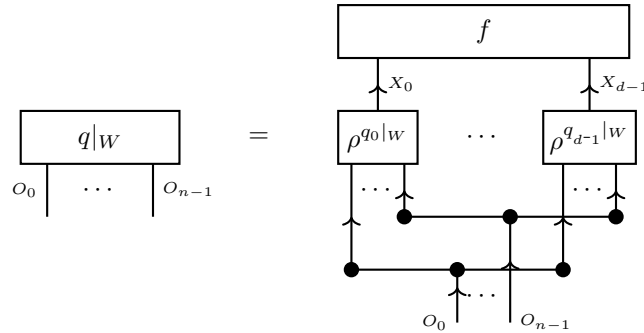
Here, variables can be shared by the connected formulas, therefore the variables in the combined formula are unions of the possible not disjoint variables of the connected formulas.

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 64. *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 W we have*

$$q|_W [O_q] = \left\langle \{ \rho^{q_k|_W} [O_{q_k}, X_k] : k \in [d] \} \cup \{ f [X_{[d]}] \} \right\rangle [O_q] \quad (53)$$

We depict the relation of a grounding tensor to a propositional formula



11.2.2 Quantifiers

Existential and universal quantifiers appear in generic first order logic. They are not representable as linear transform of the respective quantifier-free formula. Here we represent them as a compositions of contractions with the greater zero function ψ defined as

$$\psi : \mathbb{R} \rightarrow [2] \quad , \quad \psi(x) = \begin{cases} 1 & \text{if } x > 0 \\ 0 & \text{else} \end{cases} .$$

Theorem 72. *For any formula q with variables $O_{[n]}$ we have*

$$(\exists_{a_{o_0}} q)|_W [O_1, \dots, O_{n-1}] = \psi(\langle q|_W \rangle [O_1, \dots, O_{n-1}]) \quad (54)$$

and

$$(\forall_{a_{o_0}} q)|_W [O_1, \dots, O_{n-1}] = \psi(\langle q|_W \rangle [O_1, \dots, O_{n-1}] - r \cdot \mathbb{I}[O_1, \dots, O_{n-1}]) . \quad (55)$$

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, \dots, o_{n-1} we have

$$\langle q|_W \rangle [O_1 = o_1, \dots, O_{n-1} = o_{n-1}] = \sum_{o_0 \in [r_0]} q|_W [O_0 = o_0, \dots, O_{n-1} = o_{n-1}] .$$

The existential quantification with respect to the zeroth variable is satisfied, when for at least one $o_0 \in [r_0]$ we have $q|_W [O_0 = o_0, \dots, O_{n-1} = o_{n-1}] = 1$, which is equal to

$$\sum_{o_0 \in [r_0]} q|_W [O_0 = o_0, \dots, O_{n-1} = o_{n-1}] \geq 1$$

and

$$\psi \left(\sum_{o_0 \in [r_0]} q|_W [O_0 = o_0, \dots, O_{n-1} = o_{n-1}] \right) = 1 .$$

Combining both equations we get

$$(\exists_{a_{o_0}} q|_W) [O_1 = o_1, \dots, O_{n-1} = o_{n-1}] = \psi(\langle q|_W \rangle [O_1 = o_1, \dots, O_{n-1} = o_{n-1}])$$

and arrive at the first claim.

To show the second claim, we notice that the universal quantification with respect to the zeroth variable is satisfied, when for all $o_0 \in [r_0]$ we have $q|_W [O_0 = o_0, \dots, O_{n-1} = o_{n-1}] = 1$. This is equal to

$$\sum_{o_0 \in [r_0]} q|_W [O_0 = o_0, \dots, O_{n-1} = o_{n-1}] = r$$

and

$$\psi \left(\sum_{o_0 \in [r_0]} q|_W [O_0 = o_0, \dots, O_{n-1} = o_{n-1}] - r \right) = 1 .$$

It follows that

$$\begin{aligned} (\forall_{a_{o_0}} q|_W) [O_1 = o_1, \dots, O_{n-1} = o_{n-1}] \\ = \psi(\langle q|_W \rangle [O_1 = o_1, \dots, O_{n-1} = o_{n-1}] - r \cdot \mathbb{I}[O_1 = o_1, \dots, O_{n-1} = o_{n-1}]) \end{aligned}$$

which establishes the second claim. \square

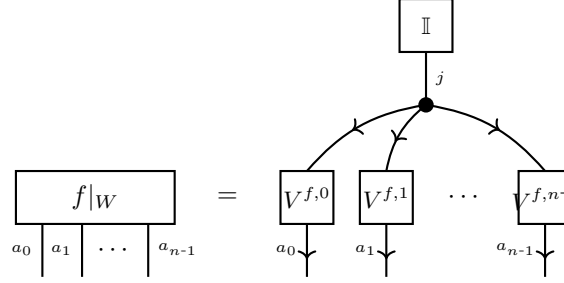
11.2.3 Basis CP Decomposition

In many situations grounding cores are sparse and representations as single tensor cores comes with a drastic overhead. We often encounter sparse grounding, that is

$$\ell_0(q|_W) \leq r^{n_q} .$$

In this case the basis CP decomposition (see Section 17.1.2) enables a representation of the grounding with significantly smaller storage demand, see Theorem 106. We depict the CP decomposition of a formula grounding in Figure 26.

Most logical syntaxes make use of the ℓ_0 -sparsity and explicitly store only the known assertions. The not specified assertions have different interpretation based on the assumptions made. Under the Closed World Assumption all unspecified assertions are assumed to be false.

Figure 26: Basis CP Decomposition of the grounding of q .

11.2.4 Example: Queries

Queries are formulas p , which are asked against a specific world W . Each argument of such formulas are called projection variables. The answer of a query is a list of solution mappings from the projection variables to objects in the world, such that the query formula is true.

Answering a query thus corresponds with finding the basis CP Decomposition (see Section 17.1.2).

There are situations, in which the CP Decomposition comes with unnecessary storage overhead compared with other formats. To be more precise, let us show how Join operations on groundings are represented by contractions.

Definition 65. Let there be a world with individual sets J_v for $v \in \mathcal{V}$ and groundings $q_l|_W$

$$q_l|_W : \bigtimes_{v \in \mathcal{V}_l} [r_v] \rightarrow [2]$$

defined on subsets of the individuals sets $\bigtimes_{v \in \mathcal{V}} J_v$. Then their Join is defined as

$$\text{Join}(q_0, \dots, q_{p-1})|_W = \prod_{l \in [p]} q_l|_W.$$

We can understand the Join of groundings by a factor graph, where each grounding is a factor connecting the variables decorating node sets \mathcal{V}_l .

Theorem 73 (Joins by Hadamard products). *For any Join of grounded formulas we have*

$$\text{Join}(q_0, \dots, q_{p-1})|_W [O_{\mathcal{V}}] = \langle q_0|_W, \dots, q_{p-1}|_W \rangle [O_{\mathcal{V}}].$$

Proof. Direct □

Thus, we can store the Join by at most linear in the number of joined formulas demand. Storing the result by solution maps on the other side can create storage overheads as we show next.

Example 15. For example take a query with many basic graph patterns with pairwise different projection variables. The global CP Decomposition would come here with an exponential storage overhead compared with storage as a tensor product of CP Decompositions to each Basic graph pattern.

Remark 29 (Distinguishing from probabilistic queries). *Let us distinguish the discussion here from those of queries in probabilistic reasoning, which have two main differences. First, we ask queries against all possible pairs of variables, instead of asking the probability of satisfaction of a specific formula. Second, since we made the epistemologic assumption of knowing possibilities and not probabilities in logics, a query is answered by a truth value. We then only output in the shape of solution mappings the variable assignments where the query formula is true. Thus, the queries here can be thought of as a batch of probabilistic queries with Boolean answers. Probabilistic queries can furthermore be understood in terms of the data extraction process described in this section. We can ask the query in probabilistic form (decomposed into atomic formulas) on the resulting empirical distribution. This results in the ratio of the worlds satisfying the query among those worlds satisfying the extraction query p .*

11.3 Representation of Knowledge Graphs

Let us now represent a specific fragment of FOL, namely Description Logics which Knowledge Bases are often referred to as Knowledge Graphs.

11.3.1 Representation as unary and binary predicates

Atomic formulas in knowledge bases following the OWL standard (thus, especially knowledge graphs) are binary (owl:ObjectProperties) and unary (owl:Class). We enumerate the predicates by $[d]$, the objects in the domain A by $[r]$, and extend the unary predicates to binaries by tensor product with $e_0 [O]$ (see discussion above) A Knowledge Graph on the set A of constants (owl:NamedIndividuals) is the tensor

$$kg : [d] \times [r] \times [r] \rightarrow [2] .$$

11.3.2 Representation as ternary predicate

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 kg with individuals A

$$\text{RDF}|_{kg} : [r] \times [r] \times [r] \rightarrow [2]$$

where

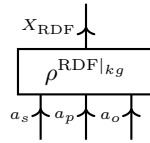
$$\text{RDF}|_{kg}(o_s, o_p, o_o) = \begin{cases} 1 & \text{if triple } \langle o_s, o_p, o_o \rangle \text{ is in Knowledge Graph } kg \\ 0 & \text{else} \end{cases} .$$

Sparse representation of the grounding tensor to a knowledge graph is of central importance, as investigated in Bigerl et al. (2020). We here do basis CP for sparse representation. Slicing the tensor ρ^{kg} along the predicate axis retrieves specific information about roles and can be efficiently be performed on these formats.

The role `rdf:type` 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 `rdf:type` as a placeholder for unitary formulas in a space of binary formulas.

Approximations of grounding tensors by decompositions leads to embeddings of the individuals such as Tucker, Complex, RESCAL (see Nickel et al. (2016)).

For our purposes of evaluating logical formulas such as SPARQL queries we use the relational encoding of the groundings, which are depicted by



11.3.3 SPARQL Queries

The SPARQL query language is a syntax to express first order logic formulas q . Given a specific world W , the execution of query is the interpretation $q|_W$, typical represented in a sparse basis CP format where each slice represents a solution mapping.

Triple Patterns Formulas which are of specific importance are triple patterns:

- Unary triple pattern with one variable, representing a formula with a single projection variable. For the example $\langle ?a_0, a, C \rangle$ see Figure 27a.
- Binary triple pattern with two variables, representing a formula with two projection variables. For the example $\langle ?a_0, R, ?a_1 \rangle$ see Figure 27b.

These triples are compositions of slicing operation, here seen as a map from single or two projection variables to the triples $(a_{o_s}, a_{o_p}, a_{o_o})$, with the RDF formula.

These slicing operations are performed by contractions with tensors $\psi^{(\dots)}$, for unary triple patterns of the shape

$$\psi : ([d] \times [r] \times [r]) \times [r] \rightarrow [2]$$

and for binary

$$\psi : ([d] \times [r] \times [r]) \times [r] \times [r] \rightarrow [2]$$

defined by identity and basis vectors depending on the type of triple.

In the example $\langle ?a_0, a, C \rangle$ we have

$$\psi^{\langle ?a_0, a, C \rangle} = \delta^{a_s, ?a_0} \otimes e_0 \otimes \delta^{a_o, ?a_1}$$

The composition $\psi(\psi^T)$ of the matrifaction of the tensor ψ is a projection. That means that applying $\psi(\psi^T)$ is the same map as applying once.

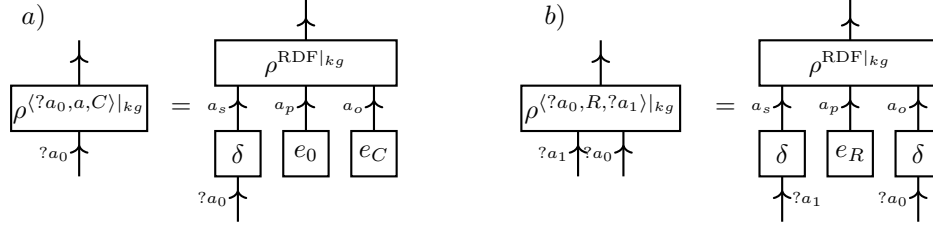


Figure 27: Triple patterns of SPARQL as tensor networks. a) Example of unary triple pattern $\langle ?a_0, a, C \rangle$ specifying whether an individual a_1 is a member of class C . b) Example of a binary triple pattern $\langle ?a_0, R, ?a_1 \rangle$ specifying whether individuals a_1 and a_2 have a relation R . By e_0, e_C, e_R we denote the one-hot encodings of the enumeration of the resources $rdf : type, C$ and R .

Basic Graph Patterns Generic SPARQL queries are compositions of triple patterns by logical connectives. Statements in SPARQL can be translated into Propositional Logics combining the triple patterns:

SPARQL	Propositional Logics
$\{f_1, f_2\}$	$f_1 \wedge f_2$
UNION $\{f_1, f_2\}$	$f_1 \vee f_2$

A SPARQL query p consistent of multiple triple pattern is the composition of the triple patterns groundings. We can represent its groundings by a contraction of the relational encodings of basic graph patterns and the respective logical connectives. Alternatively, the more direct effective calculus developed in Section 15.4 can be applied. The latter way is especially compelling, since large parts of typical SPARQL queries are mere conjunctions of the triple patterns, which is reflected in SPARQL syntax (default interpretation of instruction-free lists of patterns are conjunctions, while differing connectives need to be specified by additional instructions).

Further SPARQL features are

- FILTER $\{\cdot\}$ does not depend on triple patterns (e.g. numeric inequalities, regex functions on strings). We can regard it as another basic formula, which does not result from a slicing of the RDF grounding tensor. Besides that, we can understand it as formulas and include it in compositions.
- OPTIONAL $\{\cdot\}$ would result in \mathbb{I} leg vectors, when there is a missing variable assignment resulting.

11.4 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 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.

11.4.1 Markov Logic Networks in FOL

Following Richardson and Domingos (2006) we define Markov Logic Networks as templates for distributions, which instantiate random worlds when choosing a set of objects A . 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.

Definition 66 (Markov Logic Networks in FOL). *A Markov Logic Network is a template of probability distributions defined by a set \mathcal{Q} of FOL formulas with maximal arity n , which is weighted by a function $\theta : \mathcal{Q} \rightarrow \mathbb{R}$. The Markov*

Logic Network instantiated for a given set of objects A and a base measure ν is the random world, which is a member of the exponential family with sufficient statistics

$$\phi_l(W) = \langle q_l|_W \rangle [\emptyset]$$

and canonical parameters θ .

The statistics

$$\langle q_l|_W \rangle [\emptyset]$$

can be interpreted as the number of substitutions to a formula, such that the formula is satisfied. Each substitution satisfying a formula adds a factor of $\exp[\theta_l]$ to the probability of the respective world before normalization.

When constructing a world tensor to a theory with predicates of different order, we already argued that we extend the arity of predicates by tensor products with e_0 . To define random world tensors, we then restrict the corresponding base measure to be supported only on those worlds where the extended predicates hold only at the individual a_0 at the extended axis.

We choose extraction formulas q_k such that any formula in the FOL MLN has a propositional equivalent (see Definition 64). The statistic map is then a formula selecting tensor as in the propositional logic case contracted with the groundings of q_k .

11.4.2 Importance Formula

Analogous to a guard formula in (Koller and Friedman, 2009, Definition 6.11)!

We want to 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 objects relevant to the distributions by an 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 Definition 73)

$$\underline{p} = \sum_{j \in [m]} \left(\bigotimes_{l \in [n]} e_{o_l^j} [O_l] \right).$$

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 MLN to those world coinciding with groundings of the importance formula coinciding with \underline{p} by designing a base measure

$$\nu_{\underline{p}} = \begin{cases} 1 & \text{if } p|_W = \underline{p} \\ 0 & \text{else} \end{cases}.$$

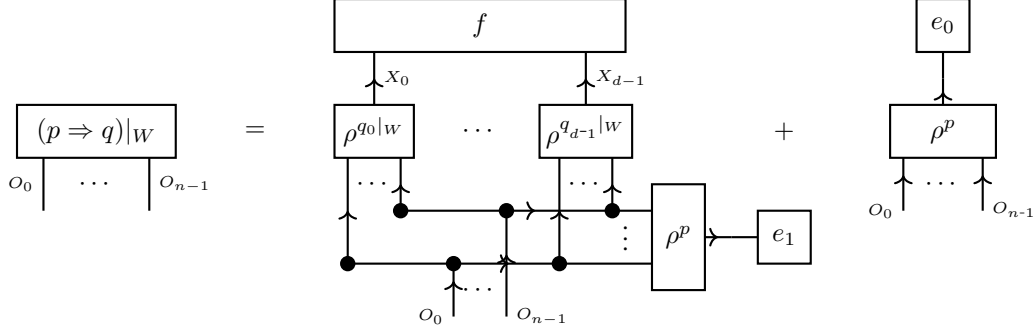
The base measure restricts the Markov Logic Network to be those worlds, where $p|_W$ is given by a fixed tensor \underline{p} . This amounts to assuming that $p|_W$ represents certain information about a FOL world, where other formulas are uncertain.

To reduce the likelihood of a world to we make the assumption that all formulas in a Markov Logic Network are of the form

$$q_l(a_{o_0}, \dots, a_{o_{n-1}}) = (p(a_{o_0}, \dots, a_{o_{n-1}}) \Rightarrow \tilde{q}(a_{o_0}, \dots, a_{o_{n-1}})) \quad (56)$$

that is a rule with the importance formula being the premise. When this assumption holds, we can think of the importance formula as a conditions on individuals to satisfy a statistical relation given by \tilde{q} .

We depict the assumption, that any formula is of the form (56) in the diagram



where the second summand depends only on the query p and therefore does not appear in the likelihood.

11.4.3 Decomposition of the log likelihood

Given a FOL MLN, the probability of a world W with $p|_W = \underline{p}$ is

$$\mathbb{P}(\mathcal{Q}|_A, \theta, \nu_{\underline{p}})[W] = \frac{1}{\mathcal{Z}(\mathcal{Q}|_A, \theta, \nu_{\underline{p}})} \exp \left[\sum_{q \in \mathcal{Q}} \theta_q \langle (p \Rightarrow \tilde{q})|_W \rangle [\emptyset] \right]$$

where the partition function is

$$\mathcal{Z}(\mathcal{Q}|_A, \theta, \nu_{\underline{p}}) = \sum_{W: p|_W = \underline{p}} \exp \left[\sum_{q \in \mathcal{Q}} \theta_q \langle (p \Rightarrow \tilde{q})|_W \rangle [\emptyset] \right]$$

Let us now decompose the statistics into constant and varying terms. We have

$$\langle (p \Rightarrow \tilde{q})|_W \rangle [\emptyset] = \langle p \wedge \tilde{q}|_W \rangle [\emptyset] + \langle \neg p|_W \rangle [\emptyset],$$

where the second term is constant among the supported worlds and the first can be enumerated by the satisfied substitutions of p , that is

$$\langle p \wedge \tilde{q}|_W \rangle [\emptyset] = \sum_{j \in [m]} \tilde{q}|_W [O_{[n]} = \sigma_{[n]}^j].$$

Using these insights we decompose a normalized log likelihood as

$$\frac{1}{m} \ln \left[\mathbb{P}(\mathcal{Q}|_A, \theta, \nu_{\underline{p}})[W] \right] = \frac{1}{m} \sum_{j \in [m]} \tilde{q}|_W [O_{[n]} = \sigma_{[n]}^j] - \frac{1}{m} \ln \left[\frac{\mathcal{Z}(\mathcal{Q}|_A, \theta, \nu_{\underline{p}})}{\exp [\langle \theta \rangle [\emptyset] \cdot \langle \neg p|_W \rangle [\emptyset]]} \right] \quad (57)$$

We notice a similarity with the likelihood in the case of MLNs in propositional logic. When we interpret each pair $a_{o_0}, \dots, a_{o_{n-1}} \in A$ satisfying $p(a_{o_0}, \dots, a_{o_{n-1}}) = 1$ as a datapoint, and choose the formulas \mathcal{F} from atomic formulas connected by propositional connectives, both approaches are equivalent. However, the partition function couples multiple samples, 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.

11.4.4 Interpretation as Likelihood of Propositional Dataset

Theorem 74. Let $\nu_{\underline{p}}$ be a base measure of worlds such that the vectors

$$\left(q_0|_{\tilde{W}} [O_0 = \sigma_0^j, \dots, O_{n-1} = \sigma_{n-1}^j], \dots, q_{d-1}|_{\tilde{W}} [O_0 = \sigma_0^j, \dots, O_{n-1} = \sigma_{n-1}^j] \right) \quad (58)$$

for $j \in [m]$ are independent and identical distributed by a distribution μ , when distributing \tilde{W} by $\nu_{\underline{p}}$.

Let there further be a set of formulas \mathcal{Q} , which formulas $q \in \mathcal{Q}$ are representable by

$$q = p \Rightarrow \tilde{q}$$

and we find a propositional formula f with

$$\tilde{q}|_W = \left\langle \{\rho^{q_k|_W} : k \in [d]\} \cup \{f\} \right\rangle [O_{[n]}] .$$

We then have for the likelihood of any by $\nu_{\underline{p}}$ supported world W that

$$\frac{1}{m} \ln \left[\mathbb{P}(\mathcal{Q}|_A, \theta, \nu_{\underline{p}})[W] \right] = \frac{1}{m} \ln \left[\mathbb{P}((\mathcal{F}, \theta), \mu)[D] \right] - \frac{1}{m} \sum_{j \in [m]} \ln [\mu[X_{[d]} = D(j)]]$$

where μ is the distribution of the random vectors (58) and by D we denote the dataset

$$D = \left\{ (q_0|_W [O_0 = o_0^j, \dots, O_{n-1} = o_{n-1}^j], \dots, q_{d-1}|_W [O_0 = o_0^j, \dots, O_{n-1} = o_{n-1}^j]) : j \in [m] \right\}$$

and \mathcal{F} is the set of propositional formulas to $q \in \mathcal{Q}$.

To show the theorem, we show first in the following lemma the factorization of the partition function of the FOL MLN.

Lemma 18. *Given the assumptions of Theorem 74, we have*

$$\frac{\mathcal{Z}(\mathcal{Q}|_A, \theta, \nu_{\underline{p}})}{\exp[\langle \theta \rangle [\emptyset] \cdot \langle \neg p|_W \rangle [\emptyset]]} = (\mathcal{Z}((\mathcal{F}, \theta), \mu))^m .$$

Proof. We have

$$\begin{aligned} \mathcal{Z}(\mathcal{Q}|_A, \theta, \nu_{\underline{p}}) &= \mathbb{E}_{W \sim \nu_{\underline{p}}} \left[\exp \left[\sum_{q \in \mathcal{Q}} \theta_q \langle (p \Rightarrow \tilde{q})|_W \rangle [\emptyset] \right] \right] \\ &= \exp[\langle \theta \rangle [\emptyset] \cdot \langle \neg p|_W \rangle [\emptyset]] \cdot \mathbb{E}_{W \sim \nu_{\underline{p}}} \left[\exp \left[\sum_{q \in \mathcal{Q}} \theta_q \sum_{j \in [m]} \tilde{q}|_W [O_{[n]} = o_{[n]}^j] \right] \right] \\ &= \exp[\langle \theta \rangle [\emptyset] \cdot \langle \neg p|_W \rangle [\emptyset]] \cdot \mathbb{E}_{W \sim \nu_{\underline{p}}} \left[\prod_{j \in [m]} \exp \left[\sum_{q \in \mathcal{Q}} \theta_q \cdot \tilde{q}|_W [O_{[n]} = o_{[n]}^j] \right] \right] \\ &= \exp[\langle \theta \rangle [\emptyset] \cdot \langle \neg p|_W \rangle [\emptyset]] \cdot \prod_{j \in [m]} \mathbb{E}_{W \sim \nu_{\underline{p}}} \left[\exp \left[\sum_{q \in \mathcal{Q}} \theta_q \cdot \tilde{q}|_W [O_{[n]} = o_{[n]}^j] \right] \right] \end{aligned}$$

Here we used, that since the substitutions of the atom formulas at the respective object tuples are independent, also the variables

$$\exp[\theta_q \cdot \tilde{q}|_W [O_{[n]} = o_{[n]}^j]]$$

for $j \in [m]$ are independent.

Since each $\tilde{q}|_W [O_{[n]}]$ depends only on the random value $x_k^j = q_k[O_{[n]}]$ we have

$$\begin{aligned} \mathbb{E}_{W \sim \nu_{\underline{p}}} \left[\exp \left[\sum_{q \in \mathcal{Q}} \theta_q \cdot \tilde{q}|_W [O_{[n]} = o_{[n]}^j] \right] \right] &= \sum_{x_{[d]} \in [2]^d} \mathbb{E}_{W \sim \nu_{\underline{p}}} [\forall k \in [d] : q_k[O_{[n]}] = x_k] \\ &\quad \cdot \exp \left[\sum_{f \in \mathcal{F}} \theta_f \cdot f[X_{[d]} = x_{[d]}] \right] \\ &= \sum_{x_{[d]} \in [2]^d} \mu[X_{[d]} = x_{[d]}] \cdot \exp \left[\sum_{f \in \mathcal{F}} \theta_f \cdot f[X_{[d]} = x_{[d]}] \right] \\ &= \mathcal{Z}((\mathcal{F}, \theta), \mu) . \end{aligned}$$

Combining the above, we arrive at the claim. \square

Proof of Theorem 74. Use Equation 57 and the Lemma 18. Note that we need to correct the likelihood by the average log base measure on the data, since that term is appearing in the likelihood of a MLN. \square

Let us now investigate, when the assumptions of independent data can be matched.

Lemma 19. *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 (58) 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

Situations where atom base μ measures are not uniform:

- 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
- extraction formulas contradicting each other, more general not independent from each other

11.4.5 Approximation by Independent Samples

In general, we cannot assume that the p, q_0, \dots, q_{d-1} do not influence each other.

We approximate the partition function into factors to each solution map of p . This amounts to the assumption, that the atoms created to each tuple are independent.

For example, when solution maps share the resources, which form the arguments of an atom extraction query, they coincide on this atom and are thus dependent.

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.

11.4.6 Extraction of Samples from FOL Knowledge Bases

The decomposition of the likelihood suggests the following approach to generate samples from groundings:

- Define for $k \in [d]$ queries q_k generating the the atoms X_k : Predicates along with assignment of variables / constants to its positions.
- Define a query formula p , which we decompose in the basis CP decomposition for later interpretation of datapoints
- Contract the groundings of each formula f^k with the grounding of p to build a data core

11.4.7 Representation by Tensor Networks

Let us now understand the extraction process as a relation between a tuple of individuals and the extracted world in the factored system of atoms X_k .

$$\mathcal{R} = \{(a_{o_0}, \dots, a_{o_{n-1}}, X_0, \dots, X_{d-1}) : p(a_{o_0}, \dots, a_{o_{n-1}}) = 1, \forall_{k \in [d]} : X_k = q_k(a_{o_0}, \dots, a_{o_{n-1}})\}$$

Towards constructing the encoding of this we enumerate the individuals in the set A and use in the following the respective one-hot encoding

$$e : A \rightarrow \mathbb{R}^{|A|}.$$

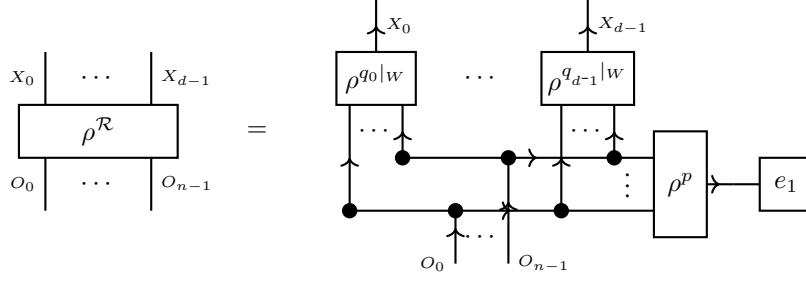
The combination of the queries p and $\{q_0, \dots, q_{d-1}\}$ by the relational encoding

$$\rho^{\mathcal{R}} \subset \left(\bigotimes_{l \in [n]} \mathbb{R}^r \right) \otimes \left(\bigotimes_{k \in [d]} \mathbb{R}^2 \right)$$

defined by

$$\rho^{\mathcal{R}} = \{(e_{a_{o_0}, \dots, a_{o_{n-1}}}, X_0, \dots, X_{d-1}) : p(a_{o_0}, \dots, a_{o_{n-1}}) = 1, \forall_{k \in [d]} : X_k = q_k(a_{o_0}, \dots, a_{o_{n-1}})\}.$$

We can represent the encoding $\rho^{\mathcal{R}}$ by the diagram



Here the contraction of ρ^p with the truth vector e_1 represents the matching condition posed by p when extracting pairs of individuals.

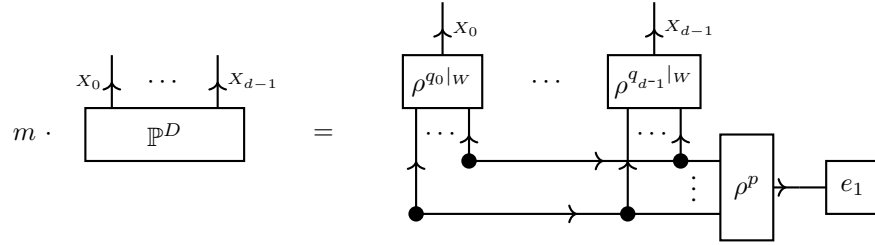
The empirical distribution is then the normalized contraction leaving only the legs to the extracted atomic formulas open, that is

$$\mathbb{P}^D = \frac{\langle \rho^{\mathcal{R}} \rangle [X_{[d]}]}{\langle \rho^{\mathcal{R}} \rangle [\emptyset]}.$$

Here the number of extracted data is the denominator

$$m = \langle \rho^{\mathcal{R}} \rangle [\emptyset].$$

We depict this by



11.4.8 Basis CP Decomposition of extracted data

To connect with the empirical distribution introduced in Section 3.8 we now show how the dataset D extracted from the interpretations of the formulas p, q_0, \dots, q_{d-1} on a FOL world W can be represented by tensor networks.

Each datacore is then a contraction with the grounding of a formula, which is contracted with the grounding of the extraction query in the basis CP decomposition,

$$\rho^{D_k} = \langle \rho^{q_k|W}, p|_W \rangle [J, X_k]$$

see Figure 28.

Theorem 75. *We have*

$$\langle \rho^{\mathcal{R}} \rangle [X_{[d]}] = \langle \{ \rho^{D_k} [J, X_k] : k \in [d] \} \rangle [X_{[d]}]$$

and thus

$$\mathbb{P}^D = \langle \{ \rho^{D_k} [J, X_k] : k \in [d] \} \rangle [X_{[d]} | \emptyset].$$

Proof. Using that $p|_W$ is binary and the core can be copied (ref to theorem in binary calculus). \square

Towards efficient calculation of the data cores, we build a basis CP decomposition of $p|_W$, where we further demand $\sigma = \mathbb{I}$. This is a collection of basis leg cores $V^{p,l}$ such that

$$\underline{p}[O_{[n]}] = \langle \{ V^{p,l} [J, O_l] : l \in [n] \} \rangle [O_{[n]}].$$

The log-likelihood of a MLN is then the contraction of the data cores with the representation of the MLN wrt the atomic formulas q_k .

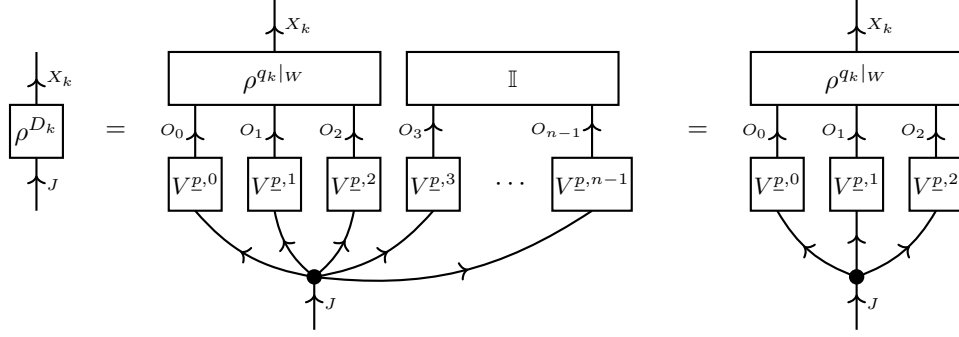


Figure 28: Generation of Datacores given a formula f , whose variables (here a_1, a_2, a_3) are selected from an extraction query p . Variables, which are not appearing in the formula f are trivialized over (here a_4, \dots, a_n). For consistency we denoted the index $x_{f(k)}$ by x_k .

11.4.9 Representation with auxiliary term variables

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 Section 8.2.4), treated as a categorical variable, with respect to the constant in the extraction query. The advantages are that we can avoid the ρ -formalims and directly model the categorical distributions.

11.4.10 Generic Tensor Network Decomposition of Extracted Data

More naturally: Formulas are compositions of predicates. When quantor free then a tensor network.

Remark 30 (Alternative Representation of empirical distributions). *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.*

11.4.11 Design of the Formulas

Most intuitive when labeling individuals by classes. Extraction formulas q_0, \dots, q_{d-1} can then be defined by subclasses of the member of a class and relations between objects of different classes. We then choose \mathcal{F} as more involved formulas decomposed into connectives acting on these atoms. The extraction query p is then designed based on class memberships to ensure, that the arguments of the formulas are always of specific classes.

We propose to

- Execute an extraction query to get pairs of individuals (the pairDf).
- Propositionalize the FOL Formulas independently on each tuple taking the individuals as a set of constant and filtering on the possible properties of each individuals. (Can understand as adding knowledge that most of the relations do not hold)
- Understand each such generated knowledge base as datapoint and average over them to get the empirical distribution to be fit.
- Fit a MLN describing the statistical relations of unseen results of the extraction query, based on likelihood maximation.

11.5 Generation of FOL worlds

Definition 67 (Reproduction of Empirical Distributions). *Given an empirical distribution $\mathbb{P}^D \in \bigotimes_{k \in [d]} \mathbb{R}^2$, we say that a triple $(W, p, q_{[d]})$ of a FOL world W an importance formula p and extraction formulas q_d reproduces \mathbb{P}^D , when*

$$\mathbb{P}^D = \left\langle \{p|_W\} \cup \{\rho^{q_k|_{k_g}} : k \in [d]\} \right\rangle [X_{[d]}|\emptyset] .$$

11.5.1 Samples by single objects

The first reproduction scheme identifies the datapoints with objects $A = [m]$

Theorem 76. *Given a dataset D the world $W[L, O]$*

$$W[L, O] = \sum_{k \in [d]} \sum_{j \in [m] : D_k(j)=1} e_k[L] \otimes e_j[O]$$

reproduces with the trivial importance query and extraction queries coinciding with the predicates the dataset D .

11.6 Samples by pairs of objects

We instantiate multiple objects for each datapoint, one for each variable of the importance formula, i.e. $A = [m] \times [n]$

11.7 Example: Generation of Knowledge Graphs

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.

Definition 68 (Reproduction of Empirical Distributions). *Given an empirical distribution $\mathbb{P}^D \in \bigotimes_{k \in [d]} \mathbb{R}^2$, we say that a tuple $(kg, p, \{q_0, \dots, q_{d-1}\})$ of a Knowledge Graph kg and queries p, q_k reproduces \mathbb{P}^D , when*

$$\mathbb{P}^D = \left\langle \{p|_{kg}\} \cup \{\rho^{q_k|_{kg}} : k \in [d]\} \right\rangle [X_{[d]}|\emptyset] .$$

Any distribution with rational coordinates is an empirical distribution, since each coordinate can be interpreted as the frequency of the respective world in the data D .

11.7.1 Samples by single resources

TBox: The categorical variables of the factored system are the coarsest classes. We define atomic formulas by the state indicators of each categorical variable as in Section 8.2.4. Each such atomic formula corresponds with a sub-class of the classes. By definition, each collection of state indicators define thus pairwise disjoint subclasses.

ABox: The samples are represented by single individuals in the Knowledge Graph. Their sub-class memberships corresponding with the categorical variables of the system are instantiated whenever the atom is true in the sample.

Theorem 77. *Let there any empirical distribution $\mathbb{P}^D \in \bigotimes_{k \in [d]} \mathbb{R}^2$ and $m \in \mathbb{N}$ such that $\text{im}(m \cdot \mathbb{P}^D) \subset \mathbb{N}$. Then the tuple $(kg, p, \{q_0, \dots, q_{d-1}\})$ defined by*

$$kg = \bigcup_{x_0, \dots, x_{d-1} \in \times_{k \in [d]} [2]} \{(s_{j, x_0, \dots, x_{d-1}} \text{ rdf:type } C) : j \in [m \cdot \mathbb{P}^D(x_0, \dots, x_{d-1})]\} \quad (59)$$

$$\bigcup_{x_0, \dots, x_{d-1} \in \times_{k \in [d]} [2]} \{(s_{j, x_0, \dots, x_{d-1}} \text{ rdf:type } C_k) : j \in [m \cdot \mathbb{P}^D(x_0, \dots, x_{d-1})], k \in [d], x_k = 1\} \quad (60)$$

$$p = \text{SELECT } \{ ?x \} \text{ WHERE } \{ ?x \text{ rdf:type } C . \}$$

$$q_k = \text{SELECT } \{ ?x \} \text{ WHERE } \{ ?x \text{ rdf:type } C_k . \}$$

reproduces \mathbb{P}^D .

Proof. With respect to any enumeration of the resources of kg we have

$$p|_{kg} = \sum_{x_0, \dots, x_{d-1} \in \times_{k \in [d]} [2]} \sum_{j \in [m \cdot \mathbb{P}^D(x_0, \dots, x_{d-1})]} e_{s_{j, x_0, \dots, x_{d-1}}} \quad (61)$$

and

$$q_k|_{kg} = \sum_{x_0, \dots, x_{d-1} \in \times_{k \in [d]} [2] : x_k = 1} \sum_{j \in [m \cdot \mathbb{P}^D(x_0, \dots, x_{d-1})]} e_{s_{j, x_0, \dots, x_{d-1}}} . \quad (62)$$

Summing over the resource variables of these tensors in a contraction we get

$$\left\langle \{p|_{kg}\} \cup \{\rho^{q_k|_{kg}} : k \in [d]\} \right\rangle [X_{[d]}] = \sum_{k \in [d]} m \cdot \mathbb{P}^D(x_0, \dots, x_{d-1}) \cdot e_{x_0, \dots, x_{d-1}} = m \cdot \mathbb{P}^D \quad (63)$$

and therefore

$$\left\langle \{p|_{kg}\} \cup \{\rho^{q_k|_{kg}} : k \in [d]\} \right\rangle [X_{[d]}|\emptyset] = \mathbb{P}^D. \quad (64)$$

□

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} .

11.7.2 Samples by pairs of resources

Remark 31 (Refinement of the Samples). *We can split each sample node into a pair of individuals. For this we need to specify, which each class membership will be encoded in a unary or binary attribute of the splitted individuals. This specification is possible based on the extraction query and the atomic formulas.*

Taking any importance query p , which has no permutation symmetries, we can instantiate each projection variable for each sample and prepare the links according to the triple patterns. When the atom queries q_0, \dots, q_{d-1} have different triple patterns compared with p , we instantiate those in cases where $x_k = 1$.

Label individuals $a_{j,l}$ by dataindex and variable index and

Theorem 78. *Let there be queries p, q_0, \dots, q_{d-1} , a Knowledge Graph kg containing individuals $a_{j,l}$ and a data map D . If*

$$p|_{kg} = \sum_j \bigotimes_{l \in [n]} e_{j,l}$$

and for any $k \in [d]$

$$q_k|_{kg} = \sum_{j: D^k(j)=1} \bigotimes_{l \in q_k} e_{j,l}.$$

Then the tuple $(kg, p, \{q_0, \dots, q_{d-1}\})$ reproduces \mathbb{P}^D .

Proof. It can be shown that the contraction is

$$\sum_j \bigotimes_{k \in [d]} e_{D^k(j)}.$$

□

Theorem 79. *The Knowledge Graph*

$$\text{RDF}|_{kg} = \sum_{t \in p} \sum_j t(a_{j,l}) + \sum_{k \in [d]} \sum_{t \in q_k \in [d]} \sum_{j: D^k(j)=1} t(a_{j,l})$$

reproduces the data, if the summed tensors are pairwise orthogonal. *Here we denote by $t(\cdot)$ the one-hot encoding of the tuple when mapping the corresponding individuals on the projection variables of the triple pattern t .*

11.7.3 General orthogonal projection approach

Given any triples, such that their composed projections to pairwise different triple patterns vanish. One reproducing KG is then the sum of the projection of any reproducing KG.

Theorem 80. *Let p, q_0, \dots, q_{d-1} contain of triple patterns such that for $t \neq \tilde{t}$*

$$P_t \circ P_{\tilde{t}} = 0.$$

For any KG reproducing the data D also the KG

$$\sum_t P_t \text{RDF}|_{kg}$$

is a reproducing KG.

Proof. By coinciding on all triple patterns t since

$$P_t \left(\sum_{\tilde{t}} P_{\tilde{t}} \text{RDF}|_{kg} \right) = P_t P_{\tilde{t}} \text{RDF}|_{kg} = P_t \text{RDF}|_{kg}.$$

Then the extraction contraction on kg and $\sum_t P_t \text{RDF}|_{kg}$ coincides. \square

From this perspective, it is obvious that any triple in the orthogonal complement of all projections can be added and the KG stays reproducing.

11.8 Discussion

Models of FOL are called Probabilistic Relational Models. Extensions are models that also handle structural uncertainty, i.e. distributions of worlds with varying A .

Statistical Models of Knowledge Graphs going beyond the typical single edge type perspective of network science Barabási (2016); ?.

Part III

Contraction Calculus

Based on the logical interpretation we often handle tensor calculus with specific tensors. Often, they are binary (that is their coordinates are in $\{0, 1\}$ corresponding with a Boolean), and sparse (that is having a decomposition with less storage demand). We investigate it in this part in more depth the properties of such tensors, which were exploited in the previous parts.

12 Coordinate Calculus

In the previous chapters, information to states has been stored in coordinates of a tensor. To distinguish from other schemes of calculus, we call this scheme of storing and retrieving information the coordinate calculus. We in this chapter investigate in more depth, which operations can be performed based on such tensors and proof the applied properties.

Add summations and further operations?

Discuss things like $\exp[T]$ as coordinatewise operations, by circles in the factor graph diagrams.

12.1 One-hot encodings as basis

Lemma 20 (Basis of tensor spaces). *The image of the one-hot encoding map is a linear basis of the tensor space $\bigotimes_{k \in [d]} \mathbb{R}^{m_k}$. Any element $f \in \bigotimes_{k \in [d]} \mathbb{R}^{m_k}$ has a decomposition*

$$f[X_0, \dots, X_{d-1}] = \sum_{x_0, \dots, x_{d-1}} f[X_0 = x_0, \dots, X_{d-1} = x_{d-1}] \cdot e_{x_0, \dots, x_{d-1}}[X_0, \dots, X_{d-1}].$$

We notice that the coordinates are the weights to the basis elements in the one-hot decomposition.

Proof. For any $\tilde{x}_0, \dots, \tilde{x}_{d-1} \in \bigtimes_{k \in [d]} [m_k]$ we have

$$f[X_{[d]} = \tilde{x}_{[d]}] = \sum_{x_0, \dots, x_{d-1}} f[X_0 = x_0, \dots, X_{d-1} = x_{d-1}] \cdot e_{x_0, \dots, x_{d-1}}[X_{[d]} = \tilde{x}_{[d]}] = f[X_{[d]} = \tilde{x}_{[d]}] \cdot e_{x_0, \dots, x_{d-1}}[X_{[d]} = \tilde{x}_{[d]}]$$

\square

Definition 69. Given any real-valued function

$$f : \bigtimes_{k \in [d]} [m_k] \rightarrow \mathbb{R}$$

we define the coordinate encoding by

$$T^f = \sum_{x_0, \dots, x_{d-1} \in \times_{k \in [d]} [m_k]} f(x_0, \dots, x_{d-1}) \cdot e_{x_0, \dots, x_{d-1}}.$$

Theorem 81 (Coordinate Calculus). *Given any tensor $T[X_0, \dots, X_{d-1}]$ can retrieve its coordinate indexed by x_0, \dots, x_{d-1} as*

$$T[X_0 = x_0, \dots, X_{d-1} = x_{d-1}] = \langle T, e_{x_0, \dots, x_{d-1}} \rangle [\emptyset].$$

Proof. We use the decomposition in Lemma 20 and have

$$\begin{aligned} \langle \{T, e_{x_0, \dots, x_{d-1}}\} \rangle [\emptyset] &= \sum_{\tilde{x}_0, \dots, \tilde{x}_{d-1}} T[X_0 = \tilde{x}_0, \dots, X_{d-1} = \tilde{x}_{d-1}] \cdot \langle \{e_{\tilde{x}_0, \dots, \tilde{x}_{d-1}}, e_{x_0, \dots, x_{d-1}}\} \rangle [X_0, \dots, X_{d-1}] \\ &= \sum_{\tilde{x}_0, \dots, \tilde{x}_{d-1}} T[X_0 = \tilde{x}_0, \dots, X_{d-1} = \tilde{x}_{d-1}] \cdot \delta_{(\tilde{x}_0, \dots, \tilde{x}_{d-1}), (x_0, \dots, x_{d-1})} \\ &= T[X_0 = x_0, \dots, X_{d-1} = x_{d-1}], \end{aligned}$$

where we used that one-hot encodings are orthonormal. \square

Coordinate calculus is the representation of real-valued functions as tensors, from which its evaluations can be retrieved by the scheme of Theorem 81.

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 contraction and in particular coordinate retrieval.

Theorem 82. *Given a tensor network $\mathcal{T}^{\mathcal{G}}$ on a hypergraph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$, disjoint subsets $A, B \subset \mathcal{V}$ and $x_B \in [m_B]$, we have*

$$\langle \mathcal{T}^{\mathcal{G}} \rangle [X_A, X_B = x_B] = \langle \{ \langle T^e \rangle [X_{e/B}, X_B = x_B] : e \in \mathcal{E} \} \rangle [X_A].$$

Proof. Using a delta tensor, which copies basis vector when contracting with one. \square

If we retrieve a single coordinate of a tensor, we have the situation $A = \emptyset, B = \mathcal{V}$. In that case, Theorem 82 shows, that the coordinate is the product of the coordinates of the cores.

12.2 Coordinatewise Transforms

This is central to coordinate calculus!

Definition 70. *Let $f : \mathbb{R}^p \rightarrow \mathbb{R}$ be a function. Then the coordinatewise transform of tensors $T^l[X_{[d]}]$, where $l \in [p]$, under f is the tensor*

$$f(T^0, \dots, T^{p-1})[X_{[d]}]$$

with coordinates

$$f(T^0, \dots, T^{p-1})[X_{[d]} = x_{[d]}] = f(T^0[X_{[d]} = x_{[d]}], \dots, T^{p-1}[X_{[d]} = x_{[d]}]).$$

To abbreviate notation, we will denote the transform by $f(T^0[X_{[d]}], \dots, T^{p-1}[X_{[d]}])$ and $f(T^0, \dots, T^{p-1})$.

Example 16 (Hadamard product). *Hadamard products of tensors are a special way of coordinate calculus, where the transform is the product and thus*

$$\cdot (T^0, \dots, T^{p-1})[X_{[d]}] = \langle \{T^l[X_{[d]}] : l \in [p]\} \rangle [X_{[d]}].$$

Applications of this are the computation of conjunctions, as we will exploit in effective calculus.

12.3 Differentiation of Contraction

We add additional variables Y selecting a coordinate of a tensor, which is varied in a differentiation.

Lemma 21. *For any tensor network $\mathcal{T}^{\mathcal{G}}$ with positive T^e we have*

$$\begin{aligned} \frac{\partial}{\partial T^e[Y_e]} \langle \mathcal{T}^{\mathcal{G}} \rangle [X_{\mathcal{V}} | \emptyset] &= \left\langle \delta[Y_e, X_e], \frac{\langle \mathcal{T}^{\mathcal{G}} \rangle [X_e]}{T^e[X_e]}, \langle \mathcal{T}^{\mathcal{G}} \rangle [X_{\mathcal{V}/e} | X_e] \right\rangle [Y_e, X_{\mathcal{V}}] \\ &\quad - \langle \mathcal{T}^{\mathcal{G}} \rangle [X_{\mathcal{V}} | \emptyset] \otimes \left\langle \frac{\langle \mathcal{T}^{\mathcal{G}} \rangle [Y_e]}{T^e[Y_e]} \right\rangle [Y_e]. \end{aligned}$$

Proof. By multilinearity of tensor network contractions we have

$$\frac{\partial}{\partial T^e[Y_e]} \langle \mathcal{T}^g \rangle [X_\nu] = \langle \{\delta[Y_e, X_e]\} \cup \{T^{\tilde{e}}[X_{\tilde{e}}] : \tilde{e} \neq e\} \rangle [Y_e, X_\nu]$$

and thus

$$\frac{\partial}{\partial T^e[Y_e]} \langle \mathcal{T}^g \rangle [\emptyset] = \langle \{\delta[Y_e, X_e]\} \cup \{T^{\tilde{e}}[X_{\tilde{e}}] : \tilde{e} \neq e\} \rangle [Y_e] .$$

Using both we get

$$\begin{aligned} \frac{\partial}{\partial T^e[Y_e]} \langle \mathcal{T}^g \rangle [X_\nu | \emptyset] &= \frac{\partial}{\partial T^e[Y_e]} \frac{\langle \mathcal{T}^g \rangle [X_\nu]}{\langle \mathcal{T}^g \rangle [\emptyset]} \\ &= \frac{\frac{\partial}{\partial T^e[Y_e]} \langle \mathcal{T}^g \rangle [X_\nu]}{\langle \mathcal{T}^g \rangle [\emptyset]} - \frac{\langle \mathcal{T}^g \rangle [X_\nu] \frac{\partial}{\partial T^e[Y_e]} \langle \mathcal{T}^g \rangle [\emptyset]}{(\langle \mathcal{T}^g \rangle [\emptyset])^2} \\ &= \frac{\langle \{\delta[Y_e, X_e]\} \cup \{T^{\tilde{e}}[X_{\tilde{e}}] : \tilde{e} \neq e\} \rangle [Y_e, X_\nu]}{\langle \mathcal{T}^g \rangle [\emptyset]} \\ &\quad - \langle \mathcal{T}^g \rangle [X_\nu | \emptyset] \cdot \frac{\langle \{\delta[Y_e, X_e]\} \cup \{T^{\tilde{e}}[X_{\tilde{e}}] : \tilde{e} \neq e\} \rangle [Y_e]}{\langle \mathcal{T}^g \rangle [\emptyset]} \end{aligned} \quad (65)$$

For the first term we get with a normation equation (see Theorem 88) that

$$\begin{aligned} \frac{\langle \{\delta[Y_e, X_e]\} \cup \{T^{\tilde{e}}[X_{\tilde{e}}] : \tilde{e} \neq e\} \rangle [Y_e, X_\nu]}{\langle \mathcal{T}^g \rangle [\emptyset]} &= \frac{\langle \{\delta[Y_e, X_e]\} \cup \{T^{\tilde{e}}[X_{\tilde{e}}] : \tilde{e} \in \mathcal{E}\} \rangle [Y_e, X_\nu]}{T^e[X_e] \cdot \langle \mathcal{T}^g \rangle [\emptyset]} \\ &= \frac{\langle \delta[Y_e, X_e], \langle \mathcal{T}^g \rangle [X_\nu | \emptyset] \rangle [Y_e, X_\nu]}{T^e[X_e]} \\ &= \frac{\langle \delta[Y_e, X_e], \langle \mathcal{T}^g \rangle [X_e | \emptyset], \langle \mathcal{T}^g \rangle [X_{\nu/e} | X_e] \rangle [Y_e, X_\nu]}{T^e[X_e]} . \end{aligned}$$

Analogously, we have

$$\frac{\langle \{\delta[Y_e, X_e]\} \cup \{T^{\tilde{e}}[X_{\tilde{e}}] : \tilde{e} \neq e\} \rangle [Y_e]}{\langle \mathcal{T}^g \rangle [\emptyset]} = \frac{\langle \delta[Y_e, X_e], \langle \mathcal{T}^g \rangle [X_e | \emptyset] \rangle [Y_e]}{T^e[X_e]} .$$

With (65), we arrive at the claim

$$\begin{aligned} \frac{\partial}{\partial T^e[Y_e]} \langle \mathcal{T}^g \rangle [X_\nu | \emptyset] &= \left\langle \delta[Y_e, X_e], \frac{\langle \mathcal{T}^g \rangle [X_e]}{T^e[X_e]}, \langle \mathcal{T}^g \rangle [X_{\nu/e} | X_e] \right\rangle [Y_e, X_\nu] \\ &\quad - \langle \mathcal{T}^g \rangle [X_\nu | \emptyset] \otimes \left\langle \frac{\langle \mathcal{T}^g \rangle [Y_e]}{T^e[Y_e]} \right\rangle [Y_e] . \end{aligned}$$

□

Lemma 22. For any function $f(T^e)[X_\nu]$ we have

$$\begin{aligned} \frac{\partial}{\partial T^e[Y_e]} \langle \langle \mathcal{T}^g \rangle [X_\nu | \emptyset], f(T^e)[X_\nu] \rangle [\emptyset] &= \frac{\langle \mathcal{T}^g \rangle [X_e = x_e | \emptyset]}{T^e[X_e = x_e]} \left(\langle \langle \mathcal{T}^g \rangle [X_{\nu/e} | X_e = x_e], f(T^e)[X_\nu, Y_e] \rangle [\emptyset] \right. \\ &\quad \left. - \langle \langle \mathcal{T}^g \rangle [X_\nu | \emptyset], f(T^e)[X_\nu] \rangle [\emptyset] \right) \\ &\quad + \left\langle \langle \mathcal{T}^g \rangle [X_\nu | \emptyset] \frac{\partial f(T^e)[X_\nu]}{\partial T^e[Y_e]} \right\rangle [\emptyset] \end{aligned}$$

Proof. By product rule of differentiation we have

$$\begin{aligned} \frac{\partial}{\partial T^e [X_e = x_e]} \langle \langle \mathcal{T}^g \rangle [X_\nu | \emptyset], f(T^e)[X_\nu] \rangle [\emptyset] &= \left\langle \frac{\partial}{\partial T^e [X_e = x_e]} \langle \mathcal{T}^g \rangle [X_\nu | \emptyset], f(T^e)[X_\nu] \right\rangle [\emptyset] \\ &+ \left\langle \langle \mathcal{T}^g \rangle [X_\nu | \emptyset], \frac{\partial}{\partial T^e [X_e = x_e]} f(T^e)[X_\nu] \right\rangle [\emptyset]. \end{aligned}$$

The claim now follows with the application of Lemma 21 on the first term. \square

12.4 Selection Encodings

Selection encodings as introduced in Definition 15 are best understood in terms of linear mapping interpretations of tensors. We will first provide by basis encodings 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.

12.4.1 Tensors as linear maps

The state sets $\times_{k \in [d]} [m_k]$ can be interpreted as an enumeration of basis elements e_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 $\times_{k \in [d]} [m_k]$ are an intermediate representation.

Definition 71. Let there be two tensor spaces V_1 and V_2 with basis by sets $\mathcal{M}_1 \subset V_1$ and $\mathcal{M}_2 \subset V_2$ of cardinality m_1 and m_2 , which are enumerated by variables O_1 and O_2 . The basis encoding of a linear map $f \in \mathbb{L}(V_1, V_2)$ is the tensor

$$\beta^f [O_1, O_2] \in \mathbb{R}^{m_1} \otimes \mathbb{R}^{m_2}.$$

Basis encodings are standard linear algebra tools, where matrices are understood as linear maps between vector spaces.

Theorem 83. 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 chosen bases in the spaces. We have

$$\beta^{F^2 \circ F^1} [O_1, O_3] = \left\langle \beta^{F^2} [O_2, O_3], \beta^{F^1} [O_1, O_2] \right\rangle [O_1, O_3].$$

Proof. By basis decompositions and representations of linear maps as contractions. \square

A typical instance is matrix multiplication, where matrices understood as representations of linear maps.

12.4.2 Selection encodings

Selection encodings (see Definition 15) are related to basis encodings of linear maps as we show in the next theorem.

Theorem 84. Given a function

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

we define a linear map $F^f \in \mathbb{L}(\bigotimes_{k \in [d]} \mathbb{R}^{m_k}, \bigotimes_{s \in [n]} \mathbb{R}^{p_s})$ by the basis elements to $x_1 \in \times_{k \in [d]} [m_k]$ and $x_2 \in \times_{s \in [n]} [p_s]$ by

$$F^f(e_{x_1}, e_{x_2}) = \langle f(e_{x_1}), e_{x_2} \rangle [\emptyset].$$

We then have

$$\gamma^f = \beta^{F^f}.$$

While relational encoding works for maps from $\times_{k \in [d]} [m_k]$ to arbitrary sets (which are enumerated), selection encodings as introduced in Definition 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

$$f(x) = \gamma^f [X = x, L] .$$

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.

Lemma 23. *Let $T [X_{\mathcal{V}}]$ be a tensor in $\otimes_{v \in \mathcal{V}} \mathbb{R}^{m_v}$ and let A, B be a disjoint partition of \mathcal{V} , that is $A \dot{\cup} B = \mathcal{V}$. Then the function*

$$f : \times_{v \in A} [m_v] \rightarrow \otimes_{v \in B} \mathbb{R}^{m_v}$$

with coordinates

$$f(x_A) = T [X_A = x_A, X_B]$$

obeys

$$\gamma^f [X_{\mathcal{V}}] = T [X_{\mathcal{V}}] .$$

Here we have renamed the selection variables L_A by the categorical variables X_A .

Proof. From Theorem 83 using the basis encoding equivalence of Theorem 84. □

Example 17 (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 Lemma 23. We can construct \mathcal{H}^T by first building the selection encoding to \mathcal{H} and then applying the construction of Lemma 23 with $A = L$ and $B = X_{[d]}$.*

We use selection encodings to express compositions of functions, based on the next theorem.

Theorem 85 (Selection Encoding for Linear Compositions). *Let ϕ be a tensor valued function from $\times_{k \in [d]} [m_k]$ to \mathbb{R}^p with image coordinates ϕ_l and let f be a tensor Then*

$$\left(\sum_{l \in [p]} f[L_{\phi} = l] \cdot \phi_l \right) [X_{[d]}] : \times_{k \in [d]} [m_k] \rightarrow \mathbb{R}$$

is represented as

$$\left(\sum_{l \in [p]} f[L_{\phi} = l] \cdot \phi_l \right) [X_{[d]}] = \langle \gamma^{\phi} [X_{[d]}, L_{\phi}], f[L_{\phi}] \rangle [X_{[d]}] .$$

Proof. The representation holds, since for any $x_{[d]} \in \times_{k \in [d]} [m_k]$ we have

$$\langle \gamma^{\phi} [X_{[d]}, L_{\phi}], f[L_{\phi}] \rangle [X_{[d]} = x_{[d]}] = \sum_{l \in [p]} f[L_{\phi} = l] \cdot \phi_l(x_{[d]}).$$

□

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 ϕ and f in Theorem 85 as basis decompositions of linear maps.

12.5 Discussion

Representations of linear maps is the typical application of tensors, reason for refering to tensor networks as multilinear algebra.

13 Directed Tensor Calculus

We in this chapter investigate the properties of tensors, which where arising in probabilistic and logical reasoning. We observed already before, that

- Conditional probability tensors are directed tensors.
- Logical formulas are binary tensors.

Thus, the set of tensors, which are both directed and binary is of much interest. We will show, that they are equal to the set of relational encodings of functions.

13.1 Directed Tensors

Directionality as defined in Definition 12 represents the constraints on the structure of tensors: Summing over outgoing trivializes the tensor.

Definition 72 (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 T^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, which an extreme example being the Dirac Delta Tensors. Further example are relational encodings of invertible functions.

Example 18 (Dirac Delta Tensors). *Given a node set e colored with a constant dimension m Diracs Delta Tensors are the tensors*

$$\delta^{e,m} \in \bigotimes_{v \in e} \mathbb{R}^m$$

with coordinates

$$\delta_{x_0, \dots, x_{d-1}}^{e,m} = \begin{cases} 1 & \text{if } x_0 = \dots = x_{d-1} \\ 0 & \text{else} \end{cases} . \quad (66)$$

The contractions with respect to subsets $\tilde{\mathcal{V}} \subset e$ are

$$\langle \{\delta^{e,m}\} \rangle [X_{\tilde{\mathcal{V}}}] = \begin{cases} m & \text{if } \mathcal{V} = \emptyset \\ \mathbb{I} & \text{if } |\tilde{\mathcal{V}}| = 1 \\ \delta^{\tilde{\mathcal{V}},m} & \text{else} \end{cases} . \quad (67)$$

Thus are directed for any orientation of the respective edge with exactly one incoming variable.

Lemma 24. *Let $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ be a hypergraph and $\mathcal{T}^{\mathcal{G}}$ 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 $\mathcal{T}^{\tilde{\mathcal{G}}}$ by*

- *Recolored Edges $\tilde{\mathcal{E}} = \{\tilde{e} : e \in \mathcal{E}\}$ where $\tilde{e} = \{v^e : v \in e\}$, which decoration tensor $T^{\tilde{e}}$ has same coordinates as T^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

$$\langle \mathcal{T}^{\mathcal{G}} \rangle [X_{\mathcal{V}}] = \langle \mathcal{T}^{\tilde{\mathcal{G}}} \rangle [X_{\mathcal{V}}] .$$

Proof. For any $x_{\mathcal{V}}$ we have

$$\begin{aligned} \langle \mathcal{T}^{\tilde{\mathcal{G}}} \rangle [X_{\mathcal{V}} = x_{\mathcal{V}}] &= \left\langle \{T^{\tilde{e}}[X_{\{v^e : v \in e\}}] : e \in \mathcal{E}\} \cup \{\delta^{\{v\} \cup \{v^e : e \ni v\}}[X_{\{v^e : e \ni v\}}, X_v = x_v] : v \in \mathcal{V}\} \right\rangle [\emptyset] \\ &= \langle \{T^{\tilde{e}}[X_{\{v^e : v \in e\}}] = x_{\{v : v \in e\}} : e \in \mathcal{E}\} \rangle [\emptyset] \\ &= \langle \mathcal{T}^{\mathcal{G}} \rangle [X_{\mathcal{V}} = x_{\mathcal{V}}] , \end{aligned}$$

which establishes the claim. \square

13.1.1 Normation

Normed tensors (see Definition 13) are directed and directed tensors invariant under normation wrt their incoming and outgoing variable, as we show next.

Theorem 86. *For any tensor network $\mathcal{T}^{\mathcal{G}}$ on variables \mathcal{V} that can be normed with respect to \mathcal{V}^{in} and \mathcal{V}^{out} , the normation is directed with \mathcal{V}^{in} incoming and \mathcal{V}^{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

$$\langle \langle \mathcal{T}^{\mathcal{G}} \rangle [\mathcal{V}^{\text{in}} | \mathcal{V}^{\text{out}}], e_{x_{\mathcal{V}^{\text{in}}}} \rangle [\emptyset] = \frac{\langle \mathcal{T}^{\mathcal{G}} \cup \{e_{x_{\mathcal{V}^{\text{in}}}}\} \rangle [\emptyset]}{\langle \mathcal{T}^{\mathcal{G}} \cup \{e_{x_{\mathcal{V}^{\text{in}}}}\} \rangle [\emptyset]}.$$

By Definition 12, $\langle \mathcal{T}^{\mathcal{G}} \rangle [\mathcal{V}^{\text{out}} | \mathcal{V}^{\text{in}}]$ is thus directed. \square

The normation operation coincides in cases of non-negative tensors with the conditioning of a Markov Network representing a probability distribution.

13.1.2 Contraction of Directed Tensors

Let us now investigate, which contractions inherit the directionality of the tensors.

Theorem 87. *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}^{in} as those nodes, which do not appear as outgoing variables. For any tensor network $\mathcal{T}^{\mathcal{G}}$ respecting the direction of \mathcal{G} we have that*

$$\langle \mathcal{T}^{\mathcal{G}} \rangle [X_{\mathcal{V}^{\text{in}}}] = \mathbb{I} [X_{\mathcal{V}^{\text{in}}}] ,$$

that is $\langle \mathcal{T}^{\mathcal{G}} \rangle [X_{\mathcal{V}}]$ is a directed tensor with \mathcal{V}^{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 Lemma 24 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. When orienting each edge

We show the theorem over induction on the number n of cores.

$n = 1$: It holds trivially, when $\mathcal{T}^{\mathcal{G}}$ consists of a single core

$n \rightarrow n - 1$: Let us assume, that the claim holds for graphs with $n - 1$ hyperedges and let $\mathcal{T}^{\mathcal{G}}$ 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 91 and get

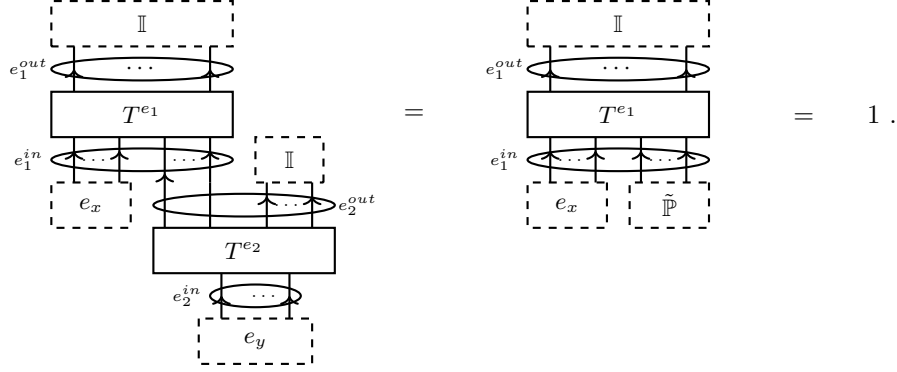
$$\begin{aligned} \langle \mathcal{T}^{\mathcal{G}} \rangle [X_{\mathcal{V}^{\text{in}}}] &= \left\langle \mathcal{T}^{(\mathcal{V}, \mathcal{E}/\{e\})} \cup \{T^e [X_{e^{\text{in}}}, X_{e^{\text{out}}}] \} \right\rangle [X_{\mathcal{V}^{\text{in}}}] \\ &= \left\langle \mathcal{T}^{(\mathcal{V}, \mathcal{E}/\{e\})} \cup \{ \langle T^e \rangle [X_{e^{\text{in}}}] \} \right\rangle [X_{\mathcal{V}^{\text{in}}}] \\ &= \left\langle \mathcal{T}^{(\mathcal{V}, \mathcal{E}/\{e\})} \cup \{ \mathbb{I} [X_{e^{\text{in}}}] \} \right\rangle [X_{\mathcal{V}^{\text{in}}}] \\ &= \left\langle \mathcal{T}^{(\mathcal{V}, \mathcal{E}/\{e\})} \right\rangle [X_{\mathcal{V}^{\text{in}}}] . \end{aligned}$$

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

$$\langle \mathcal{T}^{\mathcal{G}} \rangle [X_{\mathcal{V}^{\text{in}}}] = \left\langle \mathcal{T}^{(\mathcal{V}, \mathcal{E}/\{e\})} \right\rangle [X_{\mathcal{V}^{\text{in}}}] = \mathbb{I} [X_{\mathcal{V}^{\text{in}}}] .$$

\square

Schematically, the direction conservation property argument is depicted as:



14 Contraction Equations

We have observed, that many concepts and theorems in probability theory and logics can be understood as contraction equations. We first provide a summary of the used contraction equations.

14.1 Contraction equations in logical and probabilistic reasoning

Let us summarize the application of contractions and normation in the definition of

- Marginal probabilities (Definition 19, Theorem 1)

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

- Conditional probabilities (Definition 20, Theorem 2)

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

- The partition function of a Markov Networks

$$\mathcal{Z}(\mathcal{T}^G) = \langle \mathcal{T}^G \rangle [\emptyset]$$

- The probability distribution of a Markov Network is (Definition 23)

$$\mathbb{P}^{\mathcal{T}^G} = \langle \mathcal{T}^G \rangle [\mathcal{V}|\emptyset]$$

Further the following properties have been defined by contraction equations:

- X_0 and X_1 are independent when (Definition 21, Theorem 6)

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

- X_0 and X_1 are called independent conditioned on X_2 when (Definition 22, Theorem 7)

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

14.2 Normation Equations

Theorem 88 (Normation as a Contraction Equation). *For any on \mathcal{V}^{in} normable tensor $T[X_{\mathcal{V}}]$, where $\mathcal{V}^{in} \dot{\cup} \mathcal{V}^{out} = \mathcal{V}$, we have*

$$\langle T \rangle [X_{\mathcal{V}}] = \langle \langle T \rangle [X_{\mathcal{V}^{out}}|X_{\mathcal{V}^{in}}], \langle T \rangle [X_{\mathcal{V}^{in}}] \rangle [X_{\mathcal{V}}] .$$

Proof. Let us choose indices $x_{\mathcal{V}^{in}}$ and $x_{\mathcal{V}^{out}}$. We have that

$$\langle T \rangle [X_{\mathcal{V}^{in}} = x_{\mathcal{V}^{in}}|X_{\mathcal{V}^{out}} = x_{\mathcal{V}^{out}}] = \frac{\langle T \rangle [X_{\mathcal{V}^{in}} = x_{\mathcal{V}^{in}}, X_{\mathcal{V}^{out}} = x_{\mathcal{V}^{out}}]}{\langle T \rangle [X_{\mathcal{V}^{in}} = x_{\mathcal{V}^{in}}]}$$

and therefor

$$\langle T \rangle [X_{\mathcal{V}^{in}} = x_{\mathcal{V}^{in}}, X_{\mathcal{V}^{out}} = x_{\mathcal{V}^{out}}] = \langle T \rangle [X_{\mathcal{V}^{in}} = x_{\mathcal{V}^{in}}|X_{\mathcal{V}^{out}} = x_{\mathcal{V}^{out}}] \cdot \langle T \rangle [X_{\mathcal{V}^{in}} = x_{\mathcal{V}^{in}}]$$

Since the equation holds for arbitrary indices, the claim is established. \square

Theorem 89 (Generic Chain Rule). *For any Tensor $T[X_{\mathcal{V}}]$ and any total order \prec on the nodes \mathcal{V} we have*

$$T[X_{\mathcal{V}}] = \langle \{ \langle T \rangle [X_v | X_{\{\tilde{v} : \tilde{v} \prec v, \tilde{v} \neq v\}}] : v \in \mathcal{V} \} \rangle [X_{\mathcal{V}}]$$

Proof. We apply Theorem 88 on the tensor

$$\langle T \rangle [X_v, X_{\{\tilde{v} : v \prec \tilde{v}, \tilde{v} \neq v\}} | X_{\{\tilde{v} : \tilde{v} \prec v, \tilde{v} \neq v\}} = x_{\{\tilde{v} : \tilde{v} \prec v, \tilde{v} \neq v\}}],$$

where $v \in \mathcal{V}$ and $x_{\mathcal{V}}$ are chosen arbitrarily. For any $v \in \mathcal{V}$ we get

$$\langle T \rangle [X_v, X_{\{\tilde{v} : v \prec \tilde{v}, \tilde{v} \neq v\}} | X_{\{\tilde{v} : \tilde{v} \prec v, \tilde{v} \neq v\}}] = \langle \langle T \rangle [X_{\{\tilde{v} : v \prec \tilde{v}, \tilde{v} \neq v\}} | X_v, X_{\{\tilde{v} : \tilde{v} \prec v, \tilde{v} \neq v\}}], \langle T \rangle [X_v | X_{\{\tilde{v} : \tilde{v} \prec v, \tilde{v} \neq v\}}] \rangle [X_{\mathcal{V}}].$$

Applying this equation iteratively and making use of the commutation of contractions we get for any $v \in \mathcal{V}$

$$\langle T \rangle [X_v, X_{\{\tilde{v} : v \prec \tilde{v}, \tilde{v} \neq v\}} | X_{\{\tilde{v} : \tilde{v} \prec v, \tilde{v} \neq v\}}] = \langle \langle T \rangle [X_{\tilde{v}} | X_{\{\tilde{v} : \tilde{v} \prec v, \tilde{v} \neq v\}}] : v \prec \tilde{v} \rangle [X_{\mathcal{V}}].$$

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

14.3 Proof of Hammersley-Clifford Theorem

Different to the standard case, we show a version of Hammersley-Clifford for hypergraphs.

Lemma 25. *Let $T[X_{\mathcal{V}}]$ be a positive tensor and $y_{\mathcal{V}}$ an arbitrary index. Then we have*

$$T[X_{\mathcal{V}}] = \left\langle \left(\langle T \rangle [X_{\mathcal{V}/\bar{\mathcal{V}}}, X_{\bar{\mathcal{V}}} = y_{\bar{\mathcal{V}}}] \right)^{(-1)^{|\bar{\mathcal{V}}| - |\mathcal{V}|}} : \bar{\mathcal{V}} \subset \tilde{\mathcal{V}} \subset \mathcal{V} \right\rangle [X_{\mathcal{V}}],$$

where the exponentiation is performed coordinatewise and positivity of T ensures the well-definedness.

Proof. It suffices to show, that for an arbitrary index $x_{\mathcal{V}}$ be an arbitrary index we have

$$T[X_{\mathcal{V}} = x_{\mathcal{V}}] = \prod_{\tilde{\mathcal{V}} \subset \mathcal{V}} \prod_{\bar{\mathcal{V}} \subset \tilde{\mathcal{V}}} \left(\langle T \rangle [X_{\mathcal{V}/\bar{\mathcal{V}}}, X_{\bar{\mathcal{V}}} = y_{\bar{\mathcal{V}}}] \right)^{(-1)^{|\bar{\mathcal{V}}| - |\mathcal{V}|}}.$$

We do this by applying a logarithm on the right hand side and grouping the terms by $\bar{\mathcal{V}}$ as

$$\begin{aligned} & \ln \left[\prod_{\tilde{\mathcal{V}} \subset \mathcal{V}} \prod_{\bar{\mathcal{V}} \subset \tilde{\mathcal{V}}} \left(\langle T \rangle [X_{\mathcal{V}/\bar{\mathcal{V}}}, X_{\bar{\mathcal{V}}} = y_{\bar{\mathcal{V}}}] \right)^{(-1)^{|\bar{\mathcal{V}}| - |\mathcal{V}|}} \right] \\ &= \sum_{\tilde{\mathcal{V}} \subset \mathcal{V}} \ln [\langle T \rangle [X_{\mathcal{V}/\bar{\mathcal{V}}}, X_{\bar{\mathcal{V}}} = y_{\bar{\mathcal{V}}}] \left(\sum_{\tilde{\mathcal{V}} \subset \mathcal{V} : \bar{\mathcal{V}} \subset \tilde{\mathcal{V}}} (-1)^{|\bar{\mathcal{V}}| - |\mathcal{V}|} \right) \\ &= \sum_{\tilde{\mathcal{V}} \subset \mathcal{V}} \ln [\langle T \rangle [X_{\mathcal{V}/\bar{\mathcal{V}}}, X_{\bar{\mathcal{V}}} = y_{\bar{\mathcal{V}}}] \left(\sum_{i \in [|\mathcal{V}| - |\bar{\mathcal{V}}|]} (-1)^i \binom{|\mathcal{V}| - |\bar{\mathcal{V}}|}{i} \right) \end{aligned}$$

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

$$\begin{aligned} & \ln \left[\prod_{\tilde{\mathcal{V}} \subset \mathcal{V}} \prod_{\bar{\mathcal{V}} \subset \tilde{\mathcal{V}}} \left(\langle T \rangle [X_{\mathcal{V}/\bar{\mathcal{V}}}, X_{\bar{\mathcal{V}}} = y_{\bar{\mathcal{V}}}] \right)^{(-1)^{|\bar{\mathcal{V}}| - |\mathcal{V}|}} \right] \\ &= \ln [T[X_{\mathcal{V}} = x_{\mathcal{V}}]] \left(\sum_{i \in [0]} (-1)^i \binom{0}{i} \right) \\ &= \ln [T[X_{\mathcal{V}} = x_{\mathcal{V}}]] . \end{aligned}$$

Applying the exponential function on both sides establishes the claim. \square

Lemma 26. *Let T 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*

$$\langle T \rangle [X_{A,B} | X_{\mathcal{V}/\{A,B\}}] = \langle \langle T \rangle [X_A | X_{\mathcal{V}/\{A,B\}}], \langle T \rangle [X_B | X_{\mathcal{V}/\{A,B\}}] \rangle [X_{\tilde{\mathcal{V}}}]$$

then

$$\prod_{\tilde{\mathcal{V}} \subset \tilde{\mathcal{V}}} (\langle T \rangle [X_{\mathcal{V}/\tilde{\mathcal{V}}} = x_{\mathcal{V}/\tilde{\mathcal{V}}}, X_{\tilde{\mathcal{V}}} = y_{\tilde{\mathcal{V}}}])^{(-1)^{|\tilde{\mathcal{V}}| - |\tilde{\mathcal{V}}|}} = 1.$$

Proof. We abbreviate

$$Z_{\tilde{\mathcal{V}}} = \langle T \rangle [X_{\mathcal{V}/\tilde{\mathcal{V}}} = x_{\mathcal{V}/\tilde{\mathcal{V}}}, X_{\tilde{\mathcal{V}}} = y_{\tilde{\mathcal{V}}}] .$$

By reorganizing the sum over $\tilde{\mathcal{V}} \subset \tilde{\mathcal{V}}$ into $\tilde{\mathcal{V}} \subset \tilde{\mathcal{V}}/A \cup B$ we have

$$\prod_{\tilde{\mathcal{V}} \subset \tilde{\mathcal{V}}} (Z_{\tilde{\mathcal{V}}})^{(-1)^{|\tilde{\mathcal{V}}| - |\tilde{\mathcal{V}}|}} = \prod_{\tilde{\mathcal{V}} \subset \tilde{\mathcal{V}}/\{A,B\}} \left(\frac{Z_{\tilde{\mathcal{V}}} \cdot Z_{\tilde{\mathcal{V}} \cup \{A,B\}}}{Z_{\tilde{\mathcal{V}} \cup \{A\}} \cdot Z_{\tilde{\mathcal{V}} \cup \{B\}}} \right)^{(-1)^{|\tilde{\mathcal{V}}| - |\tilde{\mathcal{V}}|}} . \quad (68)$$

From the independence assumption it follows that for any index x

$$\begin{aligned} \langle T \rangle [X_A = x_A | X_{\mathcal{V}/\tilde{\mathcal{V}} \cup \{A,B\}} = x_{\mathcal{V}/\tilde{\mathcal{V}} \cup \{A,B\}}, X_{\tilde{\mathcal{V}}} = y_{\tilde{\mathcal{V}}}, X_B = x_B] \\ = \langle T \rangle [X_A = x_A | X_{\mathcal{V}/\tilde{\mathcal{V}} \cup \{A,B\}} = x_{\mathcal{V}/\tilde{\mathcal{V}} \cup \{A,B\}}, X_{\tilde{\mathcal{V}}} = y_{\tilde{\mathcal{V}}}] \\ = \langle T \rangle [X_A = x_A | X_{\mathcal{V}/\tilde{\mathcal{V}} \cup \{A,B\}} = x_{\mathcal{V}/\tilde{\mathcal{V}} \cup \{A,B\}}, X_{\tilde{\mathcal{V}}} = y_{\tilde{\mathcal{V}}}, X_B = y_B] \end{aligned}$$

Applying this in each squares bracket term of (68) we get

$$\begin{aligned} \frac{Z_{\tilde{\mathcal{V}}}}{Z_{\tilde{\mathcal{V}} \cup \{A\}}} &= \frac{\langle T \rangle [X_A = x_A | X_{\mathcal{V}/\tilde{\mathcal{V}} \cup \{A,B\}} = x_{\mathcal{V}/\tilde{\mathcal{V}} \cup \{A,B\}}, X_{\tilde{\mathcal{V}}} = y_{\tilde{\mathcal{V}}}, X_B = x_B]}{\langle T \rangle [X_A = y_A | X_{\mathcal{V}/\tilde{\mathcal{V}} \cup \{A,B\}} = x_{\mathcal{V}/\tilde{\mathcal{V}} \cup \{A,B\}}, X_{\tilde{\mathcal{V}}} = y_{\tilde{\mathcal{V}}}, X_B = x_B]} \\ &= \frac{\langle T \rangle [X_A = x_A | X_{\mathcal{V}/\tilde{\mathcal{V}} \cup \{A,B\}} = x_{\mathcal{V}/\tilde{\mathcal{V}} \cup \{A,B\}}, X_{\tilde{\mathcal{V}}} = y_{\tilde{\mathcal{V}}}, X_B = y_B]}{\langle T \rangle [X_A = y_A | X_{\mathcal{V}/\tilde{\mathcal{V}} \cup \{A,B\}} = x_{\mathcal{V}/\tilde{\mathcal{V}} \cup \{A,B\}}, X_{\tilde{\mathcal{V}}} = y_{\tilde{\mathcal{V}}}, X_B = y_B]} \\ &= \frac{Z_{\tilde{\mathcal{V}} \cup \{B\}}}{Z_{\tilde{\mathcal{V}} \cup \{A,B\}}} . \end{aligned}$$

Thus, each factor in (68) is trivial, which establishes the claim. \square

We are finally ready to proof Theorem 9 based on the Lemmata above.

Theorem 90 (Theorem 9). *Let $\mathbb{P}[X_{\mathcal{V}}]$ be a probability distribution and \mathcal{G} a clique-capturing hypergraph, such that for A, B, C we have that X_A is independent of X_B conditioned on X_C , when C separates A and B in the hypergraph. Then there is a Markov Network on \mathcal{G} , which distribution is equal to $\mathbb{P}[X_{\mathcal{V}}]$.*

Proof of Theorem 9. By Lemma 25 we have for any index $x_{\mathcal{V}}$

$$\mathbb{P}[X_{\mathcal{V}} = x_{\mathcal{V}}] = \prod_{\tilde{\mathcal{V}} \subset \mathcal{V}} \prod_{\tilde{\mathcal{V}} \subset \tilde{\mathcal{V}}} (\mathbb{P}[X_{\tilde{\mathcal{V}}} = x_{\tilde{\mathcal{V}}}, X_{\mathcal{V}/\tilde{\mathcal{V}}} = y_{\mathcal{V}/\tilde{\mathcal{V}}}])^{(-1)^{|\tilde{\mathcal{V}}| - |\tilde{\mathcal{V}}|}}$$

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 Lemma 26 we then have

$$\prod_{\tilde{\mathcal{V}} \subset \tilde{\mathcal{V}}} (\mathbb{P}[X_{\tilde{\mathcal{V}}} = x_{\tilde{\mathcal{V}}}, X_{\mathcal{V}/\tilde{\mathcal{V}}} = y_{\mathcal{V}/\tilde{\mathcal{V}}}])^{(-1)^{|\tilde{\mathcal{V}}| - |\tilde{\mathcal{V}}|}} = 1.$$

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

$$T^e[X_e] = \prod_{\tilde{v} : \alpha(\tilde{v})=e} \prod_{\tilde{v} \subset \tilde{v}} (\mathbb{P}[X_{\tilde{v}} = x_{\tilde{v}}, X_{\mathcal{V}/\tilde{v}} = y_{\mathcal{V}/\tilde{v}}])^{(-1)^{|\tilde{v}|-|\tilde{v}|}}.$$

and get, that

$$\begin{aligned} \mathbb{P}[X_{\mathcal{V}}] &= \langle \{T^e : e \in \mathcal{E}\} \rangle [X_{\mathcal{V}}] \\ &= \langle \{T^e : e \in \mathcal{E}\} \rangle [X_v | \emptyset]. \end{aligned}$$

We have thus constructed a Markov Network with trivial partition function, which contraction coincides with the probability distribution. \square

14.4 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 91. *Let $\mathcal{T}^{\mathcal{G}}$ 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*

$$\langle \mathcal{T}^{\mathcal{G}} \rangle [X_{\tilde{v}}] = \left\langle \mathcal{T}^{\mathcal{G}_1} \cup \{ \langle \mathcal{T}^{\mathcal{G}_2} \rangle [X_{\mathcal{V}_2 \cap (\mathcal{V}_1 \cup \tilde{v})}] \} \right\rangle [X_{\tilde{v}}].$$

Proof. For any index $x_{\tilde{v}}$ we show that

$$\langle \mathcal{T}^{\mathcal{G}} \rangle [X_{\tilde{v}} = x_{\tilde{v}}] = \left\langle \mathcal{T}^{\mathcal{G}_1} \cup \{ \langle \mathcal{T}^{\mathcal{G}_2} \rangle [X_{\mathcal{V}_2 \cap (\mathcal{V}_1 \cup \tilde{v})}] \} \right\rangle [X_{\tilde{v}} = x_{\tilde{v}}].$$

By definition we have

$$\begin{aligned} \langle \mathcal{T}^{\mathcal{G}} \rangle [X_{\tilde{v}} = x_{\tilde{v}}] &= \sum_{x_{\mathcal{V}/\tilde{v}}} \prod_{e \in \mathcal{E}} T^e[X_e = x_e] \\ &= \sum_{x_{\mathcal{V}/\tilde{v}}} \left(\prod_{e \in \mathcal{E}_1} T^e[X_e = x_e] \right) \cdot \left(\prod_{e \in \mathcal{E}_2} T^e[X_e = x_e] \right) \\ &= \sum_{x_{\mathcal{V}_1/\tilde{v}}} \sum_{x_{\mathcal{V}_2/(\tilde{v} \cup \mathcal{V}_1)}} \left(\prod_{e \in \mathcal{E}_1} T^e[X_e = x_e] \right) \cdot \left(\prod_{e \in \mathcal{E}_2} T^e[X_e = x_e] \right) \\ &= \sum_{x_{\mathcal{V}_1/\tilde{v}}} \left(\prod_{e \in \mathcal{E}_1} T^e[X_e = x_e] \right) \cdot \left(\sum_{x_{\mathcal{V}_2/(\tilde{v} \cup \mathcal{V}_1)}} \prod_{e \in \mathcal{E}_2} T^e[X_e = x_e] \right). \end{aligned}$$

When contracting the variables $X_{\mathcal{V}_2/(\tilde{v} \cup \mathcal{V}_1)}$ on $\mathcal{T}^{\mathcal{G}_2}$, the variables $X_{\mathcal{V}_2 \cap (\tilde{v} \cup \mathcal{V}_1)}$ are left open. We therefore have for any $x_{\mathcal{V}_2 \cap (\tilde{v} \cup \mathcal{V}_1)}$

$$\langle \mathcal{T}^{\mathcal{G}_2} \rangle [X_{\mathcal{V}_2 \cap (\tilde{v} \cup \mathcal{V}_1)} = x_{\mathcal{V}_2 \cap (\tilde{v} \cup \mathcal{V}_1)}] = \left(\sum_{x_{\mathcal{V}_2/(\tilde{v} \cup \mathcal{V}_1)}} \prod_{e \in \mathcal{E}_2} T^e[X_e = x_e] \right).$$

It follows with the above, that

$$\begin{aligned} \langle \mathcal{T}^{\mathcal{G}} \rangle [X_{\tilde{v}} = x_{\tilde{v}}] &= \sum_{x_{\mathcal{V}_1/\tilde{v}}} \left(\prod_{e \in \mathcal{E}_1} T^e[X_e = x_e] \right) \cdot \langle \mathcal{T}^{\mathcal{G}_2} \rangle [X_{\mathcal{V}_2 \cap (\tilde{v} \cup \mathcal{V}_1)} = x_{\mathcal{V}_2 \cap (\tilde{v} \cup \mathcal{V}_1)}] \\ &= \left\langle \mathcal{T}^{\mathcal{G}_1} \cup \{ \langle \mathcal{T}^{\mathcal{G}_2} \rangle [X_{\mathcal{V}_2 \cap (\mathcal{V}_1 \cup \tilde{v})}] \} \right\rangle [X_{\tilde{v}} = x_{\tilde{v}}]. \end{aligned}$$

\square

14.5 Support of Contractions

To state the next theorem we introduce the nonzero function $\chi : \mathbb{R} \rightarrow [2]$ by

$$\chi(x) = \begin{cases} 0 & \text{if } x = 0 \\ 1 & \text{else} \end{cases} \quad (69)$$

Applied coordinatewise on tensors it marks the nonzero coordinates by 1.

We show that adding binary tensor cores to an contraction orders the results by the partial ordering introduced in Definition 46

Theorem 92 (Monotonicity of Tensor Contractions). *Let $\mathcal{T}^{\mathcal{G}}, \mathcal{T}^{\tilde{\mathcal{G}}}$ be tensor network of non-negative tensors and $X_{\tilde{\mathcal{V}}}$ an arbitrary set of random variables. Then we have*

$$\chi \left(\left\langle \mathcal{T}^{\mathcal{G}} \cup \mathcal{T}^{\tilde{\mathcal{G}}} \right\rangle [X_{\tilde{\mathcal{V}}}] \right) \prec \chi \left(\left\langle \mathcal{T}^{\mathcal{G}} \right\rangle [X_{\tilde{\mathcal{V}}}] \right).$$

Proof. It suffices to show that for any $x_{\tilde{\mathcal{V}}}$ with

$$\chi \left(\left\langle \mathcal{T}^{\mathcal{G}} \cup \mathcal{T}^{\tilde{\mathcal{G}}} \right\rangle [X_{\tilde{\mathcal{V}}} = x_{\tilde{\mathcal{V}}}] \right) = 1$$

we also have

$$\chi \left(\left\langle \mathcal{T}^{\mathcal{G}} \right\rangle [X_{\tilde{\mathcal{V}}} = x_{\tilde{\mathcal{V}}}] \right) = 1.$$

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

$$\left\langle \mathcal{T}^{\mathcal{G}} \cup \mathcal{T}^{\tilde{\mathcal{G}}} \right\rangle [X_{\mathcal{V}} = x_{\mathcal{V}}] > 0$$

and it follows that

$$\left\langle \mathcal{T}^{\mathcal{G}} \right\rangle [X_{\mathcal{V}} = x_{\mathcal{V}}] > 0 \quad \text{and} \quad \left\langle \mathcal{T}^{\tilde{\mathcal{G}}} \right\rangle [X_{\mathcal{V}} = x_{\mathcal{V}}] > 0.$$

But this already implies, that

$$\chi \left(\left\langle \mathcal{T}^{\mathcal{G}} \right\rangle [X_{\tilde{\mathcal{V}}} = x_{\tilde{\mathcal{V}}}] \right) = 1.$$

□

Let us now state an equivalence of the contraction, when we add the result of the same contraction

Theorem 93 (Invariance under adding subcontractions). *Let $\mathcal{T}^{\mathcal{G}}$ be a tensor network of non-negative tensors with variables $X_{\mathcal{V}}$ and let $\mathcal{T}^{\tilde{\mathcal{G}}}$ be a subset. Then we have for any subset $X_{\tilde{\mathcal{V}}}$ of $X_{\mathcal{V}}$*

$$\left\langle \mathcal{T}^{\mathcal{G}} \cup \left\{ \chi \left(\left\langle \mathcal{T}^{\tilde{\mathcal{G}}} \right\rangle [X_{\tilde{\mathcal{V}}}] \right) \right\} \right\rangle [X_{\mathcal{V}}] = \left\langle \mathcal{T}^{\mathcal{G}} \right\rangle [X_{\mathcal{V}}].$$

Proof. For any $x_{\mathcal{V}}$ with

$$\left\langle \mathcal{T}^{\mathcal{G}} \right\rangle [X_{\mathcal{V}} = x_{\mathcal{V}}] = 0$$

we also have

$$\left\langle \mathcal{T}^{\mathcal{G}} \cup \left\{ \chi \left(\left\langle \mathcal{T}^{\tilde{\mathcal{G}}} \right\rangle [X_{\tilde{\mathcal{V}}}] \right) \right\} \right\rangle [X_{\mathcal{V}} = x_{\mathcal{V}}] = 0.$$

For any $x_{\mathcal{V}}$ with

$$\left\langle \mathcal{T}^{\mathcal{G}} \right\rangle [X_{\mathcal{V}} = x_{\mathcal{V}}] \neq 0$$

we have for the reduction $x_{\tilde{\mathcal{V}}}$ of the index $x_{\mathcal{V}}$ that

$$\left\langle \mathcal{T}^{\tilde{\mathcal{G}}} \right\rangle [X_{\tilde{\mathcal{V}}} = x_{\tilde{\mathcal{V}}}] \neq 0$$

and thus

$$\left\langle \mathcal{T}^{\mathcal{G}} \cup \left\{ \chi \left(\left\langle \mathcal{T}^{\tilde{\mathcal{G}}} \right\rangle [X_{\tilde{\mathcal{V}}}] \right) \right\} \right\rangle [X_{\mathcal{V}} = x_{\mathcal{V}}] = \left\langle \mathcal{T}^{\mathcal{G}} \right\rangle [X_{\mathcal{V}} = x_{\mathcal{V}}] \cdot \chi \left(\left\langle \mathcal{T}^{\tilde{\mathcal{G}}} \right\rangle [X_{\tilde{\mathcal{V}}} = x_{\tilde{\mathcal{V}}}] \right) = \left\langle \mathcal{T}^{\mathcal{G}} \right\rangle [X_{\mathcal{V}} = x_{\mathcal{V}}].$$

□

Remark 32. *Similar statements hold, when dropping the non-negativity assumption on the, but demanding that all variables are left open.*

15 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.

15.1 Encoding of Subsets and Relations

Based on the concept of one-hot encodings of states we in this chapter develop the construction of encodings to sets, relations and functions.

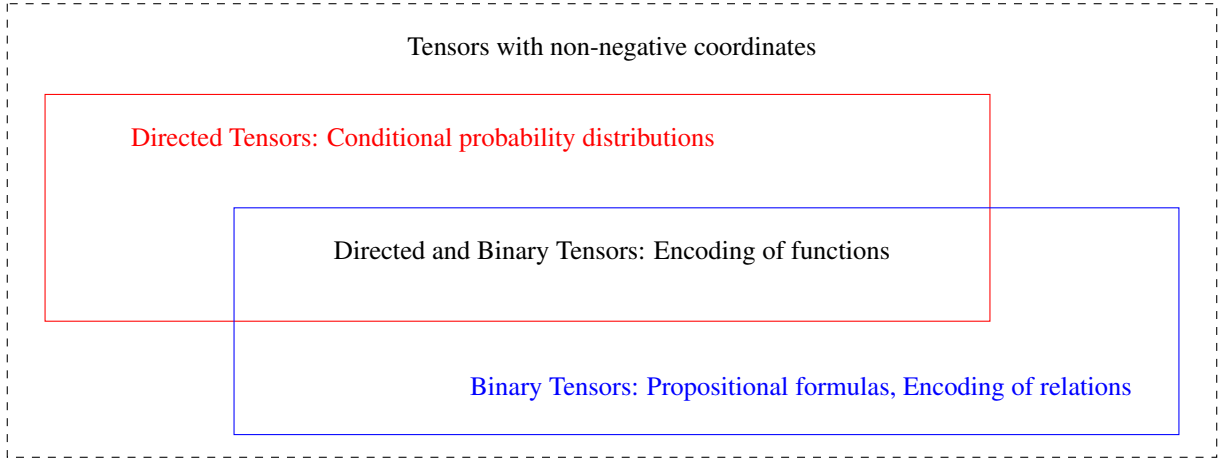


Figure 29: Sketch of the tensors with non-negative coordinates. We investigate in this chapter tensors, which are directed and binary.

Here we show how we can use binary numbers to encode the truth of set memberships.

Definition 73 (Subset Encoding). *We say that an arbitrary set \mathcal{U} is enumerated by a categorical variables $O_{\mathcal{U}}$ taking values in $[m_{\mathcal{U}}]$, when $m_{\mathcal{U}} = |\mathcal{U}|$ and there is a bijective function*

$$I : [m_{\mathcal{U}}] \rightarrow \mathcal{U} .$$

Given an set \mathcal{U} enumerated by the variable $O_{\mathcal{U}}$, any subset $\mathcal{V} \subset \mathcal{U}$ is encoded by the tensor $e_{\mathcal{V}}[O]$ defined for $x \in [|\mathcal{U}|]$ as

$$e_{\mathcal{V}}[O = x] = \begin{cases} 1 & \text{if } I(x) \in \mathcal{V} \\ 0 & \text{else} \end{cases} . \quad (70)$$

In a one-hot basis decomposition we have

$$e_{\mathcal{V}}[O] := \sum_{x \in [|\mathcal{U}|] : I(x) \in \mathcal{V}} e_x [O] .$$

Encoding of subsets as vectors: Each coordinate associated with a possible element, $\{0, 1\}$ encoding whether in subset. The encodings is thus a binary tensor. Any subset encoding is a binary tensor.

Definition 74 (Relation Encoding). *Given two finite sets $\mathcal{U}^{\text{in}}, \mathcal{U}^{\text{out}}$, a relation 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}^{in} and \mathcal{U}^{out} by the categorical variables O_{in} and O_{out} and interpretation maps $I^{\text{in}}, I^{\text{out}}$, we define the encoding of this subset as the tensor $e_{\mathcal{R}}[O_{\text{in}}, O_{\text{out}}]$ with the coordinates

$$e_{\mathcal{R}}[O_{\text{in}} = x_{\text{in}}, O_{\text{out}} = x_{\text{out}}] = \begin{cases} 1 & \text{if } (I^{\text{in}}(x_{\text{in}}), I^{\text{out}}(x_{\text{out}})) \in \mathcal{R} \\ 0 & \text{else} \end{cases} . \quad (71)$$

The relation encoding has a decomposition into one-hot encodings as

$$e_{\mathcal{R}}[O_{\text{in}}, O_{\text{out}}] = \sum_{x_{\text{in}}, x_{\text{out}} : (I^{\text{in}}(x_{\text{in}}), I^{\text{out}}(x_{\text{out}})) \in \mathcal{R}} e_{x_{\text{in}}} [X_{\text{in}}] \otimes e_{x_{\text{out}}} [X_{\text{out}}] .$$

Relations are subsets of cartesian products and encodings of relations are the encodings of subsets by vectors. They have a matrix structure by the cartesian product, which can be further folded to tensors, when the sets itself are cartesian products.

Theorem 94. *The relational encoding is a bijection between the set of relations and the set of binary tensors.*

Proof. By definition, a relational encoding is the encoding of a subset and thus a binary tensor. Any matricification of a binary tensor marks by its 1 coordinates the elements of a relation. \square

We can thus understand any matricification of a binary tensor as the encoding of a relation and vice versa.

15.1.1 Higher order relations

We can extend this contraction to relations of higher order, and arrive at encoding schemes usable for relational databases.

Definition 75. *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 encoding of the relation as the tensor $e_{\mathcal{R}}[O_{[d]}]$ with coordinates

$$e_{\mathcal{R}}[O_0 = x_0, \dots, O_{d-1} = x_{d-1}] = \begin{cases} 1 & \text{if } (I^0(x_0), \dots, I^{d-1}(x_{d-1})) \in \mathcal{R} \\ 0 & \text{else} \end{cases} . \quad (72)$$

Example 19 (Propositional Formulas). *Propositional formulas are equal to the subset encoding of their models. The sets \mathcal{U}^k are all $[2]$ and are interpreted as the possible assignments to the boolean atoms.*

Example 20 (Relational Databases). *Relational Databases can be encoded as tensors using the relation encoding scheme. Each column is thereby understood as a categorical variable, which 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 17) to represent $e_{\mathcal{R}}$.

15.2 Encoding of Functions

Real valued functions are directly tensors by definition.

15.2.1 Relational Encoding of Functions

In Definition 14 we have introduced the relational encoding of functions between the states of factored systems. We now generalize the representation scheme towards maps between arbitrary unstructured sets.

Definition 76 (Relational Encoding of Maps). *Any map*

$$f : \mathcal{U}^{\text{in}} \rightarrow \mathcal{U}^{\text{out}}$$

can be represented by a relation

$$\mathcal{R}^f := \{(x, f(x)) : x \in \mathcal{U}^{\text{in}}\} \subset \mathcal{U}^{\text{in}} \times \mathcal{U}^{\text{out}} .$$

Given an enumeration of the sets by O_{in} and O_{out} we define the relational encoding of f as the tensor

$$\rho^f [O_{\text{in}}, O_{\text{out}}] = e_{\mathcal{R}^f} [O_{\mathcal{U}^{\text{in}}}, O_{\mathcal{U}^{\text{out}}}] .$$

Remark 33 (Reduction to images). When f maps into a set of infinite cardinality, we restrict \mathcal{U}^{out} to the image of f and enumerate the image by a variable X_f . This scheme is applied, when f is itself a tensor, i.e. $\mathcal{U}^{\text{out}} = \mathbb{R}$. While the variable X_f can in general be of the same cardinality as the domain set \mathcal{U}^{in} , it will be in [2] when considering binary tensors.

We notice, that any relational 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 binary tensors is characterized by the relational encoding of functions. This is shown in the next theorem, by the claim that any binary tensor which is directed is the relational representation of a function.

Theorem 95. 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 $f : \mathcal{U}^{\text{in}} \rightarrow \mathcal{U}^{\text{out}}$ such that $\mathcal{R} = \mathcal{R}^f$, the relational encoding ρ^f is a directed tensor with O_{in} incoming and O_{out} outgoing.

Proof. When f is a function, we have for any $o_{\text{in}} \in [r_{\text{in}}]$

$$\sum_{o_{\text{out}} \in [r_{\text{out}}]} \rho^f [O_{\text{in}} = o_{\text{in}}, O_{\text{out}} = o_{\text{out}}] = \rho^f [O_{\text{in}} = o_{\text{in}}, O_{\text{out}} = f(o_{\text{in}})] = 1.$$

Conversely let there be a relation \mathcal{R} , such that $\rho^{\mathcal{R}}$ is directed. To this end, we observe that for any $o_{\text{in}} \in [r_{\text{in}}]$ the tensor

$$e_{\mathcal{R}} [O_{\text{in}} = o_{\text{in}}, O_{\text{out}}]$$

is a binary tensor with coordinate sum one and therefore a basis vector. It follows that the function $f : \mathcal{U}^{\text{in}} \rightarrow \mathcal{U}^{\text{out}}$

$$f(o_{\text{in}}) = e^{-1}(e_{\mathcal{R}} [O_{\text{in}} = o_{\text{in}}, O_{\text{out}}])$$

is well-defined. We then have by construction

$$\rho^f = \sum_{o_{\text{in}} \in [r_{\text{in}}]} e_{o_{\text{in}}} \otimes e_{f(o_{\text{in}})} = \sum_{o_{\text{in}} \in [r_{\text{in}}]} e_{o_{\text{in}}} \otimes e_{\mathcal{R}} [O_{\text{in}} = o_{\text{in}}, O_{\text{out}}] = e_{\mathcal{R}}$$

and therefore $\mathcal{R} = \mathcal{R}^f$. □

We are specially interested in sets of states of a factored system, which amounts to the case in Definition 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_{in} and d_{out} state variables can be represented by $d_{\text{in}} + d_{\text{out}}$ -ary relations and Definition 76 has an obvious generalization to this case with multiple enumeration variables.

Since the relational encoding of any map between factored systems is directed, it can be interpreted by a conditional probability tensor, as we state next.

Corollary 11. The relational encoding ρ^f (see Definition 14) of a function f between factored systems is a conditional probability tensor, where the legs to the image system are the conditions and the legs to the target system the distribution legs.

These are deterministic conditional probability tensors, in the sense that any slice with respect to the input variables is a basis tensor. Through contractions with distribution tensors (e.g. distributions in domain systems) they get stochastic. This is for example the case in the empirical distribution, which can be understood as the forwarding of the uniform distribution on the sample enumeration.

15.3 Calculus of relational encodings

15.3.1 Function Evaluation

Theorem 96 (Basis Calculus). Retrieving the value of the function at a specific state is then the contraction of the tensor representation with the one-hot encoded state. For any state indexed by (x_0, \dots, x_{d-1}) we have

$$e_{f(x_0, \dots, x_{d-1})} [X_f] = \langle \{\rho^f, e_{x_0, \dots, x_{d-1}}\} \rangle [X_f].$$

Thus, we can retrieve the function evaluation by the inverse one-hot mapping as

$$f(x_0, \dots, x_{d-1}) = e^{-1}(\langle \{\rho^f, e_{x_0, \dots, x_{d-1}}\} \rangle [X_f]).$$

This generalizes Example ?? to arbitrary maps between factored systems.

Proof. By doing the summation in tensor product notation. \square

We can thus use tensor contractions to calculate the values of functions. Since basis vectors being the one-hot encoding of the domain system are mapped to basis vectors being the encoding of the image system, we call these contraction basis calculus.

15.3.2 Composition of function

We have already used, that combination of propositional formulas by connectives can be represented by contractions. In a more general perspective, any composition of functions between factored systems is in its relational encoding the contraction of the encoded functions.

Theorem 97 (Composition of Functions). *Let there be two maps between factored systems*

$$f : \bigtimes_{v \in \mathcal{V}_1} [m_v] \rightarrow \bigtimes_{v \in \mathcal{V}_2} [m_v]$$

and

$$g : \bigtimes_{v \in \mathcal{V}_2} [m_v] \rightarrow \bigtimes_{v \in \mathcal{V}_3} [m_v]$$

with the image system of f is the domain system of g . Then the relational encoding of the composition

$$g(f) : \bigtimes_{v \in \mathcal{V}_1} [m_v] \rightarrow \bigtimes_{v \in \mathcal{V}_3} [m_v]$$

satisfies

$$\rho^{g(f)} = \langle \{\rho^f, \rho^g\} \rangle [\mathcal{V}_1 \cup \mathcal{V}_3] .$$

Proof. By definition we have

$$\rho^{g(f)} = \sum_{i_1 \in \bigtimes_{v \in \mathcal{V}_1} [m_v]} e_{i_1} \otimes e_{g(f(i_1))} ,$$

where by i_1 we denote in a slide abuse of notation the tuple indexing the state of the domain factored systems of f . On the other side, we have that

$$\begin{aligned} \langle \{\rho^f, \rho^g\} \rangle [X_{\mathcal{V}_1} X_{\mathcal{V}_3}] &= \sum_{\tilde{i}_2 \in \bigtimes_{v \in \mathcal{V}_2} [m_v]} \left(\sum_{i_1 \in \bigtimes_{v \in \mathcal{V}_1} [m_v]} e_{i_1} \cdot (e_{f(i_1)})_{\tilde{i}_2} \right) \left(\sum_{i_2 \in \bigtimes_{v \in \mathcal{V}_2} [m_v]} (e_{i_2})_{\tilde{i}_2} \cdot e_{g(i_2)} \right) \\ &= \sum_{\tilde{i}_2 \in \bigtimes_{v \in \mathcal{V}_2} [m_v]} \left(\sum_{i_1 \in \bigtimes_{v \in \mathcal{V}_1} [m_v]} e_{i_1} \cdot \delta_{f(i_1), \tilde{i}_2} \right) \left(\sum_{i_2 \in \bigtimes_{v \in \mathcal{V}_2} [m_v]} \delta_{i_2 \tilde{i}_2} \cdot e_{g(i_2)} \right) \\ &= \sum_{i_1 \in \bigtimes_{v \in \mathcal{V}_1} [m_v]} \sum_{\tilde{i}_2 \in \bigtimes_{v \in \mathcal{V}_2} [m_v]} \delta_{f(i_1), \tilde{i}_2} \cdot (e_{i_1} \otimes e_{g(i_2)}) \\ &= \sum_{i_1 \in \bigtimes_{v \in \mathcal{V}_1} [m_v]} e_{i_1} \otimes e_{g(f(i_1))} . \end{aligned}$$

Here we represented the contraction of the variables in \mathcal{V}_2 by the summation over another index i_2 . In the last equation we used that the delta tensor does not vanish only for $\tilde{i}_2 = f(i_1)$. The claim follows, since both expressions are equal. \square

We can use Theorem 97 iteratively to further decompose the function g . In this way, the relational encoding of a function consistent of multiple compositions can be represented as the contractions of all the functions. The relational encoding of propositional formulas for instance can in this way be represented as a contraction of the encodings of its logical connectives applied on the respective formula spaces.

15.3.3 Compositions with real functions

Follows from composition calculus above with the usage of

$$T = \langle \rho^T, \text{Id}_{\text{im}(T)} \rangle [X_{[d]}] .$$

We here investigate how the composition of a tensor

$$T : \bigtimes_{k \in [d]} [m_k] \rightarrow \mathbb{R}$$

with arbitrary functions

$$h : \mathbb{R} \rightarrow \mathbb{R}$$

can be represented. This is for example relevant, when representing the distributions of an exponential family.

Our main strategy is in understanding the tensor T as a map to its finite image, seen as the enumerated states of a categorical variable building a factored system. We then use the relational encoding ρ^T of this map between factored systems.

By $h|_{\mathcal{M}}$ we further denote the restriction of a real function h to an enumerated set $\mathcal{M} = \{x_i : i \in [|\mathcal{M}|]\} \subset \mathbb{R}$, i.e. the vector

$$h|_{\mathcal{M}} : [|\mathcal{M}|] \rightarrow \mathbb{R}$$

defined for $i \in [|\mathcal{M}|]$ as

$$h|_{\mathcal{M}}(i) = h(x_i) .$$

Theorem 98. *We have for any tensor T and real function h (see Figure 30)*

$$h \circ T = \langle \{\rho^T, h|_{\text{im}(T)}\} \rangle [X_{[d]}] .$$

Proof. We enumerate the image $\text{im}(T)$ by $\{x_i : i \in [|\text{im}(T)|]\}$. For arbitrary but fixed $x_0, \dots, x_{d-1} \in \bigtimes_{k \in [d]} [m_k]$ let \tilde{i} be such that $T(x_0, \dots, x_{d-1}) = x_{\tilde{i}}$. Then we have for X_T denoting the image variable of ρ^T that

$$\langle \{\rho^T, e_{x_0, \dots, x_{d-1}}\} \rangle [\{X_T\}] = e_{\tilde{i}}$$

and

$$\langle \{\rho^T, h|_{\text{im}(T)}, e_{x_0, \dots, x_{d-1}}\} \rangle [\emptyset] = \langle \{h|_{\text{im}(T)}, e_{\tilde{i}}\} \rangle [\emptyset] = h(T(x_0, \dots, x_{d-1})) .$$

Since x_0, \dots, x_{d-1} was chosen arbitrarily from $\bigtimes_{k \in [d]} [m_k]$, this shows that $h \circ T$ and $\langle \{\rho^T, h|_{\text{im}(T)}\} \rangle [X_{[d]}]$ coincide on all inputs and are thus equal. \square

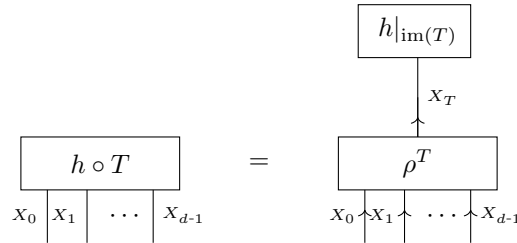


Figure 30: Representation of the composition of a tensor T with a real function h .

Corollary 12. *For any tensor $T [X_{[d]}]$ we have*

$$T [X_{[d]}] = \langle \rho^T, \text{Id}_{\text{im}(T)} \rangle [X_{[d]}] .$$

Proof. Directly by using $h = \text{Id}$. \square

Corollary 13. *For any tensor T , which is directed with $X_{[d]}$ incoming, we have*

$$\mathbb{I} [X_{[d]}] = \langle \rho^T \rangle [X_{[d]}] .$$

Proof. Directly by using $h = \mathbb{I}$. □

Corollary 14. *Let U be a directed and binary tensor with incoming variables being $X_{[d]}$, and T a tensor, which variables are the outgoing variables of U . Let further $h : \mathbb{R} \rightarrow \mathbb{R}$ be any real function. Then*

$$h \circ \langle \{U, T\} \rangle [X_{[d]}] = \langle \{U, h \circ T\} \rangle [X_{[d]}] .$$

Proof. Since U is a directed and binary tensor, we find a map

$$f : \bigtimes_{k \in [d]} [m_k] \rightarrow \bigtimes_{l \in [r]} [m_l]$$

such that $U = \rho^f$ and a map V such that $T = V|_{\text{im}(f)}$. Then Theorem 98 implies that

$$\langle \{U, T\} \rangle [X_{[d]}] = V \circ f .$$

It follows that

$$h \circ \langle \{U, T\} \rangle [X_{[d]}] = h \circ V \circ f$$

and by another application of Theorem 98 that

$$\begin{aligned} h \circ V \circ f &= \langle \rho^f, h \circ V|_{\text{im}(f)} \rangle [X_{[d]}] \\ &= \langle \{U, h \circ T\} \rangle [X_{[d]}] . \end{aligned}$$

The claim follows as a combination of both equations. □

Example 21 (Shannon entropy of empirical distribution). *The Shannon entropy of an empirical distribution can be efficiently computed by contraction of the datatensor with itself along the atom indices, then applying a coordinatewise \ln and averaging.*

This follows from commutations of contraction and coordinatewise contraction (see Corollary 14), using that the datacores are directed. To be more precise, let \tilde{D} be a copy of D with identical image space and copied domain space. Then $\rho^{\tilde{D}_k}$ are tensor cores with identical outgoing legs to ρ^{D_k} , but different incoming legs. We have that

$$\begin{aligned} \mathbb{H} [\mathbb{P}^D] &= \left\langle \left\{ \rho^{D_k} : k \in [d] \right\} \cup \left\{ \frac{1}{m} \mathbb{I} \right\} \cup \left\{ -\ln \left\langle \rho^{\tilde{D}_k} : k \in [d] \right\} \cup \left\{ \frac{1}{m} \mathbb{I} \right\} \right\rangle [X_{[d]}] \right\rangle [\emptyset] \\ &= \left\langle \left\{ \frac{1}{m} \mathbb{I}, -\ln \left\langle \rho^{D_k}, \rho^{\tilde{D}_k} : k \in [d] \right\} \cup \left\{ \frac{1}{m} \mathbb{I} \right\} \right\rangle [X_{[d]}] \right\rangle [\emptyset] \end{aligned}$$

where in the second equation we used Corollary 14.

15.3.4 Decomposition in case of structured images

Here the introduction of multiple head variables, i.e. when

$$\mathcal{U}^{\text{out}} = \bigtimes_{k \in [d]} \mathcal{U}^k$$

This is the case for the empirical distributions!

When the image admits a cartesian representation, the relational encoding can be represented by a contraction of relational encodings to each image coordinate.

Theorem 99. *Let f be a function between factored systems*

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

and denote by

$$f^k : [m] \rightarrow [m_k]$$

the restrictions of f to axes of $\bigtimes_{k \in [d]} [m_k]$. We assign the variable X to the factored system in the domain system of f and the variables X_k for $k \in [d]$ to the image system of f .

We then have

$$\rho^f [X, X_{[d]}] = \left\langle \{ \rho^{f^k} [X, X_k] : k \in [d] \} \right\rangle [X, X_{[d]}] .$$

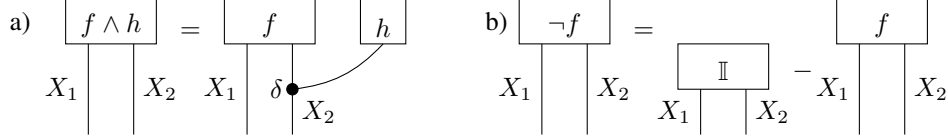


Figure 31: Decomposition schemes by effective calculus. a) Conjunction, b) Negations.

Proof. We have

$$\rho^f [X, X_{[d]}] = \left\langle \{\rho^{f^k} [X, X_k] : k \in [d]\} \right\rangle [X, X_{[d]}]$$

since for any $x \in [m]$

$$\rho^f [X = x, X_{[d]}] = \bigotimes_{k \in [d]} \rho^{f^k} [X = x, X_k] = \left\langle \{\rho^{f^k} [X, X_k] : k \in [d]\} \right\rangle [X = x, X_{[d]}] .$$

To get a representation in the basis CP format, we rename the variable X in the cores $\rho^{f^k} [X, X_k]$ by I and observe that for a trivial scalar core $\mathbb{I} [I]$

$$\left\langle \{\rho^{f^k} [X, X_k] : k \in [d]\} \right\rangle [X, X_{[d]}] = \left\langle \{\mathbb{I} [I]\} \cup \{\rho^{f^k} [I, X_k] : k \in [d]\} \cup \{\delta [X, I]\} \right\rangle [X, X_{[d]}] .$$

□

15.4 Effective Coordinate Calculus

In some situations, we can perform basis calculus by coordinate calculus (see Definition 70). This includes conjunctions, which correspond with coordinatewise multiplication, and negation, which correspond with coordinatewise subtraction from 1.

In an alternative perspective, effective calculus amounts to an contraction against the directionality of the relational encodings. For specific functions, slices of the relational encodings with respect to head variables are basis vectors. In that case, we can perform basis calculus in the inverse direction than suggested by the directions of the tensors. We exemplify this situation in the following theorem for relational encodings of logical conjunctions and negations.

Theorem 100. For any formulas f, h we have

$$\langle \rho^\wedge [X_{f \wedge h}, X_f, X_h], e_1 [X_{f \wedge h}] \rangle [X_f, X_h] = e_1 [X_f] \otimes e_1 [X_h] .$$

In particular, it holds that (see Figure 31a)

$$(f \wedge h)[X_{[d]}] = \langle f, h \rangle [X_{[d]}] .$$

Proof. We decompose

$$\rho^\wedge [X_{f \wedge h}, X_f, X_h] = e_1 [X_{f \wedge h}] \otimes e_1 [X_f] \otimes e_1 [X_h] + e_0 [X_{f \wedge h}] (\mathbb{I} [X_f, \mathbb{I} [X_h]] - e_1 [X_f] \otimes e_1 [X_h])$$

and get the first claim as

$$\begin{aligned} \langle \rho^\wedge [X_{f \wedge h}, X_f, X_h], e_1 [X_{f \wedge h}] \rangle [X_f, X_h] &= \langle e_1 [X_{f \wedge h}] \otimes e_1 [X_f] \otimes e_1 [X_h], e_1 [X_{f \wedge h}] \rangle [X_f, X_h] \\ &= e_1 [X_f] \otimes e_1 [X_h] . \end{aligned}$$

To show the second claim we use

$$\begin{aligned} (f \wedge h)[X_{[d]}] &= \langle \rho^f [X_f, X_{[d]}], \rho^h [X_h, X_{[d]}], \rho^\wedge [X_{f \wedge h}, X_f, X_h], e_1 [X_{f \wedge h}] \rangle [X_{[d]}] \\ &= \langle \rho^f [X_f, X_{[d]}], \rho^h [X_h, X_{[d]}], (e_1 [X_f] \otimes e_1 [X_h]) \rangle [X_{[d]}] \\ &= \langle f, h \rangle [X_{[d]}] . \end{aligned}$$

□

A similar decomposition holds for negations, as we show next.

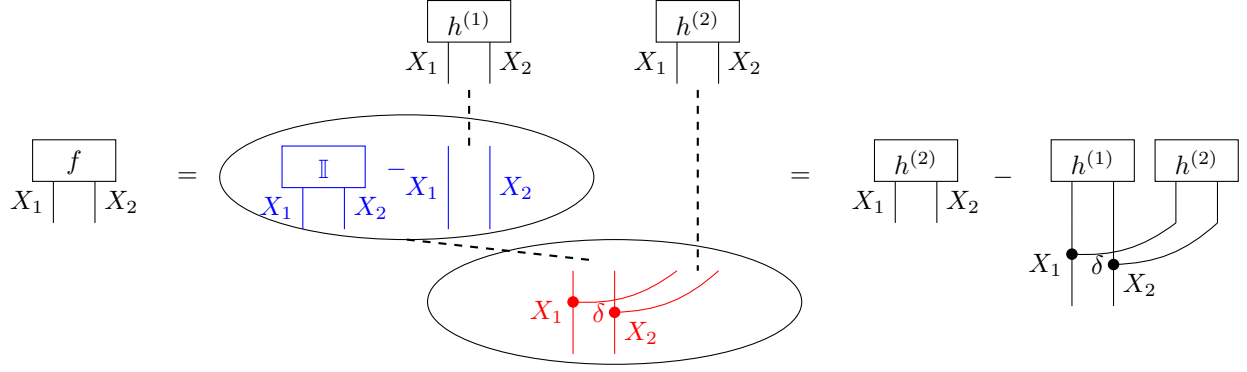


Figure 32: Example of a decomposition by effective calculus of a formula $f(X_1, X_2) = \neg h^{(1)}(X_1, X_2) \wedge h^{(2)}(X_1, X_2)$ into a sum of contractions.

Theorem 101. For any formula f we have

$$\langle \rho^\neg [X_f, X_{\neg f}], e_1 [X_{\neg f}] \rangle [X_f] = e_0 [X_f] = \mathbb{I} [X_f] - e_1 [X_f] .$$

and

$$\langle \rho^\neg [X_f, X_{\neg f}], e_0 [X_{\neg f}] \rangle [X_f] = e_1 [X_f] .$$

In particular, it holds that (see Figure 31b)

$$(\neg f)[X_{[d]}] = \mathbb{I} [X_{[d]}] - f [X_{[d]}] .$$

Proof. Using that for two dimensional variables we have $\mathbb{I} [X] = e_0 [X] + e_1 [X]$. □

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 17. In Figure 32 we provide an example of this decomposition.

15.5 Applications in Machine Learning

The neural paradigm of Machine Learning describes the relevance of sparse function to be effective models in the sense of learning and approximation.

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 binary tensor networks with basis calculus interpretation.

We have already observed in Theorem 96, that the value of discrete maps can be calculated by contractions of the directed binary relation encodings. This has been framed as Basis Calculus. What is more, tensor network decompositions into directed binary tensors correspond with representation of functions as compositions of smaller functions. We can understand each composition as marking a neuron in an architecture and thus have established a neural perspective on binary directed tensor networks.

16 Contraction Message Passing

In this chapter we introduce local contraction passed along tensor clusters to approximatively calculate global contractions. 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.

16.1 Exact Contractions

We apply Theorem 91 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.

16.1.1 Construction of Cluster Graphs

Definition 77 (Cluster Graph). *Given a tensor network $\mathcal{T}^{\mathcal{G}}$ 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 $\mathcal{V}_i = \bigcup_{e \in C_i} e$. The clusters form a graph where edges between C_i and C_j exist, when the node sets \mathcal{V}_i and \mathcal{V}_j are not disjoint. In this case, we define separation sets $S_{i,j} = C_i \cup C_j$

Theorem 102. *Given a tensor network $\mathcal{T}^{\mathcal{G}}$ 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

$$\langle \mathcal{T}^{\mathcal{G}} \rangle [X_{\tilde{\mathcal{V}}}] = \left\langle \{ \langle \mathcal{T}^{C_i} \rangle [\mathcal{V}_i \cap (\tilde{\mathcal{V}}_i \cup \tilde{\mathcal{V}})] : i \in [n] \} \right\rangle [X_{\tilde{\mathcal{V}}}] .$$

Proof. By Theorem 91 applied for each cluster seen as a subgraph. □

16.1.2 Message Passing to calculate contractions

Having a hypergraph \mathcal{G} , we iteratively apply Theorem 91 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

$$\delta_{j \rightarrow i} = \left\langle \{ \delta_{\tilde{j} \rightarrow j} : C_{\tilde{j}} \prec C_j \} \cup \mathcal{T}^{C_j} \right\rangle [X_{\mathcal{V}_i \cap \mathcal{V}_j}] . \quad (73)$$

and

$$\delta_{j \leftarrow i} = \left\langle \{ \delta_{i \leftarrow \tilde{j}} : C_i \prec C_{\tilde{j}} \} \cup \mathcal{T}^{C_i} \right\rangle [X_{\mathcal{V}_i \cap \mathcal{V}_j}] . \quad (74)$$

We note, that the messages are well defined by these recursive equations, exactly when the cluster graph is a polytree.

Lemma 27. *When the cluster graph is a tree, we have for neighbored clusters C_i and C_j with $C_j \prec C_i$*

$$\delta_{j \rightarrow i} = \left\langle \{ \mathcal{T}^{C_{\tilde{j}}} : C_{\tilde{j}} \prec C_j \} \right\rangle [X_{\mathcal{V}_i \cap \mathcal{V}_j}]$$

and

$$\delta_{j \leftarrow i} = \left\langle \{ \mathcal{T}^{C_{\tilde{j}}} : C_i \prec C_{\tilde{j}} \} \right\rangle [X_{\mathcal{V}_i \cap \mathcal{V}_j}] .$$

Proof. By induction over the cardinality of the preceding clusters.

$n = 1$: Only a single cluster before, therefore trivial.

$n + 1 \rightarrow n$: Assuming the statement holds for up to n preceding clusters, let there be $n + 1$ preceding clusters. Then Theorem 91 splits contractions into terms, which are by inductive assumption the messages. □

Theorem 103. *When the cluster graph is a tree, then we have for each cluster C_i with neighbors $N(i)$*

$$\langle \mathcal{T}^{\mathcal{G}} \rangle [\mathcal{V}_i] = \left\langle \{ \delta_{j \rightarrow i} : j \in N(i), C_j \prec C_i \} \cup \{ \delta_{i \leftarrow j} : j \in N(i), C_i \prec C_j \} \cup \mathcal{T}^{C_i} \right\rangle [\mathcal{V}_i] . \quad (75)$$

Proof. By Theorem 91 we split into contractions of the clusters up and down of the respective neighbors and apply the above lemma. □

16.1.3 Variable Elimination Cluster Graphs

Remark 34 (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,*

16.1.4 Bethe Cluster Graphs

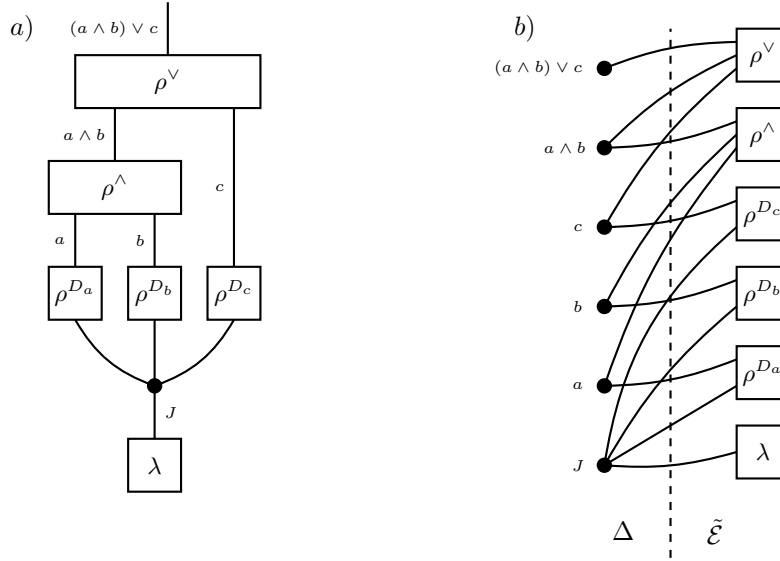


Figure 33: Example of a Bethe Cluster Graph. a) Example of a Tensor Network \mathcal{T}^G , which represents the by λ averaged evaluation of the formula $(a \wedge b) \vee c$ on data D . b) Corresponding Bethe Cluster Hypergraph, which dual is bipartite by the sets Δ and $\tilde{\mathcal{E}}$.

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 Δ^G , as defined in the following (see also Figure 33).

Definition 78. *Given a tensor network \mathcal{T}^G on a decorated hypergraph \mathcal{G} , we define the Bethe Cluster Hypergraph $\tilde{\mathcal{G}}$ as $(\tilde{\mathcal{V}}, \tilde{\mathcal{E}} \cup \Delta^G)$ 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 T^e*
- *Nodes $\tilde{\mathcal{V}} = \bigcup_{e \in \mathcal{E}} \tilde{e}$*
- *Delta Edges $\Delta^G = \{\{v^e : e \ni v\} : v \in \mathcal{V}\}$, each of which decorated by a delta tensor $\delta^{\{v^e : e \ni v\}}$*

By Lemma 24 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 Δ^G . This further makes the dual of the Bethe Cluster Hypergraph a proper graph (i.e. edges consist of node pairs).

16.1.5 Computational Complexity

Tree-width here.

Naive execution of $\langle \mathcal{T}^G \rangle [\tilde{\mathcal{V}}]$: $\prod_{v \in \mathcal{V}} m_v$ 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.

16.2 Approximate Contractions

We ignore that cluster graphs are not trees and perform contraction message passing along neighbored clusters. For the contraction of basis tensor networks, this scheme still provides the exact contraction.

16.2.1 Exact Message Passing for Directed and Binary Contractions

A Tensor Network of directed and binary cores represents the evaluation of composed functions. In a Message Passing Perspective each component (let us call them neurons) can be evaluated, when the evaluation of the ancestor neurons are known.

Lemma 28. *Required? Basis vector factorization suffices? For any collection of categorical variables $X_{[d]}$ with identical dimension and any $x_0 \in [m_0]$ we have*

$$\langle \delta[X_0, \dots, X_{d-1}] \rangle [X_0 = x_0, X_1, \dots, X_{d-1}] = \bigotimes_{k \in \{1, \dots, d-1\}} e_{x_0} [X_k] .$$

Proof. Directly by sum decomposition of delta tensors. □

Theorem 104. *Let \mathcal{T}^G be a tensor network on a directed acyclic hypergraph \mathcal{G} , such that each tensor is Boolean and directed, and such that each variable is appearing only once as an outgoing variable of a hyperedge. We build a cluster graph by storing each edge as a cluster and use the topological order on \mathcal{G} . Then*

$$\delta_{j \rightarrow i} = \langle \mathcal{T}^G \rangle [X_{V_i \cap V_j}] .$$

Proof. **Lemma above needed?** By using that each message is a basis vector (Using Theorem 87 in an induction argument) and can thus be splitted into the product of multiple copies.

Any hypercore, which is not a precessor to T^{e_i} can be omitted from the contraction by a root-stripping argument using its directionality. Therefore we have

$$\langle \mathcal{T}^G \rangle [X_{V_i \cap V_j}] = \left\langle \{T^{e_j} : e_j \prec e_i\} \right\rangle [X_{V_i \cap V_j}] .$$

We then replace each variable which is appearing more than once in outgoing legs by a delta tensor, which does not change the contraction by Lemma 24.

We then follow a leaf stripping argument and apply Theorem 91 iteratively on the remaining leaves. Along that line, the leave and its successors are contracted. The contraction is a basis vector and can therefore be represented as an outer product of basis vectors. □

When replacing e_{x_0} by an arbitrary vector in Lemma 28 we have

$$\langle \delta[X_0, \dots, X_{d-1}], V[X_0] \rangle [X_1, \dots, X_{d-1}] \neq \bigotimes_{k \in \{1, \dots, d-1\}} V[X_k] .$$

Therefore, the messages will in general differ from the exact contractions. To provide intuition of what happens in this case, let us take the following cases into account:

- $\lambda \cdot e_x$ sent multiple times: Result gets a factor of $\lambda^{\# \text{copies}}$ compared with the exact contraction.
- $e_x + e_{\tilde{x}}$: Result is the exact contraction added by the crossterms of sending e_x in one and $e_{\tilde{x}}$ in the other direction.

16.2.2 Case of Matrices

Here a toy example of cycling messages starting with \mathbb{I} . When normating the messages, the maximal singular vectors will be dominant.

We investigate the Bethe message passing for a tensor network consisting of a single matrix.

Theorem 105. *The stable messages are the linear subspace of the maximal singular values of the fixed core T .*

Proof. Having a Singular Value Decomposition of T and decompose the messages in the orthonormal system of the respective singular vectors. \square

For the propagation of a and b on binary $f(a, b)$ starting with trivial messages of a and b the above theorem implies:

- Case of single possible world: Exact $f(a, b) \in \{a \wedge b, a \wedge \neg b, \neg a \wedge b, \neg a \wedge \neg b\}$ Messages are after first iteration exact
- Case of two possible worlds and $f(a, b) \in \{a \Leftrightarrow b, a \Leftrightarrow \neg b, \neg a \Leftrightarrow b, \neg a \Leftrightarrow \neg b\}$. In this situation any start message is stable and determines the other.
- Case of two possible worlds and $f(a, b) \in \{a, b, \neg a, \neg b\}$. In this situation one stable message is determined by the specified atom and the other is always stable.
- Case of three possible worlds: Approximative (exact: $1/3, 2/3$, approximative: $1 - \text{goldenratio}, \text{goldenratio}$)
- Case of four possible worlds: Exact (\mathbb{I})

16.2.3 Case of Tensors

Let there now be a single tensor of arbitrary order. When the tensor is not ODECO, we cannot find a CP-Decomposition with leg vectors building an orthonormal system in the respective leg spaces. This prohibits direct application of the same techniques in the case of a matrix, which is always ODECO.

16.3 Basis Calculus

Message Passing of directed and binary message by relational encoding of functions can be interpreted as function evaluation. This is because any relational encoding of a function, the decomposition

$$\rho^f = \sum_{y \in \text{im}(f)} \left(\sum_{i: f(i)=y} e_i \right) \otimes e_y$$

is a SVD of the matricification of ρ^f with respect to incoming and outgoing legs.

Passing a message e_i in direction thus gives the message $e_{f(i)}$.

Remark 35 (Basis Calculus as Message Passing). *Given a tensor network of directed and binary tensor cores T^e , each representing a function f^e . When there are not directed cycles, we define the compositions of f^e to be the function f 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*

$$\langle \{T^e : e \in \mathcal{E}\} \rangle [\mathcal{V}_2] = e_{f(x_v : v \in \mathcal{V}_1)}.$$

16.4 Applications

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.

17 Sparse Tensor Representations

We in this chapter investigate, which sparsity notations enable tensors to be representable as contractions of tensor networks.

17.1 CP Formats

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 tensor of higher order, each such tensor product is over multiple vectors,

Definition 79. A CP Decomposition of rank n of a tensor $T \in \bigotimes_{k \in [d]} \mathbb{R}^{m_k}$ is a collections of tensors $\sigma[I]$ and $V^k[I, X_k]$ for $k \in [d]$, where I takes values in $[n]$, such that

$$T[X_{[d]}] = \langle \{\sigma[I]\} \cup \{V^k[I, X_k] : k \in [d]\} \rangle [X_{[d]}].$$

We say that the CP Decomposition is

- *directed*, when for each k the core V^k is directed with I incoming and X_k outgoing.
- *binary*, when for each k the core V^k is binary.
- *basis*, where we demand both properties, that is for each $k \in [d]$ and $i \in [n]$

$$V^k[J = i, X_k] \in \{e_{[x_k]}[X_k] x_k \in [m_k]\}.$$

- *basis+*, when for each $k \in [d]$ and $i \in [n]$

$$V^k[J = i, X_k] \in \{e_{[x_k]}[X_k] x_k \in [m_k]\} \cup \{\mathbb{I}[X_k]\}.$$

We denote by $\text{rank}(T)$, respectively $\text{rank}^{\text{bin}}(T)$, $\text{rank}^{\text{bas}}(T)$ and $\text{rank}^{\text{bas}+}(T)$ the minimal cardinality such that T has a CP Decomposition with directed cores, respectively binary cores, basis cores and basis+ cores.

We have by definition

$$T[X_{[d]}] = \sum_{i \in [n]} \sigma[J = i] \left(\bigotimes_{k \in [d]} V^k[J = i, X_k] \right).$$

The right side can be seen as an alternative definition of CP Decompositions by summations of elementary tensors.

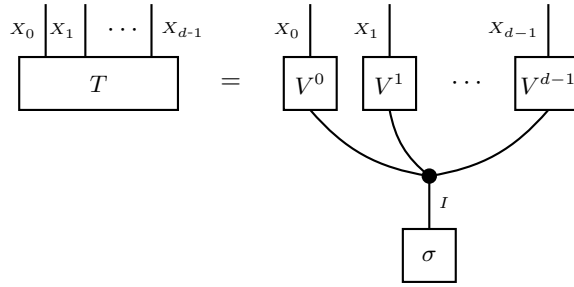


Figure 34: Tensor Network diagram of a generic CP decomposition (see Definition 79)

We introduce different notions of sparsities based on CP Decomposition with different properties of their leg cores.

17.1.1 Directed Leg Cores

This is the canonical CP Decomposition, where the vectors $V^{k,i}$ are interpreted as generalized singular vectors. Any CP decomposition can be transformed into a directed CP decomposition without enlarging the index set J , simply by dividing the vectors by their norms and multiplying it to $\sigma[J = i]$.

We then have a partially directed Tensor Network representing the decomposed tensor. The only undirected core is σ , since we do not demand it to be normed. In many applications applications, however, also the σ is directed with a single outgoing leg (see for example the empirical distributions as discussed in Section 3.8). In that case, also the decomposed tensor is directed with outgoing legs.

17.1.2 Basis Leg Cores

Directed and binary leg cores have incoming slices being basis vectors, we thus call them basis CP Decomposition. This allows the interpretation of the directed and binary CP decomposition in terms of mapping to nonzero coordinates. We start by defining the number of nonzero coordinates of tensors by the ℓ_0 -norm.

Definition 80. The ℓ_0 -norm counts the nonzero coordinates of a tensor by

$$\ell_0(T) = \#\{x_0, \dots, x_{d-1} : T_{x_0, \dots, x_{d-1}} \neq 0\}.$$

The ℓ_0 -norm is not a proper norm itself, but the limit of ℓ_p -norms (where $p \rightarrow 0$) of the flattened tensor (which are norms for $p \geq 1$).

The ℓ_0 norm is the number of nonzero coordinates. We understand the leg cores as the relational encoding of functions mapping to the slices of these coordinates given an enumeration. This is consistent with the previous analysis of Chapter 13, where we characterized binary and directed cores by the encoding of associated functions. Based on this idea, we can proof, that any tensor has a directed and binary CP decomposition with $\text{rand } \ell_0(T)$.

Theorem 106. For any tensor T we have

$$\text{rank}^{\text{bas}}(T) = \ell_0(T) .$$

Proof. We find a map

$$f : [\ell_0(T)] \rightarrow \bigtimes_{k \in [d]} [m_k] ,$$

which image is the set of nonzero coordinates of T . Denoting its image coordinate maps by f^k we have

$$T = \sum_{j \in [m]} \sigma[f(j)] \left(\bigotimes_{k \in [d]} e_{f^k(j)} \right) .$$

This is a basis CP Decomposition with rank $\ell_0(T)$. Conversely, any basis CP Decomposition of T with dimension r would have at most r coordinates different from zero and thus $\ell_0(T) \leq r$. Thus, there cannot be a CP Decomposition with a dimension $r \leq \ell_0(T)$. \square

The next theorem relates the basis CP decomposition with encodings of d -ary relations (see Definition 75).

Theorem 107. If and only if $T \in \bigotimes_{k \in [d]} \mathbb{R}^{m_k}$ has a basis decomposition with slices $\{x_{[d]}^i : i \in [n]\}$ and trivial cores, it coincides with the encoding of the d -ary relation

$$\mathcal{R} = \{x_{[d]}^i : i \in [n]\} .$$

If in addition $m_k = 2$, we can interpret basis CP decompositions as propositional formulas.

Theorem 108. If $T \in \bigotimes_{k \in [d]} \mathbb{R}^2$ has a basis decomposition with slices $\{x_{[d]}^i : i \in [n]\}$ and trivial cores, it coincides with the propositional formula

$$f[X_{[d]}] = \bigvee_{i \in [n]} Z_{x_{[d]}^i}^{\wedge} .$$

Proof. This is a generalization of Theorem 48, which follows from Theorem 33. \square

The storage demand of any CP decomposition is at most linear in the dimension 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 e , 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 store the CP Decomposition by the matrix

$$M[I, L] \in \mathbb{R}^{m \times (d+1)}$$

defined for $k \in [d]$

$$M[J = i, L = k] = e^{-1}(V^k[J = i, X_k])$$

and

$$M[J = i, L = d] = \sigma[J = i] .$$

This is a typical tabular format to store relational databases.

17.1.3 Basis+ Leg Cores

The minimal rank of CP Decomposition is closely related to polynomial sparsity of the map T , which we will define next.

Definition 81. A monomial decomposition of a tensor $T \in \bigotimes_{k \in [d]} \mathbb{R}^{m_k}$ is a set \mathcal{M} of tuples (λ, A, x_A) where $\lambda \in \mathbb{R}$, $A \subset [d]$ and $x_A \in \bigtimes_{k \in A} [m_k]$ such that

$$T[X_{[d]}] = \sum_{(\lambda, A, x_A) \in \mathcal{M}} \lambda \cdot \langle e_{x_A} \rangle [X_{[d]}] . \quad (76)$$

For any tensor $T \in \bigotimes_{k \in [d]} \mathbb{R}^{m_k}$ we define its polynomial sparsity of order r as

$$\text{rank}^r(T) = \min \left\{ |\mathcal{M}| : T[X_{[d]}] = \sum_{(\lambda, A, x_A) \in \mathcal{M}} \lambda \cdot \langle e_{x_A} \rangle [X_{[d]}] , \forall_{(\lambda, A, x_A) \in \mathcal{M}} |A| \leq r . \right\}$$

We refer to the terms in a decomposition (76) in Definition 81 as monomials of binary variables, which are enumerated by pairs (k, x_k) and indicate whether the variable X_k is in state $x_k \in [m_k]$. Such indicators are represented by the one-hot encodings

$$e_{x_k} [X_k] .$$

The monomial of multiple such binary variables indicated, whether all variables labelled by a set A are in the state X_A , which is represented by

$$e_{x_A} [X_A] = \bigotimes_{k \in A} e_{x_k} [X_k] .$$

The states of the variables labeled by $k \in [d]/A$ are not specified in the monomial and the monomial is trivially extended to

$$\langle e_{x_A} \rangle [X_{[d]}] = e_{x_A} [X_A] \otimes \mathbb{I} [X_{[d]/A}] .$$

For some $r < d$ there are tensors $T [X_{[d]}]$, which do not have a monomial decomposition of order r . In that case we the minimum is over an empty set and we define $\text{rank}^r (T) = \infty$. We characterize in the next theorem the set of tensors with monomial decompositions.

Theorem 109. *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. Any sum of tensors with a monomial decomposition of order r admits again a monomial decomposition, which is the concatenation of both. The same holds for a scalar multiplication, and thus, the sets of such tensors form a linear subspace.

The number of different tensors $\langle e_{x_A} \rangle [X_{[d]}]$ is

$$\sum_{s \in [r]} m^s \binom{d}{s}$$

Since any tensor with a monomial decomposition is a weighted sum of those,

We notice, that the set of slices is in general not linear independent, and therefore forms a frame and not a linear basis Casazza et al. (2013). The number of elements in the frame is therefore an upper bound on the dimension. \square

Theorem 109 implies, 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 $\text{rank}^r (T)$ furthermore coincides with the basis+ CP rank $\text{rank}^{\text{bas}+} (T)$.

Theorem 110. *For any tensor $T \in \bigotimes_{k \in [d]} \mathbb{R}^{m_k}$ we have*

$$\text{rank}^d (T) = \text{rank}^{\text{bas}+} (T) .$$

When $x_k = 2$ for all $k \in [d]$, we also have

$$\text{rank}^{\text{bin}} (T) = \text{rank}^d (T) .$$

Proof. To proof the first claim, we construct a basis+ CP decomposition given a monomial decomposition and vice versa. Let there be a tensor $T [X_{[d]}]$ 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

$$V^k [I, X_k] = \left(\sum_{i \in [n] : k \in A^i} e_i [I] \otimes e_{x_k^i} [X_k] \right) + \left(\sum_{i \in [n] : k \notin A^i} e_i [I] \otimes \mathbb{I} [X_k] \right)$$

and

$$\sigma[I] = \sum_{i \in [n]} \lambda^i \cdot e_i[I]$$

and notice that

$$\begin{aligned} T[X_{[d]}] &= \sum_{i \in [n]} \lambda^i \cdot \langle e_{x_A^i} \rangle [X_{[d]}] \\ &= \sum_{i \in [n]} \left(\sigma[J = i] \cdot \bigotimes_{k \in [d]} V^k[J = i, X_k] \right) \\ &= \langle \{\sigma[I]\} \cup \{V^k[I, X_k] : k \in [d]\} \rangle [X_{[d]}] . \end{aligned}$$

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

$$\text{rank}^d(T) \geq \text{rank}^{\text{bas}+}(T) .$$

Let there now be a basis+ CP decomposition we define for each $i \in [n]$

$$A^i = \{k \in [d] : V^k[J = i, X_k] \neq \mathbb{I}[X_k]\} \quad \text{and} \quad x_A^i = \{e^{-1}(V^k[J = i, X_k]) : k \in A\}$$

where by $e^{-1}(\cdot)$ we denote the inverse of the one-hot encoding.

We notice that this is a monomial decomposition of $T[X_{[d]}]$ to the tuple set

$$\mathcal{M} = \{(\sigma[J = i], A^i, x_A^i) : i \in [n]\} .$$

It follows from this that

$$\text{rank}^d(T) \leq \text{rank}^{\text{bas}+}(T)$$

and the first claim is shown.

The second claim follows from the observation, that whenever $x_k = 2$ for all $k \in [d]$ the binary CP decompositions with non-vanishing slices $V^k[J = i, X_k]$ for $k \in [d]$ and $i \in [n]$ are also basis+ CP decompositions and vice versa. \square

Example 22 (Propositional Formulas). *When all leg dimensions of a binary tensor T are 2, we can interpret T as a logical formula. We can use the binary CP decomposition of any tensor \tilde{T} with $\chi(\tilde{T}) = T$ as a CNF of T . Finding the sparsest CNF thus amounts to finding the \tilde{T} with minimal $\text{rank}^d(\tilde{T})$ such that $\chi(\tilde{T}) = T$.*

17.2 Constructive Bounds on CP Ranks

After having defined three CP Decompositions, let us investigate bounds on their ranks which proofs come with constructions of the cores.

17.2.1 Format Transformations

Especially useful, when the leg dimensions are two, where the slice decomposition shows decomposition of the tensor into monomials.

Theorem 111. *For any tensor $T[X_{[d]}] \in \bigotimes_{k \in [d]} \mathbb{R}^{m_k}$ we have*

$$\text{rank}(T) \leq \text{rank}^{\text{bin}}(T) \leq \text{rank}^{\text{bas}+}(T) \leq \text{rank}^{\text{bas}}(T) .$$

Proof. Since any CP decomposition into binary leg cores can be normed to a CP decomposition with directed leg cores, the first bound holds. The second bound holds analogously, since any CP decomposition with basis leg cores is also a CP decomposition with binary leg cores. \square

Consider for example the tensor \mathbb{I} having maximal ℓ_0 -norm being the dimension of the tensor space, but, since it is elementary, a CP decomposition with rank 1.

17.2.2 Summation of CP Decompositions

Theorem 112. For any collections of tensors $\{T^l[X_{\mathcal{V}}] : l \in [n]\}$ with identical variables and scalars $\lambda^l \in \mathbb{R}$ for $l \in [n]$ we have

$$\text{rank} \left(\sum_{l \in [n]} \lambda^l \cdot T^l \right) \leq \sum_{l \in [n]} \text{rank} (T^l) .$$

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 σ . The sum of CP Decomposition is just the combination of all slices, thus the rank is at most additive. \square

17.2.3 Contractions of CP Decompositions

More general, we can bound the sparsity of any contraction by the product of sparsities of affected tensors.

Theorem 113. For any tensor network with variables \mathcal{V} and edges \mathcal{E} we have for any subset $\tilde{\mathcal{V}} \subset \mathcal{V}$

$$\text{rank} \left(\langle \{T^e : e \in \mathcal{E}\} \rangle [\tilde{\mathcal{V}}] \right) \leq \prod_{e \in \mathcal{E} : \tilde{\mathcal{V}} \cap e \neq \emptyset} \text{rank} (T^e) .$$

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 Theorem 113 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 29. For any tensors $T^1[X_{\mathcal{V}_1}]$ and $T^2[X_{\mathcal{V}_2}]$ and any set of variables $\tilde{\mathcal{V}} \subset \mathcal{V}_1 \cup \mathcal{V}_2$ we have

$$\text{rank} \left(\langle \{T^1, T^2\} \rangle [\tilde{\mathcal{V}}] \right) \leq \text{rank} (T^1) \cdot \text{rank} (T^2) .$$

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. By connecting the cores and restoring the binary or basis properties. \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 Lemma 29 as we show next.

Lemma 30. For any tensor network consistent of two tensors $T^1[X_{\mathcal{V}_1}]$ and $T^2[X_{\mathcal{V}_2}]$ and any set $\tilde{\mathcal{V}}$ with $\tilde{\mathcal{V}} \cap \mathcal{V}_2 = \emptyset$ we have

$$\text{rank} \left(\langle \{T^1, T^2\} \rangle [\tilde{\mathcal{V}}] \right) \leq \text{rank} (T^1) .$$

The bound still holds, when we replace on both sides $\text{rank}(\cdot)$ by $\text{rank}^{\text{bin}}(\cdot)$ or by $\text{rank}^{\text{bas}}(\cdot)$.

Proof. We show the lemma by constructing a CP decomposition of $\langle \{T^1, T^2\} \rangle [\tilde{\mathcal{V}}]$ for any CP decomposition of T^1 . Let therefore take any CP decomposition of T^1 consistent of the leg cores $\{V^v : v \in \mathcal{V}_1\}$ and a scalar core σ . Then we define a new σ by

$$\tilde{\sigma} = \langle \{\sigma\} \cup \{V^v : v \in \mathcal{V}_1, v \notin \tilde{\mathcal{V}}\} \cup \{T^2\} \rangle [I] .$$

Then, the leg cores $\{V^v : v \in \tilde{\mathcal{V}}\}$ build with the scalar core $\tilde{\sigma}$ a CP decomposition of $\langle \{T^1, T^2\} \rangle [\tilde{\mathcal{V}}]$. When the CP decomposition of T^1 was binary, basis or basis+, this property is also satisfied by the constructed CP decomposition. Thus the bound also holds for the ranks $\text{rank}^{\text{bin}}(\cdot)$ or $\text{rank}^{\text{bas}}(\cdot)$. \square

Proof of Theorem 113. Use delta tensor representation to represent contractions by graphs. We then iterate through the cores and contract them to the previously contracted tensor, where we apply Lemma 29 when the tensor core has variables left open and Lemma 30 if not. \square

Example 23 (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 ρ^\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 24 (Distributions of independent variables). *Independence means factorization, conditional independence means sum over factorizations. Again, the ℓ_0 norm is bounded by the product of the ℓ_0 norm of the factors.*

17.2.4 Normations of CP Decompositions

As a theorem: If any of the above CP Decomposition is normable, the normation has the same CP ranks. Especially interesting when learning Bayesian Networks, where each core has a CP bound by the number of datapoints.

17.2.5 Sparse Encoding of Functions

We now state that the basis CP rank of relational 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 114. *For any function*

$$f : \prod_{k \in [d]} [m_k] \rightarrow \prod_{l \in [r]} [m_l]$$

between factored systems we have

$$\text{rank}^{\text{bas}}(\rho^f) = \prod_{k \in [d]} m_k.$$

Proof. With Theorem 106, the basis CP rank coincides with the number of not vanishing coordinates, which is the cardinality of the domain of f . \square

Allowing for trivial leg vectors can decrease the CP rank, as we show next.

Theorem 115. *We have*

$$\text{rank}^{\text{bas}+}(\rho^f) \leq \sum_{y \in \text{im}(f)} \text{rank}^{\text{bas}+}(\mathbb{I}_{f==y}),$$

where by $\mathbb{I}_{f==y}$ we denote the indicator, whether the function f evaluates to y .

Proof. We have

$$\rho^f = \sum_{y \in \text{im}(f)} \mathbb{I}_{f==y}[X] \otimes e_y[X_f].$$

For each $y \in \text{im}(f)$ we represent $e_y[X_f]$ in an basis+ CP format with $\text{rank}^{\text{bas}+}(\mathbb{I}_{f==y})$ summands and arrive at a basis+ CP decomposition of ρ^f with $\sum_{y \in \text{im}(f)} \text{rank}^{\text{bas}+}(\mathbb{I}_{f==y})$ summands. \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 Theorem 114, since summing the number of non zero coordinators of the indicators is the cardinality of the domain.

Example 25. *Conjunction of variables For the propositional formula $f = X_0 \wedge X_1$ we have*

$$\rho^f[X_0, X_1] = e_{1,1}[X_0, X_1] \otimes e_1[X_f] + (\mathbb{I}[X_0, X_1] - e_{1,1}[X_0, X_1]) \otimes e_0[X_f]$$

and thus

$$\text{rank}^{\text{bas}+}(\rho^f) \leq 3$$

while $\text{rank}^{\text{bas}}(\rho^f) = 4$. Especially useful for d -ary conjunctions, see Remark 8!

17.2.6 Construction by averaging the incoming legs

Basis CP Decompositions can be constructed by understanding the variable O_{in} of the relational encoding of a function $f : \mathcal{U}^{\text{in}} \rightarrow \mathcal{U}^{\text{out}}$ as the slice selection variable.

Example 26. *Empirical distributions, see Theorem 17* Let there be a data map

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

We can use Theorem 99 to find a tensor network representation for ρ^D as

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

This representation is in the CP format, when adding trivial scalar core and and delta tensor to the data index. It is furthermore in a basis CP format, since all ρ^{D^k} are directed and binary tensors. Normation to get the empirical distribution amounts to setting a slice core with coordinates $\frac{1}{m}$.

17.3 Representation by slice selection architectures

The set of slice-sparse tensors coincides with the expressivity of specific selection architecture. We first define a slice selecting tensor and then show its decomposition into a formula selecting neural network.

Definition 82. *Given a set of atomic variables $X_{[d]}$, a slice 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

$$p_{s,0} = 2, p_{s,1} = d$$

and coordinates

$$\begin{aligned} \mathcal{H}_{\wedge, d, r} [X_{[d]} = x_{[d]}, L_{0,0} = l_{0,0}, \dots, L_{r-1,0} = l_{r-1,0}, L_{0,1} = l_{0,1}, \dots, L_{r-1,1} = l_{r-1,1}] \\ = \begin{cases} 1 & \text{if } \forall_{k,s} : (l_{s,1} = k \wedge l_{s,0} \neq 2) \Rightarrow (l_{s,0} = x_k) \\ 0 & \text{otherwise} \end{cases} . \end{aligned}$$

In the next two Lemmata we first show that the defined slice selecting tensors indeed selects slices and then provide a representation as a formula selecting network.

Lemma 31. *If all input neurons with same selection index are agreeing on the connective index, the selected formula does not vanish and coincides with a slice to the set*

$$A = \{k : \exists_{s \in [n]} : l_{s,1} = k \wedge l_{s,0} \neq 2\}$$

and for $k \in A$

$$x_k = l_{s,0} \quad \text{if } l_{s,1} = k .$$

Proof. We need to show that

$$\mathcal{H}_{\wedge, d, r} [X_{[d]}, L_{0,0} = l_{0,0}, \dots, L_{r-1,0} = l_{r-1,0}, L_{0,1} = l_{0,1}, \dots, L_{r-1,1} = l_{r-1,1}] = e_{x_A} [X_{[d]}] . \quad (77)$$

If and only if an index $\tilde{x}_{[d]}$ reduced on A does not coincide with x_A , we have $e_{x_A} [X_{[d]} = \tilde{x}_{[d]}] = 0$ and otherwise $e_{x_A} [X_{[d]} = \tilde{x}_{[d]}] = 1$. Let us notice, that this condition is equivalent to

$$\forall_{k,s} : (l_{s,1} = k \wedge l_{s,0} \neq 2) \Rightarrow (l_{s,0} = x_k)$$

and thus (77) holds. □

Lemma 32. *The slice selection tensor coincides with a formula selecting neural network with neurons (see Figure 17.3):*

- unary input neuron enumerated by s , selecting one of the $X_{[d]}$ with the variable $L_{s,1}$ and selecting a connective in $\{\neg, \text{Id}, \text{True}\}$ by $L_{s,0}$
- r -ary output neuron fixed to the \wedge connective.

Proof. This can be easily checked on each coordinate. \square

It follows, that the expressivity of the slice selecting neural network coincides with the set of tensors with a bound on their slice sparsity, when $r \geq d$. For arbitrary r , the following theorem holds.

Theorem 116. Let $\mathcal{H}_{\wedge,d,r} [X_{[d]}, L]$ be a slice selecting tensor. For any parameter tensor $\theta [L]$ we have

$$\text{rank}^{\text{bas}+} (\langle \mathcal{H}_{\wedge,d,r} [X_{[d]}, L], \theta [L] \rangle [X_{[d]}]) \leq \text{rank}^{\text{bas}} (\theta [L]).$$

Proof. Each non-vanishing coordinate of θ represents by Lemma 32 a slice and their weighted sum is thus a monomial decomposition. \square

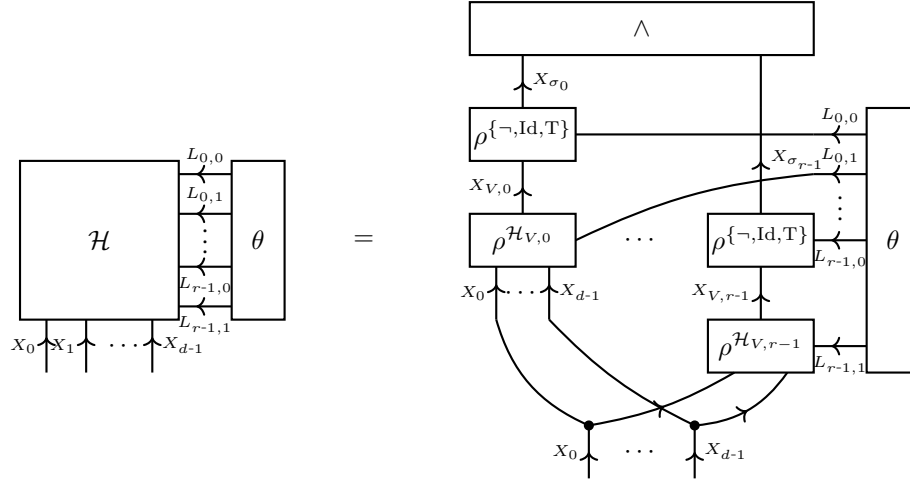


Figure 35: Representation of a basis+ Tensor by the contraction of a parameter tensor θ with a slice selecting architecture \mathcal{H} , which has a decomposition as a formula selecting neural network (see Lemma 32). The nonzero coordinates of θ represent the (see Lemma 31).

17.3.1 Applications

One application is as a parametrization scheme in the approximation of a tensor by a slice-sparse tensor, see Chapter 18.

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 36 (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 $e_0, e_1, \mathbb{I} [1]$. 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.*

17.4 Optimization of sparse tensors

Let us now study the problem of searching for the maximal coordinate in a tensor represented by a monomial decomposition.

17.4.1 Mode search in exponential families

Mode search

$$\max_{x_{[d]} \in \times_{k \in [d]} [2]} \langle \gamma^\phi [X_{[d]} = x_{[d]}, L], \theta \rangle [\emptyset] = \max_{\mu \in \mathcal{M}} \langle \mu [L], \theta [L] \rangle [\emptyset]$$

The search for maximal coordinates appears in various reasoning tasks:

- MAP query as mode search of MLN: T is the contraction of evidence with the distribution, leaving the query variables open.
- Grafting as mode search of proposal distribution: T is the contraction of the gradient of the likelihood with the relational encoding of the hypothesis.

Both tasks have been formulated as mode search problems in exponential families.

17.4.2 Higher-Order Unconstrained Binary Optimization (HUBO)

Here binary refers to the leg dimensions m_k being 2, not to binary coordinates as often referred to in this work.

Definition 83. The binary optimization of a tensor $T[X_{[d]}] \in \times_{k \in [d]}[2]$ is the problem

$$\operatorname{argmax}_{x_{[d]} \in \times_{k \in [d]}[2]} T[X_{[d]} = x_{[d]}] \quad (P_T)$$

We call Problem P_T a Higher Order Unconstrained Binary Optimization (HUBO) problem of order r and sparsity $\operatorname{rank}^r(T)$, when T has a monomial decomposition (see Definition 81) with $|A^i| \leq r$ for all $i \in [n]$, that is when $\operatorname{rank}^r(T) < \infty$.

Remark 37 (Leg dimensions larger than 2). We demanded leg dimensions $m_k = 2$ to have binary 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 Section 8.2.4).

The sparsity $\operatorname{rank}^r(T)$ is the minimal number of monomials, for which a weighted sum is equal to T . Thus we interpret Problem P_T as searching for the maximum in a polynomial consistent of $\operatorname{rank}^r(T)$ monomial terms. **Each monomial is also referred to as potential.**

17.4.3 Quadratic Unconstrained Binary Optimization (QUBO)

Quadratic Unconstrained Binary Optimization problems are HUBOs of order $r = 2$.

We refine the monomial decomposition of tensors (see Definition 81) by demanding that monomials consist of at most two variables.

Definition 84. We call a monomial decomposition \mathcal{M} of a tensor $T \in \otimes_{k \in [d]} \mathbb{R}^2$ a quadratic decomposition, if $|A| \leq 2$ for all $(\lambda, A, x_A) \in \mathcal{M}$. We denote the smallest cardinality $|\mathcal{M}|$ among quadratic decompositions of T by $\operatorname{rank}^{\text{qua}}(T)$.

If a tensor $T \in \otimes_{k \in [d]} \mathbb{R}^2$ has a quadratic decomposition, we call Problem P_T a Quadratic Unconstrained Binary Optimization (QUBO) problem of sparsity $\operatorname{rank}^{\text{qua}}(T)$.

We can transform certain HUBO problems in QUBO problems with the usage of auxiliary variables, as we show in the next lemma.

Lemma 33. 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

17.4.4 Integer Linear Programming

Let us now show how optimization problems can be represented as linear programming problems.

We define for each index tuple $x_{[d]} \in \times_{k \in [d]} [m_k]$ a vector $v_{x_{[d]}} [L] \in \mathbb{R}^d$ with coordinates

$$v_{x_{[d]}} [L = k] = x_k .$$

Definition 85. 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

$$\operatorname{argmax}_{x_{[d]} \in \times_{k \in [d]} [m_k]} \langle c[L], v_{x_{[d]}} [L] \rangle [\emptyset] \quad \text{subject to} \quad \langle M[J, L], v_{x_{[d]}} [L] \rangle [\emptyset] \prec b[J] ,$$

where by \prec we denote partial ordering of tensors (see Definition 46).

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 $T [X_{[d]}]$ 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 $c[L]$, where L takes values in $[d+n]$, as

$$c[L = l] = \begin{cases} \lambda^i & \text{if for a } i \in [n] \text{ we have } l = d + i \\ 0 & \text{else} \end{cases} . \quad (78)$$

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

$$m_i = \sum_{\tilde{i}=0}^i (|A^{\tilde{i}}| + 1)$$

enumerate by \tilde{k} the set A^i and set for $l \in [d+n]$ the coordinates

$$M^{\mathcal{M}} [J = m_i + j, L = l] = \begin{cases} 1 - 2 \cdot x_{\tilde{k}}^i & \text{if } j = l \text{ and } l = \tilde{k} \\ 1 & \text{if } j < d \text{ and } l = d + i \\ -x_{\tilde{k}}^i & \text{if } j = |A^i| + 1 \text{ and } l = \tilde{k} \\ -1 & \text{if } j = |A^i| + 1 \text{ and } l = d + i \end{cases} . \quad (79)$$

All further coordinates of $M^{\mathcal{M}} [J, L]$ not reached by this construction are set to 0. Similarly, we define $b^{\mathcal{M}} [\mathcal{M}]$ as the vector which nonvanishing coordinates are for $i \in [n]$ at

$$b^{\mathcal{M}} [\mathcal{M}] = \begin{cases} 1 & \text{if } j = \tilde{k} \text{ and } x_{\tilde{k}}^i = 0 \\ -1 + |\{\tilde{k} \in A^i : x_{\tilde{k}}^i = 1\}| & \text{if } j = |A^i| + 1 \text{ and } l = d + i \end{cases} . \quad (80)$$

Informally, we pose for each slice $|A| + 1$ 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 117. Given a monomial decomposition $\mathcal{M} = \{(\lambda^i, A^i, x_{A^i}^i) : i \in [n]\}$ of a tensor T , let $x^{ILP, \mathcal{M}}$ be a solution of the integer linear program defined by the matrix and vectors in equations (78), (79) and (80). Then we have

$$x^{ILP, \mathcal{M}}|_{[d]} \in \operatorname{argmax}_{x_{[d]} \in \times_{k \in [d]} [2]} T [X_{[d]} = x_{[d]}] ,$$

where by $x^{ILP, \mathcal{M}}|_{[d]}$ we denote the vector of the first d indices of $x^{ILP, \mathcal{M}}$.

Proof. We show that the constraints by

$$\langle M^{\mathcal{M}}[J, L], v_{x_{[d]}}[L] \rangle [\emptyset] \prec b^{\mathcal{M}}[M]$$

are satisfied, if and only if for all $i \in [n]$

$$z_i = \left(\prod_{\tilde{k} \in A^i, x_{\tilde{k}}^i = 0} (1 - x_{\tilde{k}}^i) \right) \cdot \left(\prod_{\tilde{k} \in A^i, x_{\tilde{k}}^i = 1} x_{\tilde{k}}^i \right).$$

For any $i \in [n]$ the constraints $m_i + \tilde{k}$ for $\tilde{k} \in |A^i|$ are

$$z_i \leq \begin{cases} x_{\tilde{k}}^i & \text{if } x_{\tilde{k}}^i = 1 \\ (1 - x_{\tilde{k}}^i) & \text{if } x_{\tilde{k}}^i = 0 \end{cases}$$

and thus $z_i = 0$ if we have a mismatch on one leg, which establishes

$$z_i \leq \left(\prod_{\tilde{k} \in A^i, x_{\tilde{k}}^i = 0} (1 - x_{\tilde{k}}^i) \right) \cdot \left(\prod_{\tilde{k} \in A^i, x_{\tilde{k}}^i = 1} x_{\tilde{k}}^i \right).$$

The constraint $m_i + |A^i|$ is

$$z_i \geq 1 - \left(\sum_{\tilde{k} \in A^i, x_{\tilde{k}}^i = 0} x_{\tilde{k}}^i \right) + \left(\sum_{\tilde{k} \in A^i, x_{\tilde{k}}^i = 1} (x_{\tilde{k}}^i - 1) \right)$$

which enforces

$$z_i \geq \left(\prod_{\tilde{k} \in A^i, x_{\tilde{k}}^i = 0} (1 - x_{\tilde{k}}^i) \right) \cdot \left(\prod_{\tilde{k} \in A^i, x_{\tilde{k}}^i = 1} x_{\tilde{k}}^i \right).$$

In summary we arrive at

$$z_i = \left(\prod_{\tilde{k} \in A^i, x_{\tilde{k}}^i = 0} (1 - x_{\tilde{k}}^i) \right) \cdot \left(\prod_{\tilde{k} \in A^i, x_{\tilde{k}}^i = 1} x_{\tilde{k}}^i \right)$$

if and only if the indices $y_{[d+n]}$ are feasible.

Now, for any $x_{[d]}$, there is exactly one feasible index $y_{[d+n]}$ extending $x_{[d]}$, and the objective takes for this index the value

$$\langle c[L], v_{x_{[d]}}[L] \rangle [\emptyset] = \sum_{i \in [n]} \lambda^i \cdot z_i \tag{81}$$

$$= \sum_{i \in [n]} \lambda^i \cdot \left(\prod_{\tilde{k} \in A^i, x_{\tilde{k}}^i = 0} (1 - x_{\tilde{k}}^i) \right) \cdot \left(\prod_{\tilde{k} \in A^i, x_{\tilde{k}}^i = 1} x_{\tilde{k}}^i \right) \tag{82}$$

$$= T[X_{[d]} = x_{[d]}] . \tag{83}$$

Therefore, any solution of the ILP reduced to the first d indices corresponding with the axis of T , is a solution of the binary optimization problem to T . \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 ℓ_0 -sparse, and has thus an efficient representation in a basis CP format. More precisely we have by the above construction

$$\ell_0(M^{\mathcal{M}}[J, L]) \leq \sum_{i \in [n]} 3 \cdot |A^i| + 1 .$$

17.5 Value subspaces of functions

To Basis Calculus?

We here provide a subspace perspective for the sparse representation of decomposable functions as tensor networks in the ρ relational encoding scheme of basis calculus.

Definition 86. Given any function f on $\times_{k \in [d]} [m_k]$ we define its value subspace as

$$V^f = \left\{ \langle \rho^f [X_f, X_{[d]}], W[X_f] \rangle [X_{[d]}] : W[X_f] \in \mathbb{R}^{|\text{im}(f)|} \right\}$$

Composition of functions f, h on distinct variables is then a choice of a subspace in the tensor product of the two function subspaces, that is

$$V^{f \circ h} \subset V^f \otimes V^h. \quad (84)$$

17.5.1 Propositional Formula Subspaces

Each formula defines by its relational encoding the subspace of $\bigotimes_{k \in [d]} \mathbb{R}^2$

$$V^f = \text{span}(\neg f, f) \quad (85)$$

Let us notice that the spanning vector of the subspace V^f are binary tensors summing up to the tensor of ones.

For each atom X_k we have

$$V^{X_k} = \mathbb{R}^2.$$

The tensor space carrying the factored representation of the worlds is thus

$$\bigotimes_{k \in [d]} V^{X_k}. \quad (86)$$

17.5.2 Formula Decomposition as a Subspace Choice

Given a formula $f \circ h$ composed of formulas f and h containing different atoms we have

$$V^{f \circ h} \subset V^f \otimes V^h \quad (87)$$

A connective \circ thus determines the selection of a two-dimensional subspace in the four-dimensional tensor product of subspaces to both subformulas.

Reconstruction of a formula given its formula tensor amounts to finding the HT Decomposition under the constraints of subspace choices according to the allowed logical connectives.

Given a set of positive and negative examples of a formula poses further an approximation problem of the examples by a HT Decomposition.

Advantages of this perspective are

- Given a HT the best approximation always exists (Theorem 11.58 in Hackbusch (2012)), but need to further restrict to cores given by logical connectives
- Apply Approximation algorithms: ALS or HOSVD

18 Reasoning by 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.

18.1 Approximation of Energy tensors

18.1.1 Direct Approximation

Direct approximation is the problem

$$\text{argmin}_{\theta \in \Gamma^G} \|E[X_{[d]}] - \theta[X_{[d]}\|^2.$$

18.1.2 Approximation involving Selection Architectures

Approximation involving a selection architecture \mathcal{H} is the problem

$$\operatorname{argmin}_{\theta \in \Gamma^{\mathcal{G}}} \|E - \langle \gamma^{\mathcal{H}}, \theta \rangle [X_{[d]}]\|^2.$$

In a tensor network diagram we depict this as

$$\operatorname{argmin}_{\theta \in \Gamma^{\mathcal{G}}} \left\| \begin{array}{c} \boxed{\gamma^{\mathcal{F}}} \begin{array}{c} L_{n-1} \\ \vdots \\ L_0 \end{array} \boxed{\theta} \\ X_0 \quad \dots \quad X_{d-1} \end{array} - \begin{array}{c} \boxed{Y} \\ X_0 \quad \dots \quad X_{d-1} \end{array} \right\|^2$$

Example 27 (Approximate based on a slice sparsity selecting architecture). *Use a term selecting neural network (conjunction neuron on d unary neurons selecting a variable and $\text{Id}, \neg, \text{True}$ as connective selector. Demand the parameter tensor θ 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.*

18.2 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.

18.2.1 Weighted Squares Loss Trick

Lemma 34. *Let T be a Boolean tensor, that is $\text{im}(T) \subset \{0, 1\}$. Then*

$$T[X_{[d]}] = \mathbb{I}[X_{[d]}] - (T[X_{[d]}] - \mathbb{I}[X_{[d]}])^2$$

where \mathbb{I} is a tensor with same shape as T and all coordinates being 1.

Proof. Since for each $x_{[d]} \in \times_{k \in [d]} [m_k]$ we have $T[X_{[d]} = x_{[d]}] \in \{0, 1\}$, it holds that

$$T[X_{[d]} = x_{[d]}] = 1 - (T[X_{[d]} = x_{[d]}] - 1)^2$$

and thus in coordinatewise calculus

$$T[X_{[d]}] = \mathbb{I}[X_{[d]}] - (T[X_{[d]}] - \mathbb{I}[X_{[d]}])^2.$$

□

We apply this property to reformulate optimization problems over boolean tensors into weighted least squares problems.

Theorem 118 (Weighted Squares Loss Trick). *Let Γ 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*

$$\operatorname{argmax}_{T \in \Gamma} \langle I, T \rangle [\emptyset] = \operatorname{argmin}_{T \in \Gamma} \langle I, (T[X_{[d]}] - \mathbb{I}[X_{[d]}])^2 \rangle [\emptyset] \quad (88)$$

Proof. Using the Lemma above, T is identical to $\mathbb{I}[X_{[d]}] - (T[X_{[d]}] - \mathbb{I}[X_{[d]}])^2$ and we get

$$\langle I, T \rangle [\emptyset] = \langle I, \mathbb{I}[X_{[d]}] \rangle [\emptyset] - \langle I, (T[X_{[d]}] - \mathbb{I}[X_{[d]}])^2 \rangle [\emptyset]$$

Since the first term does not depend on T , it can be dropped in the maximization problem. The (-1) factor then turns the maximization into a minimization problem. □

Theorem 118 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 Γ to differentially parametrizable sets, since then common squares risk solvers can be applied. We will call I in the Theorem 118 importance tensor, since it manipulates the relevance of each coordinate in the squares loss.

As a result, we interpret the objective

$$\langle I, (T[X_{[d]}] - \mathbb{I}[X_{[d]}])^2 \rangle [\emptyset]$$

as a weighted squares loss.

Example 28 (Proposal distribution maxima). *The Problem ?? of finding the maximal coordinate can thus be turned into*

$$\operatorname{argmax}_{l_{[n]}} \left\langle (\mathbb{P}^D - \tilde{\mathbb{P}}), \mathcal{H} \right\rangle [L_{[n]} = l_{[n]}] = \operatorname{argmin}_{l_{[n]}} \left\langle (\mathbb{P}^D - \tilde{\mathbb{P}}), (\langle \mathcal{H}, e_{l_{[n]}} [L_{[n]}] \rangle [X_{[d]}] - \mathbb{I} [X_{[d]}])^2 \right\rangle [\emptyset] .$$

18.2.2 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 $\Gamma^{\mathcal{G}, \mathbb{I}}$ for efficient solutions by alternating algorithms.

However, for any hypergraph \mathcal{G} we have $\mathbb{I} [X_{[d]}] \in \Gamma^{\mathcal{G}, \mathbb{I}}$. 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 29 (Fitting a tensor by a formula tensor). *Task: Given a tensor T , find a formula $f \in \mathcal{F}$ such that it coincides with T .*

If T 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

$$\operatorname{argmax}_{f \in \mathcal{F}} \langle f \Leftrightarrow T \rangle [\emptyset] .$$

We can use the squares risk trick and get an equivalent problem

$$\operatorname{argmin}_{f \in \mathcal{F}} \| \langle f \Leftrightarrow T \rangle [X_{[d]}] - \mathbb{I} [X_{[d]}] \|^2 .$$

18.3 Alternating Solution of Least Squares Problems

When the parameter tensor θ 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 10).

Algorithm 10 Alternating Least Squares (ALS)

for $e \in \mathcal{E}$ **do**

 Set $T^e [X_e]$ to a random element in \mathbb{R}^{p_k}

end for

while Stopping criterion is not met **do**

for $e \in \mathcal{E}$ **do**

 Set $T^e [X_e]$ to a solution of the local problem, that is

$$T^e [X_e] \leftarrow \operatorname{argmin}_{T^e [X_e]} \left\langle I, (\langle \mathcal{H}, \theta \rangle [X_{[d]}] - Y [X_{[d]}])^2 \right\rangle [\emptyset]$$

end for

end while

18.3.1 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 θ can be found by maximum calibration through a message passing algorithm, since no hidden variable has to be marginalized.

In case of skeleton expressions with many placeholders further decomposition for algorithmic efficiency are required.

- Elementary Format (EL-Format):
- CP-Format: Closest to sum of formula tensors (when all vectors are basis, then have a sum).
- TT-Format: Showed better heuristic performance in optimization

For any tensor network decomposition into cores θ^s have the derivative $\frac{\partial}{\partial \theta^s} \theta$ as the tensor network with out the core θ^s .

18.4 Regularization and Compressed Sensing

When regularizing the least squares problem by enforcing the sparsity of θ , we arrive at the compressed sensing problem

$$\operatorname{argmin}_{\theta[L]} \ell_0(\theta) \quad \text{subject to} \quad \|\langle \gamma^\phi, \theta \rangle [X_{[d]}] - E[X_{[d]}]\|_2 \leq \eta \quad (89)$$

Here, the sensing matrix is the selection tensor.

Example 30 (Formula fitting to an example). *Choosing the best formula fitting data (see Example 29) is the problem*

$$\operatorname{argmin}_{\theta[L]: \ell_0(\theta)=1} \|\langle I, \gamma^\phi, \theta \rangle [X_{[d]}] - Y\|_2 \quad (90)$$

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 31 (Propositional Formulas). *Let there be a set \mathcal{F} of formulas, then we have*

$$\langle \gamma^\mathcal{F} [X_{[d]}, L_{\text{in}}], \gamma^\mathcal{F} [X_{[d]}, L_{\text{out}}] \rangle [L_{\text{in}} = l_{\text{in}}, L_{\text{out}} = l_{\text{out}}] = \langle f_{l_{\text{in}}}, f_{l_{\text{out}}} \rangle [\emptyset] .$$

If the formulas have disjoint model sets then

$$\langle \gamma^\mathcal{F} [X_{[d]}, L_{\text{in}}], \gamma^\mathcal{F} [X_{[d]}, L_{\text{in}}] \rangle [L_{\text{in}} = l_{\text{in}}, L_{\text{out}} = l_{\text{out}}] = \begin{cases} \langle f_{l_{\text{in}}} \rangle [\emptyset] & \text{if } l_{\text{in}} = l_{\text{out}} \\ 0 & \text{else} \end{cases} .$$

Example 32 (Slice selection networks). *For the slice selection network*

$$\langle \mathcal{H}_{\wedge, d, r} [X_{[d]}, L_{\text{in}}], \mathcal{H}_{\wedge, d, r} [X_{[d]}, L_{\text{out}}] \rangle [L_{\text{in}} = l_{\text{in}}, L_{\text{out}} = l_{\text{out}}] = \begin{cases} 0 & \text{if for a } \tilde{k} \in A^{l_{\text{in}}} \cap A^{l_{\text{out}}} \text{ we have } x_{\tilde{k}}^{l_{\text{in}}} \neq x_{\tilde{k}}^{l_{\text{out}}} \\ \prod_{\tilde{k} \notin A^{l_{\text{in}}} \cup A^{l_{\text{out}}}} m_{\tilde{k}} & \text{else} \end{cases}$$

Given a fixed l_{in} , the maximum value in the respective slice is thus taken at $l_{\text{in}} = l_{\text{out}}$

19 Uniform Concentration of Random Contractions

We here derive bounds on the uniform concentration of contractions with random tensors.

The width of a vector η is the supremum of contractions with respect to a set Γ is

$$\omega_\Gamma(\eta) = \sup_{\theta \in \Gamma} \langle \theta, \eta \rangle [\emptyset] . \quad (91)$$

We are interested in the random vector

$$\eta^{\mathcal{F}, \mathbb{P}^*, D} = \langle \mathbb{P}^D, \mathcal{F} \rangle [L] - \langle \mathbb{P}^*, \mathcal{F} \rangle [L]$$

which is the difference between the mean parameters given the empirical distribution and the underlying generating distribution.

We derive bounds for the hypothesis Γ being the set of basis vectors (i.e. for feature search) and being the set of normed vectors (i.e. for feature calibration).

19.1 Naive bounds given binomial coordinates

We here investigate width bounds on random tensors η , which coordinates have marginal distributions by Binomials.

This is the case for $\eta^{\mathcal{H}, \mathbb{P}^*, m}$.

Naive means, that we do not exploit the dependencies of the coordinates on each other, as would be the case in more sophisticated chaining approaches.

19.1.1 Basis Vectors

We exploit the sub-Gaussian Norm (Def 2.5.6 in CITE Vershynin Book) to state concentration inequalities.

Definition 87 (Sub-Gaussian Norm). *The sub-Gaussian norm of a random variable X is defined as*

$$\|X\|_{\psi_2} = \inf \left\{ C > 0 : \mathbb{E} \left[\frac{X^2}{C^2} \right] \leq 2 \right\}.$$

Theorem 119. *Any coordinate of η is sub-Gaussian with*

$$\|\eta_i\|_{\psi_2} \leq \frac{C}{\sqrt{\ln 2}} \frac{1}{\sqrt{m}}$$

Proof. Centered Bernoulli is bounded and therefore sub-Gaussian. Binomial is a sum and we apply Proposition 2.6.1 in CITE Vershynin Book. \square

Theorem 120. *Let*

$$\Gamma = \{e_i : i \in [p]\}$$

and η be a random vector in \mathbb{R}^p , which coordinates have a marginal distribution being centered Binomials with a fixed $m \in \mathbb{N}$. Then

$$\|\omega_\Gamma(\eta)\|_{\psi_2} \leq C \sqrt{\frac{\ln p}{m}}.$$

where $C > 0$ is a universal constant.

Proof. Supremum of Sub-Gaussian variables. \square

19.1.2 Sphere

We first provide a Chebyshev bound on the width of the sphere.

Theorem 121. *Let η be a random tensor with marginal coordinate distributions by binomials with parameters $(p^{(x)}, m)$.*

For any $\epsilon > 0$, $\tau > 0$ and $m \in \mathbb{N}$ with probability at least $1 - \epsilon$ we have

$$\left\| \frac{\eta - \mathbb{E}[\eta]}{m} \right\|_2 \leq \tau$$

provided that

$$m \geq \frac{\sum_{x \in \times_{k \in [d]} [m_k]} p^{(x)}(1 - p^{(x)})}{\tau^2 \epsilon}.$$

Proof. Since the squared norm of the noise is the sum of squared centered and averaged Binomials, we have

$$\mathbb{E} \left[\|\eta - \mathbb{E}[\eta]\|_2^2 \right] = \sum_{x \in \times_{k \in [d]} [m_k]} \frac{p^{(x)}(1 - p^{(x)})}{m}$$

Here we used that the variance of Binomials with parameters $(p^{(x)}, m)$ is $p^{(x)}(1 - p^{(x)})m$.

It follows, that

$$\mathbb{E} \left[\left(\left\| \frac{\eta - \mathbb{E}[\eta]}{m} \right\|_2 \right)^2 \right] = \frac{\sum_{x \in \times_{k \in [d]} [m_k]} p^{(x)}(1 - p^{(x)})}{m}.$$

Then we apply a Chebyshev Bound to get for any $\tau > 0$

$$\mathbb{P} \left[\left\| \frac{\eta - \mathbb{E}[\eta]}{m} \right\|_2 > \tau \right] = \mathbb{P} \left[\left(\left\| \frac{\eta - \mathbb{E}[\eta]}{m} \right\|_2 \right)^2 > \tau^2 \right] \leq \frac{\sum_{x \in \times_{k \in [d]} [m_k]} p^{(x)}(1 - p^{(x)})}{m \cdot \tau^2} \quad (92)$$

For a $\epsilon > 0$ we choose any m with

$$m \geq \frac{\sum_{x \in \times_{k \in [d]} [m_k]} p^{(x)}(1 - p^{(x)})}{\tau^2 \epsilon}$$

and get

$$\mathbb{P} \left[\left\| \frac{\eta - \mathbb{E}[\eta]}{m} \right\|_2 > \tau \right] \leq \epsilon. \quad (93)$$

Thus, we have

$$\mathbb{P} \left[\left\| \frac{\eta - \mathbb{E}[\eta]}{m} \right\|_2 \leq \tau \right] = 1 - \mathbb{P} \left[\left\| \frac{\eta - \mathbb{E}[\eta]}{m} \right\|_2 > \tau \right] \geq 1 - \epsilon. \quad (94)$$

□

For $\phi = \delta$ the noise tensor is a rescaled and centered multinomial.

Corollary 15. *Let there be multinomial variable with parameters (p, m) where $p \in \bigotimes_{k \in [d]} \mathbb{R}^{m_k}$ a positive and normed tensor. Let D be a set of independent samples For any $\epsilon > 0$, $\tau > 0$ and $m \in \mathbb{N}$ with probability at least $1 - \epsilon$ we have*

$$\left\| \frac{\eta - \mathbb{E}[\eta]}{m} \right\|_2 \leq \tau$$

provided that

$$m \geq \frac{(1 - \langle (p)^2 \rangle [\emptyset])}{\tau^2 \epsilon}.$$

Proof. Theorem 121 with

$$\sum_{x \in \times_{k \in [d]} [m_k]} p^{(x)} (1 - p^{(x)}) = \sum_{x \in \times_{k \in [d]} [m_k]} p^{(x)} - \sum_{x \in \times_{k \in [d]} [m_k]} (p^{(x)})^2 = 1 - \langle (p)^2 \rangle [\emptyset].$$

□

19.1.3 Bounds based on the sub-gaussian norm

Unclear whether this is needed.

A faster tail decay can be achieved, when bounding sub-gaussian norms.

Theorem 122. *Let η be a random vector in \mathbb{R}^p , which coordinates have a marginal distribution being centered Binomials with a fixed $m \in \mathbb{N}$. Then*

$$\|\omega_S(\eta)\|_{\psi_2} \leq C \sqrt{\frac{p}{m}}.$$

where $C > 0$ is a universal constant.

Proof. Using that each coordinate has Sub-gaussian norm of at most 1. **Asymptotically, the binomial tends to a gaussian, which has a smaller sg norm. But the binomial has a sub-exponential regime preventing tighter sg bounds.**

Norm of a Sub-gaussian vector, another application of Proposition 2.6.1 in CITE Vershynin Book. □

19.2 Chaining bounds given binomial coordinates

To proceed with the uniform concentration investigation, we need a concentration bound on Binomials.

Theorem 123. *For any $p \in [0, 1]$ and $m \in \mathbb{N}$ any $X \sim B(p, m)$ satisfies for any $t > 0$*

$$\mathbb{P}[X - \mathbb{E}[X] > t] \leq \exp \left[-\frac{t^2}{2mp + \frac{2t}{3}} \right]$$

and

$$\mathbb{P}[\mathbb{E}[X] - X > t] \leq \exp \left[-\frac{t^2}{2mp} \right]$$

Thus

$$\mathbb{P}[|\mathbb{E}[X] - X| > t] \leq 2 \exp \left[-\frac{t^2}{2mp + \frac{2t}{3}} \right]$$

Proof. See e.g. <https://mathweb.ucsd.edu/~fan/wp/concen.pdf> □

The binomial thus has a sub-gaussian and a sub-exponential regime.

Theorem 124. For any $p \in [0, 1]$ and $m \in \mathbb{N}$ any $X \sim B(p, m)$ satisfies for any $t > 0$

$$\mathbb{P} \left[|\mathbb{E}[X] - X| > \sqrt{4mpt} + 2t \right] \leq 2\exp[-t]$$

Proof. For any $s > 0$ we choose $t > 0$ such that

$$s = -\frac{t^2}{2mp + \frac{t}{3}}$$

and observe

$$\min \left(\frac{t^2}{4mp}, \frac{t^2}{2t} \right)$$

and

$$t \leq \max(\sqrt{4mps}, 2s) \leq \sqrt{4mps} + 2s.$$

With the above bound it holds, that

$$\mathbb{P} \left[|\mathbb{E}[X] - X| > \sqrt{4mps} + 2s \right] \leq \mathbb{P} [|\mathbb{E}[X] - X| > t] \leq 2\exp \left[-\frac{t^2}{2mp + \frac{t}{3}} \right] \leq 2\exp[-s].$$

□

We apply this on the variable $\langle \theta, \eta \rangle$.

Theorem 125. Let $\theta = \sum_{f \in \mathcal{F}} \theta_f f$, then

$$\mathbb{P} \left[m |\langle \theta, \eta \rangle| \geq \sqrt{4m} \cdot \left(\sum_{f \in \mathcal{F}} |\theta_f| \sqrt{p^{(f)}} \right) \sqrt{t} + 2 \cdot \left(\sum_{f \in \mathcal{F}} |\theta_f| \right) t \right] \leq 2|\mathcal{F}| \cdot \exp[-t]$$

Proof. We can not assume independence of the $\langle f, \eta \rangle$ (in that case we could use a Bernstein inequality) and instead take the naive bound over all formulas in \mathcal{F}

$$\mathbb{P} \left[|\langle \theta, \eta \rangle| \geq \sqrt{4m} \cdot \sum_{f \in \mathcal{F}} |\theta_f| \sqrt{p^{(f)}} \sqrt{t} + 2 \cdot \sum_{f \in \mathcal{F}} |\theta_f| t \right] \tag{95}$$

$$\leq \mathbb{P} \left[\exists f \in \mathcal{F} : |\langle f, \eta \rangle| \geq \sqrt{4m} |\theta_f| \sqrt{p^{(f)}} \sqrt{t} + 2|\theta_f| t \right] \tag{96}$$

$$\leq \sum_{f \in \mathcal{F}} \mathbb{P} \left[|\langle f, \eta \rangle| \geq \sqrt{4m} |\theta_f| \sqrt{p^{(f)}} \sqrt{t} + 2|\theta_f| t \right] \tag{97}$$

$$\leq 2|\mathcal{F}| \cdot \exp[-t]. \tag{98}$$

□

We can thus proof uniform concentration bounds given covering bounds of a hypothesis set in ℓ_1 norm (and the reweighted one).

Remark 38 (Small Formula Probabilities). Directions with small $p^{(f)}$ will require larger covering sets and thus have large contributions to the bounds. Intuitively, they correspond with exceptional special cases, which need many samples to be observed. *Strange: Should also $1 - p^{(f)}$ small intuitively be an issue?*

19.2.1 Generic Width Bounds

We define the maps

$$\nu_2^p(\theta) = \inf_{(\mathcal{F}, \theta): \sum_{f \in \mathcal{F}} \theta_f f = \theta} \left(\sum_{f \in \mathcal{F}} |\theta_f| \sqrt{p^{(f)}} \right)$$

and

$$\nu_1(\theta) = \inf_{(\mathcal{F}, \theta): \sum_{f \in \mathcal{F}} \theta_f f = \theta} \left(\sum_{f \in \mathcal{F}} |\theta_f| \sqrt{p^{(f)}} \right)$$

corresponding to the sub-gaussian and sub-exponential regimes.

Theorem 126. *Let \mathcal{F} be a set of formulas and $W \subset \mathbb{R}^{|\mathcal{F}|}$ a set of weight vectors. Then with probability at least $1 - |\mathcal{F}| C_1 \exp \left[-\frac{u^2}{2} \right]$ we have for the set*

$$\Gamma = \left\{ \sum_{f \in \mathcal{F}} \theta_f f : (\theta_f)_{f \in \mathcal{F}} \in W \right\}$$

the bound

$$\omega_\Gamma(\eta) \leq \frac{C_2 \gamma_2(\Gamma, \nu_2^p)}{\sqrt{m}} + \frac{C_3 \gamma_1(\Gamma, \nu_1)}{m}.$$

A Implementation in the `tnreason` package

We here document the implementation of the discussed concepts in the python package `tnreason`.

`tnreason` is an abbreviation of **tensor network reasoning**, 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 <https://github.com/EnexaProject/enexa-tensor-reasoning> or by

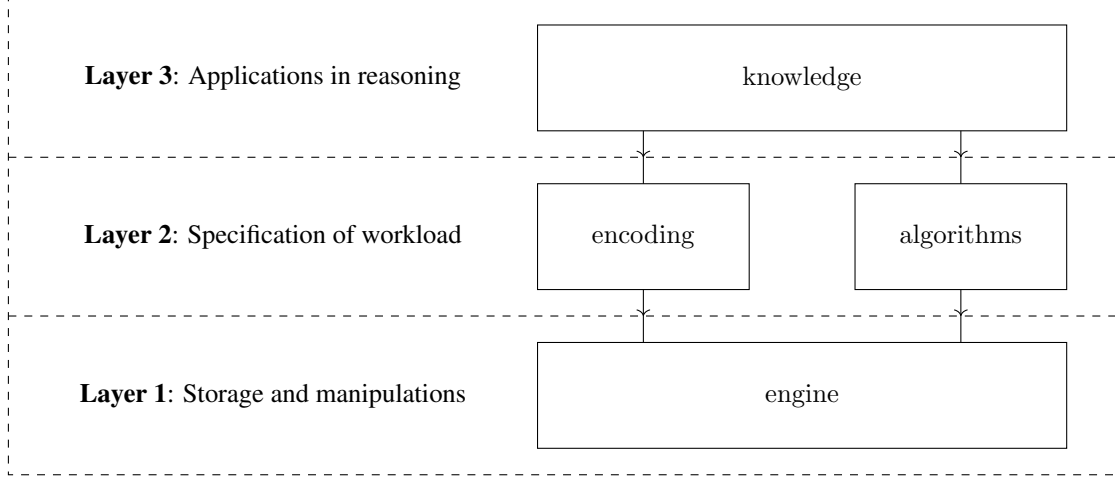
```
!pip install tnreason
```

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 encoding specific for storage, subpackage algorithms specific for manipulations
- Layer 3: Applications in reasoning, by subpackage `knowledge`, "Reasoning" -> building "reason" of `tnreason`

We sketch this structure by



A.2 Subpackage engine

The engine subpackage is for the storage and numerical manipulation of tensors and tensor networks.

We think of it as the lowest layer, specializing in storage of Tensor Networks and performing the contractions.

A.2.1 Contraction Calculus

We have described two main encoding schemes of functions, by a direct interpretation of functions as tensors or a more relational encoding. Both come with a different calculus scheme, which we have framed coordinate calculus and basis calculus.

A.2.2 Cores and Contractions

Cores

Each Tensor core has attributes

- values (array-like): storing the value of the coordinates
- colors (list of str): specifying the name of the variables represented by its axes
- name (str): to distinguish from other cores

The implemented core types differ in the values argument. Cores are instantiated by

```
engine.getCore(coreType)(coreValues, coreColors, coreName)
```

Polynomial Cores Polynomial Cores are implementations of the monomial decomposition or basis+ (see Definition 81). Here the each tuple (λ, A, X_A) is stored as a tuple of the scalar λ and a dictionary with A as keys and X_A as values.

The supported cores are

coreType	Package	Explanation
"NumpyTensorCore"	numpy	Numpy array storing the values
"PolynomialCore"	numpy	Storing the values in a binary CP Decomposition

Binary CP Decomposition

Based on the monomial decomposition $\text{rank}^d(\cdot)$ as specified in Definition 81. To store the values of a tensor we store the slices of tensors by the indices x_A .

Contractions can be performed by partially contracting the cores of the decomposition. In this way, one can avoid coordinatewise storages of high-order tensors, which can be intractable.

Tensor Networks

Tensor networks $\mathcal{T}^{\mathcal{G}}$ are defined by hypergraphs with hyperedges decorated by tensor cores. We store them by dictionaries with values being tensor cores and keys coinciding with the name of each tensor core.

Contractions

Reflected in the notation

$$\langle \mathcal{T}^{\mathcal{G}} \rangle [X_{\mathcal{V}}]$$

a contraction is defined by

- Tensor Network $\mathcal{T}^{\mathcal{G}}$, i.e. a dictionary of tensor cores
- Open Variables \mathcal{V}

Contraction calls are done by

```
engine.contract(contractionMethod, coreDict, openColors, dimensionDict, evidenceColorDict)
```

Where

- contractionMethod: str, chooses one of the contraction providers
- coreDict: Dictionary of TensorCores (of the above formats), representing the Tensor Network $\mathcal{T}^{\mathcal{G}}$
- 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

The supported contraction methods are

contractionMethod (str)	Package	Explanation
"NumpyEinsum"	numpy	Einstein summation of numpy arrays
"TensorFlowEinsum"	tensorflow	Einstein summation of tensorflow tensors
"TorchEinsum"	torch	Einstein summation of torch tensors
"TentrisEinsum"	tentris	Einstein summation of tentris hypertries
"PgmpyVariableEliminator"	pgmpy	Variable Elimination of DiscreteFactors in pgmpy
"PolynomialContractor"	numpy	Contraction of CP Decompositions stored in numpy arrays

Contractions represented as Einstein summation, as implemented in:

- numpy
- tensorflow
- pytorch
- tentris

Contractions can be executed by variable elimination as implemented in:

- pgmpy

Manipulation of Binary CP Decomposition Contraction of tensors in Binary CP Decomposition as in Section ??.

Coordinate Calculus

Main function

```
engine.coordinate_transform(coresList, transformFunction)
```

Basis Calculus

Main function

```
engine.relational_encoding()
```

basis calculus then based on contractions

A.3 Subpackage encoding

In the encoding subpackage we encode maps. Here the relational encodings ρ^f of various maps f are created. The maps are either specified by the script language (logical formulas or neuro-symbolic architectures), categorical constraints or data. 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.

We arrange the encoding subpackage into the second layer of the tnreason architecture, since it specifies tensor cores which formats are specified in engine .

A.3.1 Script Language

To specify propositional sentences, neuro-symbolic architectures and Markov Logic Networks, we developed a script language.

Propositional Sentences by Nested Lists

Are those of Propositional Logics, but instead of brackets we nest the symbols into lists.

Connectives are represented by strings, where the following are supported (see Definition ??):

Unary connective \circ	$S(\circ)$	Binary connective \circ	$S(\circ)$
\neg	"not"	\wedge	"and"
$()$	"id"	\vee	"or"
		\Rightarrow	"imp"
		\oplus	"xor"
		\Leftrightarrow	"eq"

Besides these specific connectives we exploit a generic representation scheme of propositional formulas by the so-called Wolfram code originally designed for the classification of cellular automaton rules Wolfram (1983) and popularized in the book Wolfram (2002). Along this, the coordinate encodings of connectives \circ with differing arity are flattened and interpreted as a binary number, which is transformed into a decimal number and represented as a string $S(\circ)$. We then choose a prefix to encode the arity by

- "u" for unary
- "b" for binary
- "t" for ternary
- "q" for quaternary

connectives. Together, the connective is represented by the string concatenation

$$S(\circ) = S(d) + S(\circ) .$$

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 $f_1 \circ, f_2$ are represented by

$$[S(\circ), S(f_1), S(f_2)]$$

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 applied grammar in Backus-Naur form is

Unary Connective	"not" "id"
Binary Connective	"and" "or" "imp" "xor" "eq"
Atomic Formula	Set of strings not in Connectives
Complex Formula	Atomic Formula [Unary Connective, Complex Formula] [Binary Connective, Complex Formula, Complex Formula]

Example 33 (Encoding of the Wet Street example). *For example we have*

- Atomic variable Rained by

$$S(\text{Rained}) = \text{"Rained"}$$

- Negative literal $\neg\text{Rained}$ 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 θ_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 7 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 Definition 48 and ??).

A connective selector (see Definition 48) is encoded by the list

$$S(\circ) = [S(\circ_0), \dots, S(\circ_{p-1})]$$

and a formula selector (see Definition ??) by

$$S(\mathcal{H}) = [S(\circ_0), \dots, S(\circ_{p-1})]$$

A logical neuron of order n (see Definition 52), defined by a connective selector \circ , and a formula selector \mathcal{H}_k on each argument $k \in [n]$, is encoded by

$$S(\sigma) = [S(\circ), S(\mathcal{H}_0), \dots, S(\mathcal{H}_{n-1})]$$

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 $S(\sigma)$ is the encoding of a neuron, then any element $x \in S(\sigma)$ 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

$$S(\mathcal{A}) = \{\text{key}(\sigma) : S(\sigma) \text{ for } \sigma \in \mathcal{A}\}$$

where $\text{key}(\sigma)$ 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 well-definedness of the architecture.

In order to represent neuro-symbolic architectures, the grammar of $S(\cdot)$ in Backus-Naur Form is extended by the production rules

Unary Connectives	[Unary Connective] [Unary Connective] + Unary Connectives
Binary Connectives	[Unary Connective] [Binary Connective] + Binary Connectives
Dependency Choice	Atomic Formula Neuron
Dependency Choices	[Dependency Choice] [Dependency Choice] + Dependency Choices
Neuron	[Unary Connectives, Dependency Choices] [Binary Connectives, Dependency Choices, Dependency Choices]

Example 34 (Neuro-Symbolic Architecture for the Wet Street). *Following the wet street example, we can define a neuron by*

$$S(\sigma) = [["imp", "eq"], ["Wet", "Sprinkler"], ["Street"]]$$

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

$$\begin{aligned} S(\mathcal{A}) = \{ &"neur1": [["imp", "eq"], ["neur2"], ["Street"]], \\ &"neur2": [["lnot", "id"], ["Wet", "Sprinkler"], ["Street"]] \} \end{aligned}$$

the expressivity increases. In this case, the further neuron provides the flexibility of the first atoms to be replaced by its negation.

A.3.2 Core Nomenclature

In encoding.suffixes we defined suffixes for the names of cores and colors, which highlight their origin and purpose.

Cores are named with suffixes based on their functionality

- "_conCore": logical connectives (relational encoding of the connective map)
- "_headCore": two-dimensional vectors representing of the activation core to a formula
- "_dataCore": Representing (relational encoding of the data map)
- "_catCore": Categorical constraint cores
- **Add: Formula selecting cores**

A.3.3 Color Nomenclature

Represent the three appearing types of variables:

- Categorical variables X .: cVar / aVar
- Selection variables L .: sVar
- Term variables O .: tVar

A.3.4 Relational encoding of formulas

Propositional formulas f are represented in three schemes:

- Script language $S(f)$ by nested lists (see Section A.3.1). Most practical to choose a formula from a neuro-symbolic architecture.
- Strings specifying the categorical variables X_f .
- Representation of formulas by tensor networks being contracted to ρ^f

Conversions of the formats:

- $S(f)$ to color by

```
encoding.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

```
encoding.create_raw_cores( $S(f)$ )
```

This creates the connective cores for the semantic representation of ρ^f . We encode them by

When encoding formulas with hard interpretation, we furthermore add a head core of type "truthEvaluation" since we have

$$f = \langle \rho^f, e_1 \rangle [X_f] .$$

A.3.5 Representation of MLNs

Structure Cores are binary cores relating the variables in a predefined way, which is not changing during reasoning.

- Logical interpretation: Cores ρ° **Structure Cores are those of the Bayesian Propositional Network**
- Categorical constraints: Cores ρ^Z

Activation Cores encode the weights of the formulas in a Markov Logic Network. For proper MLN only have unary cores, which we call headCores. Head cores with suffix "headCore" in name.

They are modified during reasoning: Selection of activation cores in structure learning, assigning a weight in parameter estimation.

A.3.6 Formula Selecting Maps

Encoding of Neurons according to Definition 52:

- Activation selection core with suffix "actCore" in name. Selection by variable with suffix "actVar"
- Selection of neurons as arguments with suffix "selCore" in name. Each argument of each neuron comes with a control variable with suffix "selVar".

Encoding of Formula Selecting Neural Networks (Definition 52) by creating all formula selecting neurons.

Skeleton expression (Definition ??) are stored with placeholderkeys and the candidatelists by dictionaries with the placeholderkeys and values being the possible symbols.

A.4 Subpackage algorithms

The algorithms subpackage implements basic tensor network algorithms with calls of specific execution in engine . As the encoding 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.4.1 Alternating Least Squares

- Tensor Network of Structure Cores
- Tensor Network of Parameter Cores
- List of importance cores

```
algorithms.ALS
```

A.4.2 Gibbs Sampling

- Tensor Network of Structure Cores
- Parameter cores: Variable tensor network cores representing basis vectors.
- List of importance cores

```
algorithms.Gibbs
```


A.4.3 Knowledge Propagation

algorithms.ConstraintPropagator

A.4.4 Energy-based Algorithms

algorithms.NaiveMeanField

algorithms.GenericMeanField

algorithms.EnergyBasedGibbs

A.5 Subpackage knowledge

With the knowledge subpackage we provide an interface for reasoning workload. It builds a third layer, since it used encoding to represent knowledge by tensor networks and algorithms in the execution of reasoning tasks.

A.5.1 Distributions

We encode Markov Networks by specifying a set of tensor cores. Each distribution needs to have a routine

.create_cores()

creating the factor cores and

.get_partition_function()

calculating the partition function. Although the partition function can be calculated by the contraction of all cores, we separate the method since there are situations where a faster calculation can be performed.

Empirical Distributions are distributions of sample data. We represent the values as a CP Format of data cores as specified in Section 3.8

knowledge.EmpiricalDistribution

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.

knowledge.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.5.2 Inference

By

knowledge.InferenceProvider

taking a distribution from the above as argument.

Probabilistic queries as specified Definition 33) by

.query(variableList, evidenceDict)

MAP queries by

.exact_map_query()

or by

`.annealed_sample()`

using Simulated Annealing (see Remark ??) to find an approximate maximum. The second method circumvents the creation of the coordinatewise representation of the distribution and circumvents therefore, at the expense of potentially approximative solutions, a bottleneck in case of many query variables.

Entailment from the distribution (Definition ??) 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.5.3 Parameter Estimation

EntropyMaximizer implements Algorithm 8, which is motivated by the maximum entropy principle (see Section 4.5.3) to optimize Markov Logic Networks. The class

`knowledge.EntropyMaximizer`

is initialized with the arguments

- `expressionsDict`: Dictionary of formulas in the format $S(f)$
- `satisfactionDict`: Dictionary of the satisfaction rates (mean parameters) to be matched by the optimal distribution

The optimization is then performed by

`.alternating_optimization(sweepNum, updateKeys)`

method, where the iteration in Algorithm 8 through the `updateKeys` is performed `sweepNum` times.

A.5.4 Structure Learning

Formula Booster chooses a formula given a formula selecting map.

`knowledge.FormulaBooster`

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

References

- 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>.
- Albert-László Barabási. *Network Science*. Cambridge University Press, Cambridge, illustrated edition edition, July 2016. ISBN 978-1-107-07626-6.
- Alexander Bigerl, Felix Conrads, Charlotte Behning, Mohamed Sherif, Muhammad Saleem, and Axel-Cyrille Ngonga Ngomo. Tetriz -A Tensor-Based Triple Store. September 2020.

- 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.
- 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.
- 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>.
- 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>.
- 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>.
- 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>.
- Patrick Geiß, 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>.
- 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.
- 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.
- Sepp Hochreiter. Toward a broad AI. *Communications of the ACM*, 65(4):56–57, April 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.
- 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. Series Title: Lecture Notes in Computer Science.
- 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.
- 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>.
- 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.
- 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/>.
- 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.
- 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/iay009.
- Ralph Tyrell Rockafellar. *Convex Analysis*. Princeton University Press, Princeton, reprint edition edition, January 1997. ISBN 978-0-691-01586-6.
- 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. Series Title: Lecture Notes in Computer Science.
- 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/CED3EEB903D9D8A16843CFA5AC4D577>. Publisher: Cambridge University Press.
- Theo Trouillon. Knowledge Graph Completion via Complex Tensor Factorization. 2017.
- 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].
- 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 Polytopes*. Springer, New York, 1995th edition edition, October 2013. ISBN 978-0-387-94365-7.