
MAXIMUM ENTROPY DISTRIBUTIONS AS COMPUTATION ACTIVATION NETWORKS

RESEARCH NOTES IN THE ENEXA PROJECT

July 31, 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.

1 Motivation

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.

2 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

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

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 $x_{[d]} \neq \tilde{x}_{[d]}$, $\mathcal{S}(x_{[d]}) = \mathcal{S}(\tilde{x}_{[d]})$, $\nu[X_{[d]} = x_{[d]}] \neq 0$, $\nu[X_{[d]} = \tilde{x}_{[d]}]$ 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

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

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 $\Lambda^{\mathcal{S}, \mathcal{G}^{\max}}$, where \mathcal{G}^{\max} is the maximal hypergraph $\mathcal{G}^{\max} = ([p], \{[p]\})$. 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 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 The mean polytope

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

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

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

4.1 Convex hull and faces

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.

Faces of the mean polytope are itself mean polytopes with respect to refined base measures.

4.2 Maximum entropy on the interior

A classical result states, that the maximum entropy distribution is in the exponential family $\Gamma^{\mathcal{S}, \nu}$.

Theorem 2. *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 [Wainright and Jordan, 2008], since by assumption

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

there is a canonical parameter θ with

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

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

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

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

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

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

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

4.3 Mean parameter on faces

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

5 Characterization for boolean statistics

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 $[0, 1]^p$.

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.