

Exponential families

CPSC 440/550: Advanced Machine Learning

`cs.ubc.ca/~dsuth/440/24w2`

University of British Columbia, on unceded Musqueam land

2024-25 Winter Term 2 (Jan–Apr 2025)

Previously: Density Estimation with Categorical/Gaussian Distributions

- We have discussed density estimation with **categorical and Gaussian** distribution
 - Bernoulli is a special case of categorical (up to notation changes)
- These distributions have a lot of **nice properties** for learning/inference
 - NLL is convex, and MLE has closed-form (statistics in training data)
 - A conjugate prior exists, so posterior is prior with “updated hyper-parameters”
- But these distributions make **restrictive assumptions**:
 - Categorical assumes categories are unordered, non-hierarchical, and finite
 - Gaussian assumes symmetry, full support, no outliers, uni-modal
- Many alternatives to categorical/Gaussian exist (examples later)
 - Alternatives that are in the **exponential family** maintain nice properties

Exponential Family: Definition

- General form of **exponential family** likelihood for data x with parameters θ is

$$p(x | \theta) = \frac{h(x) \exp(\eta(\theta)^\top s(x))}{Z(\theta)}$$

- The value $s(x)$ is the vector of **sufficient statistics**
 - $s(x)$ tells us everything that is relevant to θ about the data point x
- The **parameter function** η controls how parameters θ interact with the statistics
 - We'll focus on $\eta(\theta) = \theta$, which is called the **canonical form**
- The **support function** h contains terms that don't depend on θ
 - Also called the **base measure**
- The **normalizing constant** Z ensures it sums/integrates to 1 over x
 - Also called the **partition function**

Bernoulli as Exponential Family

- Is **Bernoulli** in the exponential family for some parameters w ?

$$p(x | \theta) = \theta^x (1 - \theta)^{1-x} \mathbb{1}(x \in \{0, 1\}) \stackrel{?}{=} \frac{h(x) \exp(\eta(\theta)^\top F(x))}{Z(\theta)}$$

- To introduce an exponential, also introduce a log so they cancel out:

$$\begin{aligned} p(x | \theta) &= \theta^x (1 - \theta)^{1-x} \mathbb{1}(x \in \{0, 1\}) \\ &= \exp(\log(\theta^x (1 - \theta)^{1-x})) \mathbb{1}(x \in \{0, 1\}) \\ &= \exp(x \log \theta + (1 - x) \log(1 - \theta)) \mathbb{1}(x \in \{0, 1\}) \\ &= (1 - \theta) \exp\left(x \log\left(\frac{\theta}{1 - \theta}\right)\right) \mathbb{1}(x \in \{0, 1\}) \end{aligned}$$

- The **sufficient statistic** is $s(x) = x$; normalizing constant is $Z(\theta) = 1/(1 - \theta)$
- The **parameter function** is $\eta(\theta) = \log(\theta/(1 - \theta))$ (the **log odds**)
 - Not in canonical form. Canonical form would use log odds directly as the parameter
- The **support function** is $h(x) = \mathbb{1}(x \in \{0, 1\})$ – says if we're “in the support”
- There are also **other ways to write Bernoulli as an exponential family**

Gaussian as Exponential Family

- One way to write univariate Gaussian as an exponential family:

$$\begin{aligned} p(x | \mu, \sigma^2) &= \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right) \\ &= \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{x^2}{2\sigma^2} + \frac{\mu x}{\sigma^2} - \frac{\mu^2}{2\sigma^2}\right) \\ &= \frac{1}{\sqrt{2\pi}} \frac{\exp\left(-\frac{\mu^2}{2\sigma^2}\right)}{\sigma} \exp\left(\begin{bmatrix} \mu/\sigma^2 \\ -1/2\sigma^2 \end{bmatrix}^\top \begin{bmatrix} x \\ x^2 \end{bmatrix}\right) \end{aligned}$$

- The sufficient statistics are x and x^2 , and parameters are μ/σ^2 and $-1/2\sigma^2$
- The normalizing constant is $\sigma \exp(\mu^2/2\sigma^2)$, and support is $1/\sqrt{2\pi}$

- Multivariate Gaussian looks roughly the same (with vec to flatten a matrix):

$$p(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \underbrace{\frac{1}{(2\pi)^{d/2}}}_{h(\mathbf{x})} \underbrace{\frac{\exp(-\frac{1}{2}\boldsymbol{\mu}^\top \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu})}{\log |\boldsymbol{\Sigma}|}}_{1/Z(\theta)} \exp\left(\underbrace{\begin{bmatrix} \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu} \\ \text{vec}(-\frac{1}{2}\boldsymbol{\Sigma}^{-1}) \end{bmatrix}^\top}_{\boldsymbol{\eta}(\theta)} \underbrace{\begin{bmatrix} \mathbf{x} \\ \text{vec}(\mathbf{x}\mathbf{x}^\top) \end{bmatrix}}_{s(\mathbf{x})}\right)$$

Learning with Exponential Families

- With n IID examples and canonical parameters η , the **likelihood** is

$$\begin{aligned} p(\mathbf{X} \mid \eta) &= \prod_{i=1}^n h(x^{(i)}) \frac{\exp(\eta^\top s(x^{(i)}))}{Z(\eta)} \\ &= \frac{1}{Z(\eta)^n} \exp\left(\eta^\top \sum_{i=1}^n s(x^{(i)})\right) \prod_{j=1}^n h(x^{(j)}) \\ &= \frac{\exp(\eta^\top s(\mathbf{X}))}{Z(\eta)^n} \prod_{j=1}^n h(x^{(j)}), \end{aligned}$$

defining sufficient statistics $s(\mathbf{X}) = \sum_{i=1}^n s(x^{(i)})$

- $s(\mathbf{X})$ contains **everything relevant for learning** – can **throw away the actual data**
 - For Gaussians, only knowledge of data we need is $\sum_{i=1}^n x^{(i)}$ and $\sum_{i=1}^n (x^{(i)})^2$
 - **No point in using SGD**: just compute s on each example **once**
 - Exponential families are the *only* class of distributions with a finite sufficient statistic

Learning with Exponential Families

- With iid data, canonical params η , **NLL** is $f(\eta) = -\eta^\top s(\mathbf{X}) + n \log Z(\eta) + \text{const}$
- The j th partial derivative of the NLL, divided by n , is

$$\begin{aligned}\frac{1}{n} \frac{\partial}{\partial \eta_j} f(\eta) &= -\frac{1}{n} s_j(\mathbf{X}) + \frac{1}{Z(\eta)} \frac{\partial}{\partial \eta_j} Z(\eta) \\ &= -\frac{1}{n} s_j(\mathbf{X}) + \frac{1}{Z(\eta)} \frac{\partial}{\partial \eta_j} \int h(x) \exp(\eta^\top s(x)) dx \quad (\text{use } \sum \text{ for discrete } x) \\ &= -\frac{1}{n} s_j(\mathbf{X}) + \int h(x) \frac{\exp(\eta^\top s(x))}{Z(\eta)} s_j(x) dx \quad (\text{w/ conditions}) \\ &= -\frac{1}{n} s_j(\mathbf{X}) + \int p(x | \eta) s_j(x) dx \\ &= -\mathbb{E}_{X \sim \text{data}} [s_j(X)] + \mathbb{E}_{X \sim \text{model } p_\eta} [s_j(X)]\end{aligned}$$

- The stationary points where $\nabla f(\eta) = 0$ correspond to **moment matching**:
 - Set parameters so that **expected sufficient statistics equal to statistics in data**
 - This is the source of the **simple/intuitive closed-form MLEs** we've seen so far

- If you take the second derivative of the NLL you get

$$\nabla^2 f(\eta) = \text{Cov}[s(X)],$$

the covariance of the sufficient statistics

- Covariances are positive semi-definite, $\text{Cov}[s(X)] \succeq 0$, so **NLL is convex**
- “Set the gradient to zero and solve” gives the MLE. . . for canonical params
- The NLL might *not* be convex in other parameterizations
 - e.g. multivariate Gaussians in terms of Σ
- Higher-order derivatives give higher-order moments
 - We call $\log(Z)$ the **cumulant function**
- Can show MLE **maximizes entropy over all distributions that match moments**
 - Entropy is a measure of “how random” a distribution is
 - So Gaussian is “most random” distribution that fits means and covariance of data
 - Or you can think of this as Gaussian makes “least assumptions”
 - Details for special case of $h(x) = 1$ in bonus slides

Conjugate Priors in Exponential Family

- Exponential families in canonical form are **guaranteed to have conjugate priors**
- For example, we could choose a prior like

$$p(\eta \mid \alpha) \propto \frac{\exp(\eta^\top \alpha)}{Z(\eta)^k}$$

- α is “**pseudo-counts**” for the sufficient statistics
- k **modifies the strength** of the prior (Z above is the likelihood's normalizer)
- Rewriting as $\exp(\eta^\top \alpha - k \log Z(\eta))$ shows this is itself an exponential family: canonical parameters (α, k) and sufficient stats $s(\eta) = (\eta, -\log Z(\eta))$
- Then the posterior has the same form,

$$p(\eta \mid \mathbf{X}, \alpha) \propto \frac{\exp(\eta^\top (s(\mathbf{X}) + \alpha))}{Z(\eta)^{n+k}}$$

- **Prior's normalizing constant** (some $\zeta_k(\alpha)$, **not** $Z(\eta)$) useful for Bayesian inference:
 - e.g. can derive, like before, that $p(\mathbf{X} \mid \alpha) = \zeta_{n+k}(s(\mathbf{X}) + \alpha) / \zeta_k(\alpha) \cdot \prod_{i=1}^n h(x^{(i)})$

Discriminative Models and the Exponential Family

- Going from an exponential family to a discriminative supervised learning:
 - Usual way is to **set canonical parameter to $w^T x$**
 - Gives a convex NLL, where MLE tries to match data/model's conditional statistics
 - Called **generalized linear model (GLM)** – see Stat 538A, Generalized Linear Models :)
- For example, consider Gaussian with **fixed variance** for y
 - Can write this with canonical parameter μ ; **setting $\mu = w^T x$ gives least squares**
- If we start with Bernoulli for y , we get **logistic regression**
 - Canonical parameter is log-odds, $\log(\theta)/\log(1 - \theta)$
 - Setting $w^T x = \log(\theta/(1 - \theta))$ and solving for θ gives $\theta = \sigma(w^T x)$
 - Gives a reason (sort of) for using the logistic sigmoid $\sigma(t) = 1/(1 + \exp(-t))$
- You can obtain regression models for other settings using this kind of approach
 - Set **canonical parameters to $f_\theta(x)$** , the output of a neural network
 - Use a **different exponential family** to handle a different type of data

Examples of Exponential Families

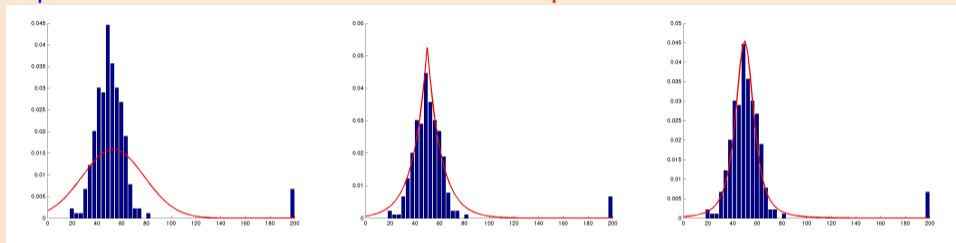
bonus!

- Bernoulli: distribution on $\{0, 1\}$
- Categorical: distribution on $\{1, 2, \dots, k\}$
- Multivariate Gaussian: distribution on \mathbb{R}^d
- Beta: distribution on $[0, 1]$ (including uniform)
- Dirichlet: distribution on discrete probabilities
- Wishart: distribution on positive-definite matrices
- Poisson: distribution on non-negative integers
- Gamma: distribution on positive real numbers
- Many, many others: [Wikipedia has a big table](#)
- ... can even have infinite-dimensional statistics via [kernel exponential families](#)

Non-Examples of Exponential Families

bonus!

- Laplace and student t distribution are **not exponential families**



- “Heavy-tailed”: have larger probability that data is far from mean
- **More robust** to outliers than Gaussian
- Ordinal logistic regression is **not in exponential family**
 - Can be used for categorical variables where **ordering matters**
- In these cases, we may not have nice properties:
 - **MLE may not be intuitive or closed-form, NLL may not be convex**
 - **May not have conjugate prior**, so need approximation

Summary

- Exponential families:
 - Have sufficient statistics and canonical parameters
 - Maximum likelihood becomes moment matching; always have conjugate priors
 - Can build discriminative models by using canonical parameter $s(x) = w^T x$
 - Many things (but not everything!) are exponential families

- Next time: mixing things up

- The **convex conjugate** of a function A is given by

$$A^*(\mu) = \sup_{w \in \mathcal{W}} \{\mu^\top w - A(w)\}$$

- For logistic regression, consider:

$$A(w) = \log(1 + \exp(w)),$$

then $A^*(\mu)$ satisfies $w = \log(\mu) / \log(1 - \mu)$

- When $0 < \mu < 1$ we have

$$\begin{aligned} A^*(\mu) &= \mu \log(\mu) + (1 - \mu) \log(1 - \mu) \\ &= -H(p_\mu), \end{aligned}$$

negative entropy of the Bernoulli distribution with mean μ

- If μ does not satisfy boundary constraint, sup is ∞

Convex Conjugate and Entropy

bonus!

- More generally, if $A(w) = \log(Z(w))$ for an exponential family then

$$A^*(\mu) = -H(p_\mu),$$

subject to boundary constraints on μ and the constraint

$$\mu = \nabla A(w) = \mathbb{E}[s(X)]$$

- Convex set satisfying these is called **marginal polytope** \mathcal{M}
- If A is convex (and lower semi-continuous), $A^{**} = A$. Then

$$A(w) = \sup_{\mu \in \mathcal{U}} \{w^\top \mu - A^*(\mu)\}$$

and when $A(w) = \log(Z(w))$ we have

$$\log(Z(w)) = \sup_{\mu \in \mathcal{M}} \{w^\top \mu + H(p_\mu)\}$$

- This can be used to derive variational methods, since we have written computing $\log(Z)$ as a convex optimization problem

- The **maximum likelihood** parameters w in exponential family satisfy:

$$\begin{aligned} & \min_{w \in \mathbb{R}^d} -w^\top s(\mathbf{X}) + \log(Z(w)) \\ &= \min_{w \in \mathbb{R}^d} -w^\top s(\mathbf{X}) + \sup_{\mu \in \mathcal{M}} \{w^\top \mu + H(p_\mu)\} \quad (\text{convex conjugate}) \\ &= \min_{w \in \mathbb{R}^d} \sup_{\mu \in \mathcal{M}} \{-w^\top s(\mathbf{X}) + w^\top \mu + H(p_\mu)\} \\ &= \sup_{\mu \in \mathcal{M}} \left\{ \min_{w \in \mathbb{R}^d} -w^\top s(\mathbf{X}) + w^\top \mu + H(p_\mu) \right\} \quad (\text{convex/concave}) \end{aligned}$$

which is $-\infty$ unless $s(\mathbf{X}) = \mu$ (e.g., maximum likelihood w), so we have

$$\min_{w \in \mathbb{R}^d} -w^\top s(\mathbf{X}) + \log(Z(w)) = \max_{\mu \in \mathcal{M}} H(p_\mu)$$

subject to $s(\mathbf{X}) = \mu$.

- **Maximum likelihood \Rightarrow maximum entropy + moment constraints**
- **Converse: MaxEnt + fit feature frequencies \Rightarrow ML(log-linear)**