
REPRESENTATION OF DISCRETE MAXIMUM ENTROPY DISTRIBUTIONS AS TENSOR NETWORKS

RESEARCH NOTES IN THE ENEXA AND QROM PROJECTS

November 21, 2025

ABSTRACT

We here summarize results of the main report on maximum entropy distributions. The principle of maximum entropy serves as motivation for Computation Activation Networks. We then restrict to these distributions and study relative entropy minimization problems.

In the report, first discussion on the mean polytope are in Chapter 3, and the general maximum entropy problem in Chapter 4. Chapter 8 contains the discussion on maximum entropy distributions in case of boolean statistics.

Contents

1	Contents	2
2	Motivation	2
3	The Maximum Entropy Problem	2
3.1	Outlook	3
4	Tensor Notation	3
4.1	Computation-Activation Networks	3
4.2	CP decompositions	3
5	The mean polytope	4
5.1	Convex hull	4
5.2	Faces	4
6	Main results: Tensor network representation of maximum entropy distributions	8
6.1	Main result	8
6.2	Maximum entropy on the interior	9
6.3	Mean parameter on faces	9
7	Characterization for boolean statistics	10
7.1	Set of maximum entropy distributions	10

7.2 Example	11
8 Generic Base Measures	12
9 Mean as a Statistic	12

1 Contents

We in this paper provide tensor network representations

- Representation of any distribution with a sufficient statistics: Generic activation tensors.
- Representation of distributions with maximum entropy, in case of positive realizability: Elementary activation tensors
- Representation of generic maximum entropy distributions: CP activation tensors.

Now, we want to characterize the CP rank of the activation tensors

- Depends on the face of the mean polytope, which contains the mean parameter
- We have thus a well-defined "CP rank" of faces
- Largest faces and vertices have always CP rank of 1, intermediate faces can have larger CP rank

For boolean statistics we further provide insights for boolean statistics (see Chapter 8.5):

- Example of independent statistics (see Exa. 8.28): Always elementary activation tensors (hypercubes)
- Example of partition statistics (see Exa. 8.30):
- Generic criterion for elementary activation: "Cube-like" polytopes (see Def. 8.29)

2 Motivation

Maximum Entropy in Physics: E.g. Maxwell-Boltzmann distributions.

Maximum Entropy in Learning: Consider a learning problem where we want to estimate a model based on observed data. The maximum entropy problem principle approaches this problem by designing statistics of the data, which means shall be reproduced in the model, and choosing the model reproducing the means of the statistic with least structure. The entropy of a distribution quantifies the degree of structureless in a distribution and is therefore maximized to solve the learning task.

3 The Maximum Entropy Problem

The mean parameter of a distribution $\mathbb{P}[X_{[d]}]$ to a statistic $\mathcal{S} : \times_{k \in [d]} [m_k] \rightarrow \times_{s \in [n]} [p_s]$ is the vector $\mu[L] \in \mathbb{R}^p$ with the coordinates

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

We express the computation of the mean parameter in the contraction of the selection encoding $\sigma^{\mathcal{S}}[X_{[d]}, L]$ of \mathcal{S}

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

The maximum entropy problem given a mean parameter $\mu^*[L]$ is

$$\max_{\mathbb{P}[X_{[d]}] \in \Lambda^{\delta, \text{MAX}, \nu}} \mathbb{H}[\mathbb{P}[X_{[d]}]] \quad \text{subject to} \quad \langle \mathbb{P}[X_{[d]}], \sigma^{\mathcal{S}}[X_{[d]}, L] \rangle [L] = \mu^*[L] \quad (\text{P}_{\mathcal{S}, \mu, \nu})$$

where $\Lambda^{\delta, \text{MAX}, \nu}$ is the set of distributions, which are representable with respect to the base measure ν .

A quick argument shows, that maximum entropy distributions always have \mathcal{S} as a sufficient statistics.

Theorem 1. *Any maximum entropy distribution with respect to a moment constraint on \mathcal{S} and a base measure ν has the sufficient statistic \mathcal{S} .*

Proof. Let $\mathbb{P}[X_{[d]}]$ be a feasible distribution for the maximum entropy problem, which does not have a sufficient statistic \mathcal{S} . Then we find $x_{[d]}, \tilde{x}_{[d]} \in \times_{k \in [d]} [m_k]$ with $x_{[d]} \neq \tilde{x}_{[d]}$, $\mathcal{S}(x_{[d]}) = \mathcal{S}(\tilde{x}_{[d]})$, $\nu[X_{[d]} = x_{[d]}] \neq 0$, $\nu[X_{[d]} = \tilde{x}_{[d]}] \neq 0$ and $\mathbb{P}[X_{[d]} = x_{[d]}] \neq \mathbb{P}[X_{[d]} = \tilde{x}_{[d]}]$. We then define a distribution $\tilde{\mathbb{P}}[X_{[d]}]$ coinciding with $\mathbb{P}[X_{[d]}]$ except for the coordinates $x_{[d]}, \tilde{x}_{[d]}$, where we set

$$\tilde{\mathbb{P}}[X_{[d]} = x_{[d]}] = \tilde{\mathbb{P}}[X_{[d]} = \tilde{x}_{[d]}] = \frac{\mathbb{P}[X_{[d]} = x_{[d]}] + \mathbb{P}[X_{[d]} = \tilde{x}_{[d]}]}{2}$$

We notice that $\tilde{\mathbb{P}}[X_{[d]}]$ is also a feasible distribution with an larger entropy than $\mathbb{P}[X_{[d]}]$. Therefore, a distribution which does not have the sufficient statistic \mathcal{S} cannot be a maximum entropy distribution. \square

This shows that any maximum entropy distribution is in $\Lambda^{\mathcal{S}, \text{MAX}}$, where MAX is the maximal hypergraph $\text{MAX} = ([p], \{[p]\})$. We search for sparse representations of the corresponding activation tensors and investigate in which cases the maximum entropy distribution is also in $\Lambda^{\mathcal{S}, \mathcal{G}}$ for sparser hypergraphs \mathcal{G} .

3.1 Outlook

To prepare for the presentation of our main results we introduce

- Computation-Activation Networks: A tensor network architecture, which will be used to represent maximum entropy distributions
- Mean polytopes: Polytopes, which contain all realizable mean parameter vectors.

We will then show, that dependent on the position of the mean parameter in the mean polytope, we can characterize the corresponding maximum entropy distribution by a Computation-Activation Network.

4 Tensor Notation

4.1 Computation-Activation Networks

Given a statistic $\mathcal{S} : \times_{k \in [d]} [m_k] \rightarrow \times_{s \in [n]} [p_s]$ we build its basis encoding tensor

$$\beta^{\mathcal{S}}[Y_{[p]}, X_{[d]}] = \sum_{x_{[d]} \in \times_{k \in [d]} [m_k]} \epsilon_{\mathcal{S}(x_{[d]})}[Y_{[p]}] \otimes \epsilon_{x_{[d]}}[X_{[d]}] .$$

A computation network is any representation of $\beta^{\mathcal{S}}[Y_{[p]}, X_{[d]}]$ as a tensor network. These can be constructed in the case statistics being a composition of connective functions.

An activation tensor is $\tau[Y_{[p]}]$ and the Computation Activation Network of \mathcal{S} and τ the tensor

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

We are interested in decomposition formats of $\tau[Y_{[p]}]$, where we use sets of tensor networks $\mathcal{T}^{\mathcal{G}}$ on a hypergraph \mathcal{G} . The family of by \mathcal{S} and a \mathcal{G} computable distributions are

$$\Lambda^{\mathcal{S}, \mathcal{G}} = \{ \langle \beta^{\mathcal{S}}[Y_{[p]}, X_{[d]}], \tau[Y_{\mathcal{V}}] \rangle [X_{[d]} | \emptyset] : \tau[Y_{\mathcal{V}}] \in \mathcal{T}^{\mathcal{G}} \} .$$

4.2 CP decompositions

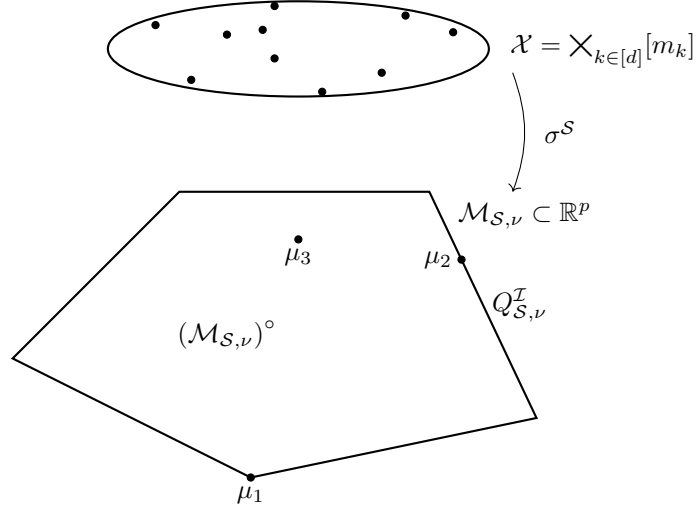
We here introduce the CP decomposition of tensors and the restriction to bas+. This will be used to represent face measures as computation activation networks.

5 The mean polytope

The mean polytope is the set of mean parameters to any distribution. We define it

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

where we denote by $\Lambda^{\delta, \text{MAX}, \nu}$ the set of all probability distributions representable with respect to ν .



5.1 Convex hull

The mean polytope is the convex hull

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

It is thus a convex polytope, inherited by the convex polytope of distributions (the standard simplex). We can characterize the maximum entropy distribution based on the position of the mean parameter in the mean polytope. To be more precise, any polytope decomposes into effective interiors of its faces and we characterize the maximum entropy distribution depending on the face to the mean parameter.

5.2 Faces

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

Definition 1. Given a mean parameter polytope $\mathcal{M}_{\mathcal{S},\nu}$ in the half space representation of Thm. ??, and any subset $\mathcal{I} \subset [n]$ we say that the set

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

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

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

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

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

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

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

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

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

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

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

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

Lemma 1. *For each face $Q_{S,\nu}^{\mathcal{I}}$ we have*

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

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

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

Lem. 1 implies in particular, that faces are mean parameter polytopes with respect to refined base measures. For reference in later chapters, we define these refined base measures next as face measures.

Definition 2. *The face measure to the face $Q_{S,\nu}^{\mathcal{I}}$ of $\mathcal{M}_{S,\nu}$ is the boolean tensor $\nu^{S,\mathcal{I}} [X_{[d]}]$ with coordinates to $x_{[d]} \in \times_{k \in [d]} [m_k]$ by*

$$\nu^{S,\mathcal{I}} [X_{[d]} = x_{[d]}] = \begin{cases} 1 & \text{if } \gamma^{x_{[d]}} \in Q_{S,\nu}^{\mathcal{I}} \\ 0 & \text{else} \end{cases}.$$

We now specify the mean parameter polytope to any face using the face measure as a refinement of the base measure.

Lemma 2. *For any face $Q_{S,\nu}^{\mathcal{I}}$ of $\mathcal{M}_{S,\nu}$, we have with the refined base measure*

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

that

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

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

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

In combination with Lem. 1 we then get

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

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

Lemma 3. *If and only if for a distribution $\mathbb{P} [X_{[d]}]$ and a face \mathcal{I} we have*

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

then $\mathbb{P} [X_{[d]}]$ is representable with respect to the base measure

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

Proof. We have

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

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

For members of exponential families, we can make a stronger statement than Lem. 3. If for any $\mathbb{P}[X[d]] \in \Gamma^{S, \nu}$ and a face \mathcal{I} we have $\langle \mathbb{P}[X[d]], \sigma^S[X[d], L] \rangle [L] \in (Q_{S, \nu}^{\mathcal{I}})^\circ$ then $\mathbb{P}[X[d]]$ is positive with respect to the base measure

$$\langle \nu[X[d]], \nu^{S, \mathcal{I}}[X[d]] \rangle [X[d]].$$

Let us now investigate tensor network representations of face measures, based on the basis encoding β^S of a statistic. Vertices of $\mathcal{M}_{S, \nu}$ are faces with single elements, that is $\{\mu[L]\}$. By Lem. 1 there must be μ must lie in the image of σ^S , since otherwise $\mathcal{M}_{S, \nu}$ would be empty. The vertex measure is then

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

Here we use that each $\mu \in Q_{S, \nu}^{\mathcal{I}} \cap \text{im}(\sigma^S)$ has integer-valued coordinates and denote

$$\epsilon_\mu[Y_{[p]}] = \bigotimes_{l \in [p]} \epsilon_{\mu[L=l]}[Y_l].$$

Theorem 2 (Face measure representation). *For any face $Q_{S, \nu}^{\mathcal{I}}$ of \mathcal{M} we have*

$$\nu^{S, \mathcal{I}}[X[d]] = \langle \beta^S[Y_{[p]}, X[d]], \kappa^{\mathcal{I}}[Y_{[p]}] \rangle [X[d]]$$

where

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

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

$$\tau^\mu[X[d]] = \langle \beta^S[Y_{[p]}, X[d]], \epsilon_\mu[Y_{[p]}] \rangle [X[d]]$$

is the indicator of the preimage of μ under σ^S . Since preimages the elements in $Q_{S, \nu}^{\mathcal{I}} \cap \text{im}(\sigma^S)$ are disjoint, the support of $\tau^\mu[X[d]]$ is disjoint and their sum

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

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

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

Motivated from the face measure representation, we define a CP rank for faces and show that normalized face measures are computable with respect to a corresponding CP format.

Definition 3. *The bas+ CP rank of a face is*

$$\text{rank}^{\text{bas}+}(Q_{S, \nu}^{\mathcal{I}}) = \min_{\{y_{[p]} : \epsilon_{y_{[p]}} \in Q_{S, \nu}^{\mathcal{I}}\} \subset \mathcal{U} \subset \epsilon(\times_{t \in [p]} [n_t]), \mathcal{U} \cup \{y_{[p]} : \epsilon_{y_{[p]}} \in \mathcal{M}_{S, \nu} / Q_{S, \nu}^{\mathcal{I}}\} = \emptyset} \text{rank}^{\text{bas}+} \left(\sum_{v \in \mathcal{U}} \epsilon_v[Y_{[p]}] \right).$$

The face measures are contraction of the vertex subset encodings with the computation. They are Computation-Activation Networks, when choosing the CP graph with rank at least $\text{rank}^{\text{bas}+}(Q_{\mathcal{S},\nu}^{\mathcal{I}})$.

Lemma 4. *For each face of the mean polytope we have*

$$\nu^{\mathcal{S},\mathcal{I}}[X_{[d]}|\emptyset] \in \Lambda^{\mathcal{S},\text{CP}^{\text{rank}^{\text{bas}+}(Q_{\mathcal{S},\nu}^{\mathcal{I}})}},$$

where $\text{CP}^{\text{rank}^{\text{bas}+}(Q_{\mathcal{S},\nu}^{\mathcal{I}})}$ is the CP graph with a hidden variable of dimension $\text{rank}^{\text{bas}+}(Q_{\mathcal{S},\nu}^{\mathcal{I}})$.

Proof. We find by definition a set \mathcal{U} of basis vectors containing the vertices of the face $\nu^{\mathcal{S},\mathcal{I}}$ but no further vertices, which has a bas+ CP rank of $\text{rank}^{\text{bas}+}(Q_{\mathcal{S},\nu}^{\mathcal{I}})$. We have therefore, that $\langle \beta^{\mathcal{S}}[Y_{[p]}, X_{[d]}], \epsilon_{\mathcal{U}}[Y_{[p]}] \rangle [X_{[d]}|\emptyset]$ is in $\Lambda^{\mathcal{S},\text{CP}^{\text{rank}^{\text{bas}+}(Q_{\mathcal{S},\nu}^{\mathcal{I}})}}$. Further it holds that

$$\begin{aligned} \langle \beta^{\mathcal{S}}[Y_{[p]}, X_{[d]}], \epsilon_{\mathcal{U}}[Y_{[p]}] \rangle [X_{[d]}] &= \sum_{y_{[p]} \in \mathcal{U} : \epsilon_{y_{[p]}} \in Q_{\mathcal{S},\nu}^{\mathcal{I}}} \langle \beta^{\mathcal{S}}[Y_{[p]}, X_{[d]}], \epsilon_{y_{[p]}}[Y_{[p]}] \rangle [X_{[d]}] \\ &\quad + \sum_{y_{[p]} \in \mathcal{U} : \epsilon_{y_{[p]}} \notin Q_{\mathcal{S},\nu}^{\mathcal{I}}} \langle \beta^{\mathcal{S}}[Y_{[p]}, X_{[d]}], \epsilon_v[Y_{[p]}] \rangle [X_{[d]}] \\ &= \sum_{y_{[p]} \in \mathcal{U} : \epsilon_{y_{[p]}} \in Q_{\mathcal{S},\nu}^{\mathcal{I}}} \langle \beta^{\mathcal{S}}[Y_{[p]}, X_{[d]}], \epsilon_{y_{[p]}}[Y_{[p]}] \rangle [X_{[d]}] \\ &= \nu^{\mathcal{S},\mathcal{I}}[X_{[d]}]. \end{aligned}$$

Here we used, that for $y_{[p]} \in \mathcal{U}$ with $\epsilon_{y_{[p]}} \notin Q_{\mathcal{S},\nu}^{\mathcal{I}}$ is not in the image of \mathcal{S} and therefore the contraction of its one-hot encoding with the computation cores vanishes. Thus, $\langle \beta^{\mathcal{S}}[Y_{[p]}, X_{[d]}], \epsilon_{\mathcal{U}}[Y_{[p]}] \rangle [X_{[d]}|\emptyset]$ coincides with the normalized face measure, which is therefore in $\Lambda^{\mathcal{S},\text{CP}^{\text{rank}^{\text{bas}+}(Q_{\mathcal{S},\nu}^{\mathcal{I}})}}$. \square

Let us now investigate, which normalized face measures can be computed using \mathcal{S} and a hypergraph \mathcal{G} .

Example 1 (Vertices). *Vertices $Q_{\mathcal{S},\nu}^{\mathcal{I}}$ are proper faces of affine dimension 0, that is they consist in single vectors. Since all vertices are in the image $\sigma^{\mathcal{S}}(\mathcal{X})$, there exists an index tuple $x_{[d]} \in \mathcal{X}$ such that $\nu[X_{[d]} = x_{[d]}] = 1$ and*

$$Q_{\mathcal{S},\nu}^{\mathcal{I}} = \{\sigma^{\mathcal{S}}[X_{[d]} = x_{[d]}, L]\}.$$

Then $\kappa^{\mathcal{I}}[Y_{[p]}]$ is the one-hot encoding of the by an interpretation map I assigned index to $\sigma^{\mathcal{S}}[X_{[d]} = x_{[d]}, L]$, that is

$$\kappa^{\mathcal{I}}[Y_{[p]}] = \epsilon_{I^{-1}(\sigma^{\mathcal{S}}[X_{[d]} = x_{[d]}, L])}[Y_{[p]}].$$

In particular, the activation core is elementary and the face measure to any vertex is in $\Lambda^{\mathcal{S},\text{EL},\nu}$.

Extending Example 1, we can provide a coarse estimation of the hypergraph \mathcal{G} required to decompose $\kappa^{\mathcal{I}}$ for generic faces $Q_{\mathcal{S},\nu}^{\mathcal{I}}$. We notice that $\kappa^{\mathcal{I}}[Y_{[p]}]$ in Thm. 2 is a sparse tensor with basis CP rank $|Q_{\mathcal{S},\nu}^{\mathcal{I}} \cap \text{im}(\sigma^{\mathcal{S}})|$ (see Chapter ??).

$$\ell_0(\kappa^{\mathcal{I}}[Y_{[p]}]) = |\sigma^{\mathcal{S}}(\nu^{\mathcal{S},\mathcal{I}})|$$

where $|\sigma^{\mathcal{S}}(\nu^{\mathcal{S},\mathcal{I}})|$ is the number of different statistic encoding vectors to the support of the face measure $\nu^{\mathcal{S},\mathcal{I}}$. Using the formalism of sparse CP decompositions Chapter ??, this characterize the basis CP rank of $\kappa^{\mathcal{I}}[Y_{[p]}]$. However, the basis CP rank is only an upper bound to generic CP rank, which can be loose. By Example 2 we provide with the maximal face an example, where the basis CP rank is given by the tensor space dimension, whereas the generic CP rank is one and the normalized face measure is thus still in $\Lambda^{\mathcal{S},\text{EL}}$.

Example 2 (Maximal face). *The maximal face $Q_{\mathcal{S},\nu}^{\emptyset} = \mathcal{M}_{\mathcal{S},\nu}$ coincides with the mean parameter itself. In this case the corresponding activation tensor to the face measure is trivial, that is*

$$\kappa^{\emptyset}[Y_{[p]}] = \mathbb{I}[Y_{[p]}].$$

κ^{\emptyset} is elementary and the normalized face measure $\nu^{\mathcal{S},\emptyset}$ to the maximal face is in $\Lambda^{\mathcal{S},\text{EL}}$.

Let us now introduce effective interiors, which enables us to find disjoint partitions of the mean polytope.

Definition 4 (Effective Interior). *Let $\mathcal{U} \subset \mathbb{R}^p$ be an arbitrary set and \mathcal{L} the minimal affine subspace of \mathbb{R}^p containing \mathcal{U} . Then the effective interior, denoted $(\mathcal{U})^\circ$ is the interior of \mathcal{U} in the space \mathcal{L} .*

Lemma 5. *Any polytope is a disjoint union of the effective interiors of its faces, that is*

$$\mathcal{M}_{\mathcal{S},\nu} = \bigcup_{\mathcal{I} \subset [n]} (Q_{\mathcal{S},\nu}^{\mathcal{I}})^\circ.$$

Proof. For any $\mu \in \mathcal{M}_{\mathcal{S},\nu}$ we find a face such that $\mu \in Q_{\mathcal{S},\nu}^{\mathcal{I}}$. If $\mu \notin (Q_{\mathcal{S},\nu}^{\mathcal{I}})^\circ$, then there is a face $Q_{\mathcal{S},\nu}^{\tilde{\mathcal{I}}} \subset Q_{\mathcal{S},\nu}^{\mathcal{I}}$ of smaller affine dimension such that $\mu \in Q_{\mathcal{S},\nu}^{\tilde{\mathcal{I}}}$. When continuing this process we reach a face such that $\mu \notin (Q_{\mathcal{S},\nu}^{\mathcal{I}})^\circ$, since the faces with affine dimension 0 are vertices and they coincide with their effective interior because they contain a single vector. \square

In this way, we find to each $\mu \in \mathcal{M}_{\mathcal{S},\nu}$ a unique exponential family with statistics \mathcal{S} and base measure by a face measure, such that μ is reproduced by an element of that exponential family. We will show in Chapter ??, that these reproducing distributions maximize the entropy among any other reproducing distribution.

6 Main results: Tensor network representation of maximum entropy distributions

Given the mean polytope discussion we now characterize the tensor network representation of maximum entropy distributions.

6.1 Main result

Theorem 3 (Generic characterization of Maximum Entropy Solutions). *Let \mathcal{S} be a statistic and ν a base measure. For any $\mu [L]$ the maximum entropy problem has a feasible distribution, if and only if $\mu [L] \in \mathcal{M}_{\mathcal{S},\nu}$. In case $\mu [L] \in \mathcal{M}_{\mathcal{S},\nu}$ there is a unique face $Q_{\mathcal{S},\nu}^{\mathcal{I}}$ such that μ is in the effective interior of $Q_{\mathcal{S},\nu}^{\mathcal{I}}$. Then the solution of the maximum entropy problem is the member*

$$\mathbb{P}(\mathcal{S}, B^{\mathcal{S}, \nu^{\mathcal{S}, \mathcal{I}}}(\mu), \nu^{\mathcal{S}, \mathcal{I}})$$

of the exponential family $\Gamma^{\mathcal{S}, \nu^{\mathcal{S}, \mathcal{I}}}$, where $\nu^{\mathcal{S}, \mathcal{I}}$ is the . If ν is an elementary Computation-Activation Network, then $\mathbb{P}(\mathcal{S}, B^{\mathcal{S}, \nu^{\mathcal{S}, \mathcal{I}}}(\mu), \nu^{\mathcal{S}, \mathcal{I}})$ is a Computation-Activation Network with respect to the CP graph of rank $\text{rank}(Q_{\mathcal{S},\nu}^{\mathcal{I}})$.

Proof. Feasibility Claim: If and only if $\mu [L] \in \mathcal{M}_{\mathcal{S},\nu}$ then there is by definition a by ν representable $\mathbb{P}[X_{[d]}]$ reproducing $\mu [L]$. Thus if and only if $\mu [L] \in \mathcal{M}_{\mathcal{S},\nu}$ there is a feasible distribution for the maximum entropy problem.

Characterization Claim: We use the following argumentation to show the second claim:

- By Lem. 5 for any μ we find a unique face $Q_{\mathcal{S},\nu}^{\mathcal{I}}$.
- By Lem. 3 all feasible distributions are representable by the with the face measure refined base measure. The maximum entropy solution is thus the same as for the $(\mathcal{S}, \mu, \langle \nu, \nu^{\mathcal{S}, \mathcal{I}} \rangle [X_{[d]}])$ instance, which we characterize in the following.
- By Lem. 2 the face $Q_{\mathcal{S},\nu}^{\mathcal{I}}$ coincides with the polytope $\mathcal{M}_{\mathcal{S}, \langle \nu, \nu^{\mathcal{S}, \mathcal{I}} \rangle [X_{[d]}]}$ and in particular μ is in the effective interior of that polytope.
- We can now apply Thm. 4 and get a characterization of the maximum entropy solution as a member of the exponential family.

Representation Claim: By Thm. 2 we can represent the face measure as a CP Computation-Activation Network. Since both the ν (by assumption) and the soft activation (always) is elementary, they do not change the CP rank when contracting to the face activating tensor of minimal rank. \square

6.2 Maximum entropy on the interior

A classical result states, that the maximum entropy distribution is in the exponential family $\Gamma^{\mathcal{S},\nu}$ (see e.g. Koller and Friedman).

Theorem 4. *If and only if μ^* is in the effective interior of $\mathcal{M}_{\mathcal{S},\nu}$, then the unique solution of the maximum entropy problem is the distribution*

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

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

Proof. By The 3.3 in Wainwright and Jordan, since by assumption

$$\mu[L] \in (\mathcal{M}_{\mathcal{S},\nu})^\circ,$$

there is a canonical parameter θ with

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

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

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

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

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

Therefore, if $\tilde{\mathbb{P}}$ does not coincide with $\mathbb{P}^{(\mathcal{S},\hat{\theta},\nu)}$, it is not a maximum entropy distribution. \square

Exponential families are in $\Lambda^{\mathcal{S},\text{EL}}$, if and only if $\langle \nu \rangle [X_{[d]}|\emptyset] \in \Lambda^{\mathcal{S},\text{EL}}$. If $\langle \nu \rangle [X_{[d]}|\emptyset] \in \Lambda^{\mathcal{S},\text{EL}}$ and $\mu[L] \in (\mathcal{M}_{\mathcal{S},\nu})^\circ$ we therefore have a sparse representation of the maximum entropy distribution with elementary activation tensors.

6.3 Mean parameter on faces

We always find a unique face of the polytope with the mean parameter being in the interior (see Lem. 5). Any distribution reproducing the mean parameter is realizable with respect to the face measure of that face (see Lem. ??). We conclude that the maximum entropy distribution of $\mu^*[L]$ with respect to \mathcal{S}, ν is also the maximum entropy distribution

Theorem 5. *Given \mathcal{S} and $\mu[L] \in \mathcal{M}_{\mathcal{S},\nu}$, let \mathcal{I} be the smallest face of $\mathcal{M}_{\mathcal{S},\nu}$ such that*

$$\mu[L] \in Q_{\mathcal{S},\nu}^{\mathcal{I}}.$$

Then the corresponding maximum entropy distribution is in $\Lambda^{\mathcal{S},\mathcal{G},\nu}$ if and only if the face measure (see Def. 2)

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

is in $\mathcal{T}^{\mathcal{G}}$.

Proof. By Thm. ?? the maximum entropy distribution is an element of the exponential family with by the face measure refined base measure $\tilde{\nu}$. Let $\theta[L]$ be a canonical parameter such that

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

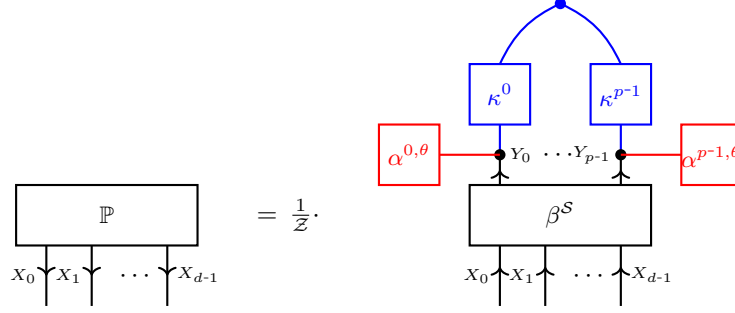


Figure 1: Tensor network decomposition of maximum entropy distributions to the constraint $\mu[L] = \langle \mathbb{P}, \sigma^S \rangle [L]$. Blue: Constraint activation cores κ^l in a CP decomposition, representing the face measure to the minimal face, such that $\mu \in Q_{S, \nu}^{\mathcal{I}}$. Red: Probabilistic activation cores $\alpha^{l, \theta} [Y_l]$ in an elementary decomposition, where each leg core is a scaled exponentials evaluated on the enumerated image $\text{im}(s_l)$.

that is $\mathbb{P}^{S, \theta, \bar{\nu}} [X_{[d]}]$ is the maximum entropy distribution. We apply Thm. 2 to represent the face measure by

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

Then for the tensor

$$\tau [Y_{[p]}] = \langle \{ \alpha^{l, \theta} [Y_l] : l \in [p] \} \cup \{ \kappa^{\mathcal{I}} [Y_{[p]}] \} \rangle [Y_{[p]}]$$

we have

$$\mathbb{P}^{S, \theta, \bar{\nu}} [X_{[d]}] = \langle \beta^S [Y_{[p]}, X_{[d]}], \tau [Y_{[p]}], \nu [X_{[d]}] \rangle [X_{[d]} | \emptyset] .$$

Thus, the maximum entropy distribution is in $\Lambda^{S, \mathcal{G}, \nu}$, if τ admits a tensor network decomposition with respect to \mathcal{G} . Since the hard activation cores are elementary, this is the case when $\kappa^{\mathcal{I}}$ admits a tensor network decomposition with respect to \mathcal{G} . \square

7 Characterization for boolean statistics

We here study the face CP ranks in case of boolean statistics. We further show that any elementary Computation-Activation Network to boolean statistics is a maximum entropy distribution.

For boolean statistics $\mathcal{F} : \times_{k \in [d]} [m_k] \rightarrow \times_{l \in [p]} [2]$ the mean polytope is a subset of the cube $[0, 1]^p$. In this case, any boolean vector in $\mathcal{M}_{\mathcal{F}, \nu}$ is a vertex. It follows, that any distribution reproducing a mean parameter $\mu [L]$ on the effective interior of $\mathcal{M}_{\mathcal{F}, \nu}$ is positive with respect to ν .

We apply the exponential distribution characterization of the maximum entropy distribution and get that the maximum entropy distribution is in $\Lambda^{S, \text{EL}}$, if and only if the face measure is in $\Lambda^{S, \text{EL}}$. This is exactly the case, when the face is an intersection of the mean polytope with a face of the cupe $[0, 1]^p$.

7.1 Set of maximum entropy distributions

Theorem 6. Any distribution in $\Lambda^{\mathcal{F}, \text{EL}}$ is a maximum entropy distribution. Any maximum entropy distribution is realized by $\Lambda^{\mathcal{F}, \text{EL}}$ if and only if the mean parameter is on an effective interior of a cube-like face.

Proof. First claim by decomposing any elementary tensor into exponential and hard activation core. Second claim by characterization of elementary faces by cube-likeness. \square

The mean parameters, which can be realized by a distribution in $\Lambda^{\mathcal{F}, \text{EL}}$ are those, which are on the effective interior of the intersection of the mean polytope with a face of the cube.

7.2 Example

Add examples: Partition Statistics - Simplices, (Logically) Independent formulas - Hypercubes

Example 3 (Maximum entropy distribution with non-elementary activation cores). *Consider two atomic variables X_0 and X_1 and a statistic \mathcal{F} consisting in the formulas*

$$f_0 = (X_0 \wedge X_1) \quad , \quad f_1 = (X_0 \Rightarrow X_1)$$

with the coordinatewise expressions

$$f_0 = \begin{bmatrix} 0 & 0 \\ 0 & 1 \end{bmatrix} \quad , \quad f_1 = \begin{bmatrix} 1 & 1 \\ 0 & 1 \end{bmatrix} .$$

We can think of X_0 as a feature on an invoice, and X_1 as a feature on the accounting proposal.

From this we have

$$\begin{aligned} \beta^{(f_0, f_1)} [Y_0 = 0, Y_1 = 0, X_0, X_1] &= \begin{bmatrix} 0 & 0 \\ 1 & 0 \end{bmatrix} \quad , \quad \beta^{(f_0, f_1)} [Y_0 = 0, Y_1 = 1, X_0, X_1] = \begin{bmatrix} 1 & 1 \\ 0 & 0 \end{bmatrix} \quad , \\ \beta^{(f_0, f_1)} [Y_0 = 1, Y_1 = 0, X_0, X_1] &= \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix} \quad \text{and} \quad \beta^{(f_0, f_1)} [Y_0 = 1, Y_1 = 1, X_0, X_1] = \begin{bmatrix} 0 & 0 \\ 0 & 1 \end{bmatrix} . \end{aligned}$$

Since the only vanishing slice of $\beta^{\mathcal{F}}$ with respect to the head variables is that to $y_{0,1} = (1, 0)$, the vertices of the mean polytope are the vectors to the other head indices. The mean polytope is the convex hull of these vertices

$$\mathcal{M}_{(f_0, f_1)} = \text{conv} \left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \end{bmatrix}, \begin{bmatrix} 1 \\ 1 \end{bmatrix} \right) .$$

This polytope has a non cube-like face (sketched blue in Figure 2), which is the convex hull of the vertices $[0\ 0]^T$, $[1\ 1]^T$. This face is parametrized by the (CP-rank 2) hard activation core

$$\kappa^{(0,0),(1,1)} [Y_0, Y_1] = \epsilon_{(0,0)} [Y_0, Y_1] + \epsilon_{(1,1)} [Y_0, Y_1] = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

and has the face measure

$$\left\langle \kappa^{(0,0),(1,1)} [Y_0, Y_1], \beta^{\mathcal{F}} [Y_0, Y_1, X_0, X_1] \right\rangle [X_0, X_1] = \begin{bmatrix} 0 & 0 \\ 1 & 1 \end{bmatrix} .$$

Any mean parameter μ on the interior of that face can be parametrized by a scalar $\lambda \in (0, 1)$

$$\mu_\lambda [L] = [\lambda \quad \lambda]^T .$$

With the canonical parameters $\theta [L] \in \mathbb{R}^2$ of the maximum entropy distributions on this face by

$$\mathbb{P} [X_0, X_1] = \frac{1}{1 + \exp [\theta [L = 0] + \theta [L = 1]]} \begin{bmatrix} 0 & 0 \\ 1 & \exp [\theta [L = 0] + \theta [L = 1]] \end{bmatrix}$$

we get the correspondence by the sigmoid

$$\lambda = \frac{1}{1 + \exp [-(\theta [L = 0] + \theta [L = 1])]} .$$

Note, that the hard activation core $\kappa^{(0,0),(1,1)} [Y_0, Y_1]$ to the blue face is the only non-elementary activation core. While the vertices have always elementary cores, the further non-vertex faces have elementary activation cores

$$\begin{aligned} \kappa^{(0,0),(1,0),(1,1)} [Y_0, Y_1] &= \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix} = \mathbb{I} [Y_0] \otimes \mathbb{I} [Y_1] \quad , \quad \kappa^{(0,0),(1,0)} [Y_0, Y_1] = \begin{bmatrix} 1 & 0 \\ 1 & 0 \end{bmatrix} = \mathbb{I} [Y_0] \otimes \epsilon_0 [Y_1] \quad , \\ \kappa^{(1,0),(1,1)} [Y_0, Y_1] &= \begin{bmatrix} 1 & 1 \\ 0 & 0 \end{bmatrix} = \epsilon_0 [Y_0] \otimes \mathbb{I} [Y_1] . \end{aligned}$$

The maximum entropy distributions to mean parameters on the interior of all other faces than the blue face are represented by Computation-Activation Networks with only elementary activation cores.

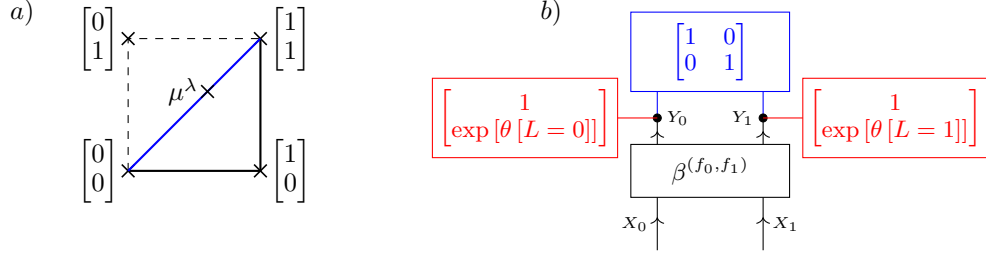


Figure 2: a) Mean polytope of the statistic $\mathcal{F} = (X_0 \wedge X_1, X_0 \Rightarrow X_1)$ (thick), as a subset of the cube $[0, 1]^2$ (dashed). The blue line is the face of the polytope, which is not cube like, that is not an intersection of the polytope with the faces of the polytope. We further define for $\lambda \in (0, 1)$ a mean parameter $\mu_\lambda [L] = [\lambda \lambda]^T$ which is on the interior of the blue face. b) Corresponding Computation-Activation Network being the maximum entropy distribution reproducing $\mu_\lambda [L]$, when λ is the sigmoid of $\theta [L = 0] + \theta [L = 1]$.

8 Generic Base Measures

The definition of the Shannon entropy is dependent on the chosen base measure. We now allow for generic non-Boolean base measures $\nu [X_{[d]}]$, which are non-negative and non-vanishing tensors. We define

- The set of by ν representable distributions

$$\Gamma^\nu = \{\mathbb{P} [X_{[d]}] : \forall x_{[d]} \in \bigtimes_{k \in [d]} [m_k] : \mathbb{P} [X_{[d]} = x_{[d]}] > 0, \langle \mathbb{P} [X_{[d]}], \nu [X_{[d]}] \rangle [\emptyset] = 1, \langle \mathbb{P} [X_{[d]}], \mathbb{I} [X_{[d]}] - \nu [X_{[d]}] \rangle [\emptyset] = 0\}$$

- The entropy of $\mathbb{P} [X_{[d]}] \in \Gamma^\nu$ by

$$\mathbb{H}^\nu [\mathbb{P}] = \langle \mathbb{P} [X_{[d]}], \ln [\mathbb{P} [X_{[d]}]] , \nu [X_{[d]}] \rangle [\emptyset]$$

- The mean parameter of $\mathbb{P} [X_{[d]}] \in \Gamma^\nu$ with respect to the statistic \mathcal{S} by

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

- Polytope of mean parameters by

$$\mathcal{M}_{\mathcal{S}, \nu} = \text{conv} (\sigma^\mathcal{S} X_{[d]} = x_{[d]}, L : \nu [X_{[d]} = x_{[d]}] \neq 0)$$

- The maximum entropy problem by

$$\text{argmax}_{\mathbb{P} \in \Gamma^\nu} \mathbb{H}^\nu [\mathbb{P}] \quad \text{subject to} \quad \mu [L] = \langle \mathbb{P} [X_{[d]}], \sigma^\mathcal{S} [X_{[d]}, L] , \nu [X_{[d]}] \rangle [L]$$

The main theorem of this work generalizes to this situation, with minor changes in the proofs.

9 Mean as a Statistic

The mean parameter μ_D given a dataset can be understood as a statistic of the dataset. We here show that for the family of maximum entropy distributions this statistic is a minimal sufficient statistic.

The family of maximum entropy distributions is the set

$$\{\mathbb{P}^\mu [X_{[d]}] : \mu \in \mathcal{M}_{\mathcal{S}, \nu}\}$$

which has been characterized above by a union of exponential families with respect to face measures.

Taking a frequentist perspective we now understand datasets by random variables $X_{[d] \times [m]}$, where for $j \in [m]$ the variables $X_{[d], j}$ are drawn i.i.d. from a maximum entropy distribution. The mean statistic is then a tensor

$$\mu_D [X_{[d] \times [m]}, L]$$

with coordinates

$$\mu_D [X_{[d] \times [m]} = x_{[d] \times [m]}, L] = \frac{1}{m} \sum_{j \in [m]} \sigma^\mathcal{S} [X_{[d], j} = x_{[d], j}, L] .$$

Theorem 7. *The mean statistic is sufficient for the family of maximum entropy distributions (\mathcal{S}, ν) .*

Proof. It suffices to show that the likelihood is a function of μ_D . Let us choose a face \mathcal{I} of $\mathcal{M}_{\mathcal{S}, \nu}$, then the likelihood is different from 0 if and only if the empirical distribution is representable with respect to the face measure. This is the case if and only if μ_D is on the face. In case that μ_D is on the face, then the likelihood of any distribution on that face exponential family is

$$\exp \left[m \cdot \left(\langle \mu_D [X_{[d] \times [m]}, L], \theta [L] \rangle [\emptyset] - A^{(\mathcal{S}, \nu)}(\theta) \right) \right]$$

We have thus shown that the likelihood is always a function of μ_D . □

References

- Daphne Koller and Nir Friedman. *Probabilistic Graphical Models: Principles and Techniques*. The MIT Press, 1. edition edition. ISBN 978-0-262-01319-2.
- Martin J. Wainwright and Michael Irwin Jordan. *Graphical Models, Exponential Families, and Variational Inference*. Now Publishers Inc. ISBN 978-1-60198-184-4.
- Günter M. Ziegler. *Lectures on Polytopes*. Springer, 1995th edition edition. ISBN 978-0-387-94365-7.