

Exponential Families
oooooooooooo

Properties of Exponential Families
oooo

Learning in Exponential Families
oooooooooooo

HW 2: today
HW 3: wed 10/30

Quiz 6: next Fri, exp. families

COMPSCI 688: Probabilistic Graphical Models

Lecture 12: Learning in Exponential Families

Dan Sheldon

Manning College of Information and Computer Sciences
University of Massachusetts Amherst

Partially based on materials by Benjamin M. Marlin (marlin@cs.umass.edu) and Justin Domke (domke@cs.umass.edu)

1 / 27

Exponential Families
oooooooooooo

Properties of Exponential Families
oooo

Learning in Exponential Families
oooooooooooo

Exponential Families

2 / 27

Exponential Families
oooooooooooo

Properties of Exponential Families
oooo

Learning in Exponential Families
oooooooooooo

Exponential Families

An exponential family defines a set of distributions with densities of the form

$$p_\theta(x) = h(x) \exp(\theta^\top T(x) - A(\theta))$$

- θ : "(natural) parameters" $\in \mathbb{R}^d$
- $T(x)$: "sufficient statistics" $\in \mathbb{R}^d$
- $A(\theta)$: "log-partition function"
- $h(x)$: "base measure" (we'll usually ignore)

$h(x) \in I$

3 / 27

Exponential Families
oooooooooooo

Properties of Exponential Families
oooo

Learning in Exponential Families
oooooooooooo

Interpretation ($h(x) = 1$)

$x \mapsto T(x) = (x, x^2)$

$\theta = (\theta_1, \theta_2)$

$\theta_1 x + \theta_2 x^2$

$p_\theta(x) = \exp(\theta^\top T(x) - A(\theta))$

$\theta^\top T(x)$ is a real-valued "score" (positive or negative), defined in terms of "features" $T(x)$ and parameters θ

$\exp(\theta^\top T(x))$ is an unnormalized probability

The log-partition function $A(\theta) = \log Z(\theta)$ ensures normalization

$$p_\theta(x) = \frac{\exp(\theta^\top T(x))}{\exp(A(\theta))}, \quad A(\theta) = \log Z(\theta) = \log \int \exp(\theta^\top T(x)) dx$$

Valid parameters are the ones for which the integral for $A(\theta)$ is finite.

4 / 27

Applications and Importance

- ▶ We can get *many* different families of distributions by selecting different “features” $T(x)$ for a variable x in some sample space:
 - ▶ Bernoulli, Binomial, Multinomial, Beta, Gaussian, Poisson, MRFs, ...
- ▶ There is a general theory that covers learning and other properties of all of these distributions!
- ▶ A good trick to seeing that a distribution belongs to an exponential family is to match its log-density to

$$p_\theta(x) = h(x) \exp(\theta^\top T(x) - A(\theta))$$

$$\log p_\theta(x) = \log h(x) + \theta^\top T(x) - A(\theta)$$

5 / 27

Example: Bernoulli Distribution

$$x \mapsto T(x) = (\mathbb{I}[x=1], \mathbb{I}[x=0])$$

$$\theta = (\theta_1, \theta_2) \in \mathbb{R}^2$$

The Bernoulli distribution with parameter $\mu \in [0, 1]$ has density (pmf)

$$p_\mu(x) = \begin{cases} \mu & x = 1 \\ 1 - \mu & x = 0 \end{cases}$$

One way to write the log-density is

$$\log p_\mu(x) = \mathbb{I}[x=1] \log \mu + \mathbb{I}[x=0] \log(1 - \mu)$$

To match this to an exponential family

$$\log p_\theta(x) = \log h(x) + \theta^\top T(x) - A(\theta),$$

7 / 27

Preview: Graphical Models

For some intuition why exponential families could be relevant for graphical models, observe that the unnormalized probability factors over “simpler” functions, just like graphical models:

$$\exp(\theta^\top T(x)) = \exp \sum_i \theta_i T_i(x) = \prod_i \exp(\theta_i T_i(x))$$

(Think: what could $T(x)$ look like to recover a graphical model?)

6 / 27

This works (and is an interesting exercise), but uses two parameters where one would suffice. Instead...

8 / 27

Example: Bernoulli, Single Parameter

To write the Bernoulli as a single-parameter exponential family, rewrite the log-density as

$$\log p_\mu(x) = \log(1-\mu) + \mathbb{I}[x=1] \log \mu + \mathbb{I}[x=0] \log(1-\mu)$$

$$= \log(1-\mu) + \mathbb{I}[x=1] (\log \mu - \log(1-\mu))$$

$$\log p_\mu(x) = \log(1-\mu) + x \log \frac{\mu}{1-\mu}$$

$$- A(\theta) \quad T(x) \quad \theta$$

$T(x) = x$
 $\theta \in \mathbb{R}$ represent $\log \frac{\mu}{1-\mu}$ "log odds"

$$\exp(\theta \cdot x) = \begin{cases} e^\theta & x=1 \\ 1 & x=0 \end{cases}$$

$$A(\theta) = \log(1 + e^\theta) = \dots = \log(1-\mu) \text{ if } \theta = \log \frac{\mu}{1-\mu}$$

Review: Bernoulli, Single Parameter

- $h(x) = 1$
- $T(x) = \mathbb{I}[x=1] = x$
- $\theta = \log \frac{\mu}{1-\mu}$
- $\exp(\theta^\top x) = \begin{cases} e^\theta & x=1 \\ 1 & x=0 \end{cases}$
- $A(\theta) = \log(1 + e^\theta)$
- It's easy to check that $\log(1 + e^\theta) = -\log(1 - \mu)$ when $\theta = \log \frac{\mu}{1-\mu}$

Example: Normal Distribution



$$p_{\mu, \sigma^2}(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(x-\mu)^2\right)$$

$$= \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(x^2 - 2x\mu + \mu^2)\right)$$

$$\log p_{\mu, \sigma^2}(x) = x^2 \cdot \frac{1}{2\sigma^2} + x \cdot \frac{\mu}{\sigma^2} - \frac{\mu^2}{2\sigma^2} - \log \sqrt{2\pi\sigma^2}$$

$$T_1(x) \quad \theta_1 \quad T_2(x) \quad \theta_2 \quad -A(\theta)$$

$$T(x) = (x^2, x)$$

$$\theta = (\theta_1, \theta_2) \in \mathbb{R}^2 \text{ represent } \left(-\frac{1}{2\sigma^2}, \frac{\mu}{\sigma^2}\right)$$

$$A(\theta) = \log \int \exp(x^2 \theta_1 + x \theta_2) dx = \dots = \frac{\mu^2}{2\sigma^2} + \log \sqrt{2\pi\sigma^2}$$

$$\text{Need } \theta_1 < 0$$

Review: Normal Distribution

$$p_{\mu, \sigma^2}(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(x-\mu)^2\right)$$

$$= \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(x^2 - 2x\mu + \mu^2)\right)$$

$$\log p_{\mu, \sigma^2}(x) = x^2 \cdot \frac{-1}{2\sigma^2} + x \cdot \frac{\mu}{\sigma^2} - \frac{\mu^2}{2\sigma^2} - \log(\sqrt{2\pi\sigma^2})$$

$$\bullet h(x) = 1$$

$$\bullet T(x) = (x^2, x)$$

$$\bullet \theta = \left(\frac{-1}{2\sigma^2}, \frac{\mu}{\sigma^2}\right)$$

$$\bullet A(\theta) = \log \int \exp(x^2 \theta_1 + x \theta_2) dx = \dots = \frac{\mu^2}{2\sigma^2} + \log(\sqrt{2\pi\sigma^2})$$

Note: we need $\theta_1 < 0$; why?

Properties of Exponential Families

13 / 27

Properties of Log-Partition Function

14 / 27

The log-partition function $A(\theta)$ has two critical properties that relate its derivatives to moments (expectations) of the sufficient statistics $T(X)$.

$$\text{derivatives of } A(\theta) \Leftrightarrow \mathbb{E}[\text{function of } T(X)]$$

First Derivative of $A(\theta) \equiv$ First Moment of $T(X)$

$$\frac{\partial}{\partial \theta} \log \sum_x \exp(\theta^T T(x)) = \frac{1}{\sum_x \exp(\theta^T T(x))} \cdot \frac{\partial}{\partial \theta} \sum_x \exp(\theta^T T(x))$$

$$= \frac{1}{Z(\theta)} \sum_x \exp(\theta^T T(x)) \cdot \frac{\partial}{\partial \theta} \theta^T T(x)$$

$X \sim p_\theta$
compute $T(X)$
take mean

Proof: (assume $h(x) \equiv 1$)

$$\begin{aligned} \frac{\partial}{\partial \theta} \log \sum_x \exp(\theta^T T(x)) &= \frac{1}{\sum_x \exp(\theta^T T(x))} \cdot \frac{\partial}{\partial \theta} \sum_x \exp(\theta^T T(x)) \\ &= \frac{1}{Z(\theta)} \sum_x \exp(\theta^T T(x)) \cdot \frac{\partial}{\partial \theta} \theta^T T(x) \\ &= \sum_x \frac{\exp(\theta^T T(x))}{Z(\theta)} \cdot T(x) \\ &= \sum_x p_\theta(x) \cdot T(x) \\ &= \mathbb{E}_{p_\theta}[T(X)] \end{aligned}$$

15 / 27

Second Derivative of $A(\theta) \equiv$ Second Moment of $T(X)$

$$\begin{aligned} T(X) &= (T_1(x), \dots, T_d(x)) \\ \theta &= (\theta_1, \dots, \theta_d) \end{aligned}$$

(Variance)
dxd matrices
Hessian

$$A: \mathbb{R}^d \rightarrow \mathbb{R}$$

convex

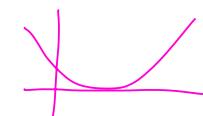
$$\frac{\partial^2}{\partial \theta \partial \theta^\top} A(\theta) = \text{Var}_{p_\theta}[T(X)]$$

Notation: $\frac{\partial^2}{\partial \theta \partial \theta^\top} A(\theta)$ is the Hessian matrix of $A(\theta)$. The (i, j) th entry is $\frac{\partial^2}{\partial \theta_i \partial \theta_j} A(\theta)$.

Proof: algebra

Important consequence: $A(\theta)$ is convex

- Variance is PSD \implies Hessian is PSD \implies A convex



16 / 27

Learning in Exponential Families

$$\text{Log-Likelihood} \quad x^{(1)}, \dots, x^{(N)} \quad \log p_\theta(x) = \log h(x) + \theta^\top T(x) - A(\theta)$$

The average log-likelihood in an exponential family is

$$\begin{aligned} \mathcal{L}(\theta) &= \frac{1}{N} \sum_{n=1}^N \log p_\theta(x^{(n)}) \\ &= \frac{1}{N} \sum_{n=1}^N (\theta^\top T(x^{(n)}) - A(\theta) + \log h(x^{(n)})) \\ &= \theta^\top \underbrace{\left(\frac{1}{N} \sum_{n=1}^N T(x^{(n)}) \right)}_{(\text{avg.}) \text{ sufficient}} - A(\theta) + \text{const.} \end{aligned}$$

- All we need to know about the data for estimation is the average value of $T(x^{(n)})$, i.e., the “sufficient statistics”

Moment-Matching

At the maximum-likelihood parameters, $\frac{\partial}{\partial \theta} \mathcal{L}(\theta) = 0$

$$\begin{aligned} 0 &= \frac{\partial}{\partial \theta} \mathcal{L}(\theta) = \frac{\partial}{\partial \theta} \left(\theta^\top \left(\frac{1}{N} \sum_{n=1}^N T(x^{(n)}) \right) - A(\theta) + \text{const} \right) \\ &= \frac{1}{N} \sum_{n=1}^N T(x^{(n)}) - \mathbb{E}_{p_\theta}[T(X)] \end{aligned}$$

⇒ at maximum-likelihood parameters, we have the *moment-matching conditions*:

$$\mathbb{E}_{p_\theta}[T(X)] = \frac{1}{N} \sum_{n=1}^N T(x^{(n)}) =: \hat{\mathbb{E}}[T(X)]$$

- “model expectation equals data expectation”
- sometimes we can easily solve for the maximum-likelihood parameters; other times numerical routines are needed

Concavity of Log-Likelihood

$$\frac{\partial^2}{\partial \theta \partial \theta^\top} \mathcal{L}(\theta) = -\frac{\partial^2}{\partial \theta \partial \theta^\top} A(\theta)$$

$$\mathcal{L}(\theta) = \underbrace{\theta^\top \left(\frac{1}{N} \sum_{n=1}^N T(x^{(n)}) \right)}_{\text{linear in } \theta} - A(\theta) + \text{const}$$



The log-likelihood is concave

- ⇒ every zero-gradient point is a global optimum
- ⇒ the moment-matching conditions are necessary and sufficient for optimality

Summary So Far

- ▶ $p_\theta(x) = h(x) \exp(\theta^\top T(\mathbf{x}) - A(\theta))$
- ▶ Bernoulli, normal, Poisson, MRF, ...
- ▶ First property: $\frac{\partial}{\partial \theta} A(\theta) = \mathbb{E}_{p_\theta}[T(X)]$
- ▶ Second property: $\frac{\partial^2}{\partial \theta \partial \theta^\top} A(\theta) = \text{Var}_{p_\theta}[T(X)]$
- ▶ Likelihood: $\mathcal{L}(\theta) = \theta^\top \bar{T} - A(\theta) + \text{const}$ where $\bar{T} = \frac{1}{N} \sum_{n=1}^N T(x^{(n)})$ are the average sufficient statistics over the data
- ▶ $\mathcal{L}(\theta)$ is concave
- ▶ Moment-matching conditions are necessary and sufficient for parameters θ to maximize the likelihood: $\mathbb{E}_{p_\theta}[T(X)] = \bar{T} = \hat{\mathbb{E}}[T(X)]$

21 / 27

23 / 27

Pairwise MRFs as an Exponential Family

Consider the chain model on $x_1, x_2, x_3, x_4 \in \{0, 1\}$:

$$p(\mathbf{x}) = \frac{\phi_{1,2}(x_1, x_2)\phi_{2,3}(x_2, x_3)\phi_{3,4}(x_3, x_4)}{Z}$$

23 / 27

Pairwise MRFs as an Exponential Family: Review

The log-density is

$$\begin{aligned} \log p(\mathbf{x}) &= \log \phi_{1,2}(x_1, x_2) + \log \phi_{2,3}(x_2, x_3) + \log \phi_{3,4}(x_3, x_4) - \log Z \\ &= \log \phi_{1,2}(0, 0) \cdot \mathbb{I}[x_1 = 0, x_2 = 0] + \log \phi_{1,2}(0, 1) \cdot \mathbb{I}[x_1 = 0, x_2 = 1] \\ &\quad + \log \phi_{1,2}(1, 0) \cdot \mathbb{I}[x_1 = 1, x_2 = 0] + \log \phi_{1,2}(1, 1) \cdot \mathbb{I}[x_1 = 1, x_2 = 1] \\ &\quad + \log \phi_{2,3}(0, 0) \cdot \mathbb{I}[x_2 = 0, x_3 = 0] + \dots \\ &\quad + \log \phi_{3,4}(0, 0) \cdot \mathbb{I}[x_3 = 0, x_4 = 0] + \dots \\ &\quad - \log Z \end{aligned}$$

24 / 27

24 / 27

This is an exponential family with

$$T(\mathbf{x}) = (\mathbb{I}[x_1 = 0, x_2 = 0], \dots, \mathbb{I}[x_1 = 1, x_2 = 1],$$

$$\mathbb{I}[x_2 = 0, x_3 = 0], \dots, \mathbb{I}[x_2 = 1, x_3 = 1],$$

$$\mathbb{I}[x_3 = 0, x_4 = 0], \dots, \mathbb{I}[x_3 = 1, x_4 = 1])$$

$$T(\mathbf{x}) = (\mathbb{I}[x_i = a, x_j = b])_{(i,j) \in E, a \in \text{Val}(X_i), b \in \text{Val}(X_j)}$$

$$\theta = (\theta_{ij}^{ab})_{(i,j) \in E, a \in \text{Val}(X_i), b \in \text{Val}(X_j)}$$

$$\log p_{\theta}(\mathbf{x}) = \theta^T \mathbf{x} - A(\theta) = \left(\sum_{(i,j) \in E} \sum_{a \in \text{Val}(X_i)} \sum_{b \in \text{Val}(X_j)} \theta_{ij}^{ab} \cdot \mathbb{I}[x_i = a, x_j = b] \right) - A(\theta)$$

The final three lines are accurate for general pairwise MRFs.

Moment-Matching for Pairwise-MRFs

If we apply the moment-matching conditions to pairwise MRFs, we recover our previous result. At the maximum-likelihood parameters:

$$\mathbb{E}_{p_{\theta}}[T(X)] = \hat{\mathbb{E}}[T(X)],$$

$$\mathbb{E}_{p_{\theta}}[\mathbb{I}[X_i = a, X_j = b]] = \hat{\mathbb{E}}[\mathbb{I}[X_i = a, X_j = b]] \quad \forall (i, j) \in E, a, b,$$

$$P_{\theta}(X_i = a, X_j = b) = \frac{\#(X_i = a, X_j = b)}{N} \quad \forall (i, j) \in E, a, b,$$

(we still have to solve for θ numerically; recall that the RHS minus the LHS is the gradient of $\mathcal{L}(\theta)$)

Moment-Matching for Gaussians

$x^{(1)}, \dots, x^{(n)}$

For a normal distribution, we had $T(x) = (x^2, x)$

$$\log p_{\mu, \sigma^2}(x) = x^2 \cdot \frac{-1}{2\sigma^2} + x \cdot \frac{\mu}{\sigma^2} - \frac{\mu^2}{2\sigma^2} - \log(\sqrt{2\pi\sigma^2})$$

We know $\mathbb{E}_{p_{\theta}}[X] = \mu$ and $\mathbb{E}_{p_{\theta}}[X^2] = \mu^2 + \sigma^2$. $\sigma^2 = \mathbb{E}[X^2] - \mu^2$

Moment-matching says the max-likelihood parameters satisfy:

$$\begin{aligned} \mathbb{E}_{p_{\theta}}[X] &= \hat{\mathbb{E}}[X] \implies \mu = \hat{\mathbb{E}}[X] \\ \mathbb{E}_{p_{\theta}}[X^2] &= \hat{\mathbb{E}}[X^2] \implies \mu^2 + \sigma^2 = \hat{\mathbb{E}}[X^2] \\ &\implies \sigma^2 = \hat{\mathbb{E}}[X^2] - \mu^2 \end{aligned}$$

We can easily solve for the maximum-likelihood μ, σ^2 .