

# Graphical Models

Exponential Families

# Families of Distributions

- A Gaussian distribution vs. the set of all Gaussian distributions
- A multinomial distribution over  $K$  outcomes vs. the set of all multinomials over the same  $K$  outcomes
- A graphical model vs. the set of all graphical models with the same structure and CPD parameterization (but different parameters)

# An Exponential Family

$$P_{\boldsymbol{\theta}}(\boldsymbol{x}) = \frac{1}{Z(\boldsymbol{\theta})} A(\boldsymbol{x}) \exp(\mathbf{t}(\boldsymbol{\theta})^\top \boldsymbol{\tau}(\boldsymbol{x}))$$

sufficient statistics function from  $\text{Val}(\mathbf{X})$  to  $\mathbb{R}^k$



# An Exponential Family

$$P_{\boldsymbol{\theta}}(\boldsymbol{x}) = \frac{1}{Z(\boldsymbol{\theta})} A(\boldsymbol{x}) \exp(\mathbf{t}(\boldsymbol{\theta})^\top \boldsymbol{\tau}(\boldsymbol{x}))$$

parameters in convex set  $\Theta \subseteq \mathbb{R}^M$



# An Exponential Family

$$P_{\boldsymbol{\theta}}(\boldsymbol{x}) = \frac{1}{Z(\boldsymbol{\theta})} A(\boldsymbol{x}) \exp(\mathbf{t}(\boldsymbol{\theta})^\top \boldsymbol{\tau}(\boldsymbol{x}))$$

natural parameter function from  $\mathbb{R}^M$  to  $\mathbb{R}^K$



parameters in convex set  $\Theta \subseteq \mathbb{R}^M$

# An Exponential Family

$$P_{\boldsymbol{\theta}}(\boldsymbol{x}) = \frac{1}{Z(\boldsymbol{\theta})} A(\boldsymbol{x}) \exp(\mathbf{t}(\boldsymbol{\theta})^\top \boldsymbol{\tau}(\boldsymbol{x}))$$

finite partition function



auxiliary measure



# An Exponential Family

$$P_{\boldsymbol{\theta}}(\boldsymbol{x}) = \frac{1}{Z(\boldsymbol{\theta})} A(\boldsymbol{x}) \exp(\mathbf{t}(\boldsymbol{\theta})^\top \boldsymbol{\tau}(\boldsymbol{x}))$$

$$Z(\boldsymbol{\theta}) = \sum_{\boldsymbol{x} \in \text{Val}(\mathbf{X})} A(\boldsymbol{x}) \exp(\mathbf{t}(\boldsymbol{\theta})^\top \boldsymbol{\tau}(\boldsymbol{x}))$$

# An Exponential Family

$$P_{\boldsymbol{\theta}}(\boldsymbol{x}) = \frac{1}{Z(\boldsymbol{\theta})} A(\boldsymbol{x}) \exp(\mathbf{t}(\boldsymbol{\theta})^\top \boldsymbol{\tau}(\boldsymbol{x}))$$

$$Z(\boldsymbol{\theta}) = \sum_{\boldsymbol{x} \in \text{Val}(\mathbf{X})} A(\boldsymbol{x}) \exp(\mathbf{t}(\boldsymbol{\theta})^\top \boldsymbol{\tau}(\boldsymbol{x}))$$

$$\mathcal{P} = \{P_{\boldsymbol{\theta}} \mid \boldsymbol{\theta} \in \Theta\}$$

parametric family



# Example: Bernoulli

$$P_{\boldsymbol{\theta}}(\mathbf{x}) = \frac{1}{Z(\boldsymbol{\theta})} A(\mathbf{x}) \exp(\mathbf{t}(\boldsymbol{\theta})^\top \boldsymbol{\tau}(\mathbf{x}))$$
$$Z(\boldsymbol{\theta}) = \sum_{\mathbf{x} \in \text{Val}(\mathbf{X})} A(\mathbf{x}) \exp(\mathbf{t}(\boldsymbol{\theta})^\top \boldsymbol{\tau}(\mathbf{x}))$$
$$\mathcal{P} = \{P_{\boldsymbol{\theta}} \mid \boldsymbol{\theta} \in \Theta\}$$

- $\Theta = [0, 1]$
- $\boldsymbol{\tau}(\mathbf{x}) = \langle \mathbf{x} == \text{heads}, \mathbf{x} == \text{tails} \rangle$
- $\mathbf{t}(\boldsymbol{\theta}) = \langle \ln \theta, \ln(1 - \theta) \rangle$ 
  - (note: image not convex!)
- $\exp(\mathbf{t}(\boldsymbol{\theta})^\top \boldsymbol{\tau}(\mathbf{x})) = \text{either } \theta \text{ or } 1 - \theta$
- $Z(\boldsymbol{\theta}) = 1$

# Example: Gaussian

$$P_{\boldsymbol{\theta}}(\mathbf{x}) = \frac{1}{Z(\boldsymbol{\theta})} A(\mathbf{x}) \exp(\mathbf{t}(\boldsymbol{\theta})^\top \boldsymbol{\tau}(\mathbf{x}))$$

$$Z(\boldsymbol{\theta}) = \sum_{\mathbf{x} \in \text{Val}(\mathbf{X})} A(\mathbf{x}) \exp(\mathbf{t}(\boldsymbol{\theta})^\top \boldsymbol{\tau}(\mathbf{x}))$$

$$\mathcal{P} = \{P_{\boldsymbol{\theta}} \mid \boldsymbol{\theta} \in \Theta\}$$

- $\Theta = \langle \mathbb{R}, \mathbb{R}^+ \rangle$
- $\boldsymbol{\tau}(\mathbf{x}) = \langle \mathbf{x}, \mathbf{x}^2 \rangle$
- $\mathbf{t}(\mu, \sigma^2) = \langle \mu / \sigma^2, -1 / 2\sigma^2 \rangle$

# Many Others!

- Poisson, exponential, geometric, Gamma, ...
- Sometimes more than one family can encode the same class of distributions.
  - Want  $\Theta$  to be convex and open
  - Want  $P_{\theta} = P_{\theta'}$  to imply that  $\theta = \theta'$  (nonredundant)
    - Equivalently,  $\mathbf{t}$  is invertible over  $\Theta$ .
- Counter examples?

# When $\mathbf{t}$ is the Identity Function

- “Natural parameters” for the sufficient statistic function  $\boldsymbol{\tau}$
- Should be familiar: **linear exponential family**

$$P_{\boldsymbol{\theta}}(\mathbf{x}) = \frac{1}{Z(\boldsymbol{\theta})} \exp\left(\boldsymbol{\theta}^{\top} \boldsymbol{\tau}(\mathbf{x})\right)$$
$$Z(\boldsymbol{\theta}) = \sum_{\mathbf{x} \in \text{Val}(\mathbf{X})} \exp\left(\boldsymbol{\theta}^{\top} \boldsymbol{\tau}(\mathbf{x})\right)$$

# Example: Natural Parameters for $\langle x, x^2 \rangle$

- $\boldsymbol{\eta} = \langle \mu / \sigma^2, -1 / 2\sigma^2 \rangle$

$$P_{\boldsymbol{\theta}}(\mathbf{x}) = \frac{1}{Z(\boldsymbol{\theta})} \exp\left(\boldsymbol{\theta}^\top \boldsymbol{\tau}(\mathbf{x})\right)$$
$$Z(\boldsymbol{\theta}) = \sum_{\mathbf{x} \in \text{Val}(\mathbf{X})} \exp\left(\boldsymbol{\theta}^\top \boldsymbol{\tau}(\mathbf{x})\right)$$
$$= \int_{-\infty}^{\infty} \exp(\eta_1 x + \eta_2 x^2) dx$$

# Example: Natural Parameters for $\langle x, x^2 \rangle$

- $\boldsymbol{\eta} = \langle \mu / \sigma^2, -1 / 2\sigma^2 \rangle$
- Z is only finite when  $\eta_2$  is negative!

$$\begin{aligned} P_{\boldsymbol{\theta}}(\mathbf{x}) &= \frac{1}{Z(\boldsymbol{\theta})} \exp\left(\boldsymbol{\theta}^\top \boldsymbol{\tau}(\mathbf{x})\right) \\ Z(\boldsymbol{\theta}) &= \sum_{\mathbf{x} \in \text{Val}(\mathbf{X})} \exp\left(\boldsymbol{\theta}^\top \boldsymbol{\tau}(\mathbf{x})\right) \\ &= \int_{-\infty}^{\infty} \exp(\eta_1 x + \eta_2 x^2) dx \end{aligned}$$

# Linear Exponential Families

- Basically boils down to the sufficient statistic function  $\tau$ 
  - Natural parameter space  $\Theta$  is whatever it needs to be to make  $Z$  finite

# Natural Parameter Space for the Bernoulli

$$P_{\boldsymbol{\theta}}(\mathbf{x}) = \frac{1}{Z(\boldsymbol{\theta})} A(\mathbf{x}) \exp(\mathbf{t}(\boldsymbol{\theta})^\top \boldsymbol{\tau}(\mathbf{x}))$$
$$Z(\boldsymbol{\theta}) = \sum_{\mathbf{x} \in \text{Val}(\mathbf{X})} A(\mathbf{x}) \exp(\mathbf{t}(\boldsymbol{\theta})^\top \boldsymbol{\tau}(\mathbf{x}))$$
$$\mathcal{P} = \{P_{\boldsymbol{\theta}} \mid \boldsymbol{\theta} \in \Theta\}$$

- $\boldsymbol{\tau}(\mathbf{x}) = \langle \mathbf{x} == \text{heads}, \mathbf{x} == \text{tails} \rangle$
- $\Theta = \mathbb{R}^2$
- $Z(\boldsymbol{\theta}) = \exp(\theta_1) + \exp(\theta_2)$
- Problem: overparameterization (add  $\mathbf{c}\mathbf{1}$  to  $\boldsymbol{\theta}$ )



# Natural Parameter Space for the Bernoulli

$$P_{\boldsymbol{\theta}}(\mathbf{x}) = \frac{1}{Z(\boldsymbol{\theta})} A(\mathbf{x}) \exp(\mathbf{t}(\boldsymbol{\theta})^\top \boldsymbol{\tau}(\mathbf{x}))$$
$$Z(\boldsymbol{\theta}) = \sum_{\mathbf{x} \in \text{Val}(\mathbf{X})} A(\mathbf{x}) \exp(\mathbf{t}(\boldsymbol{\theta})^\top \boldsymbol{\tau}(\mathbf{x}))$$
$$\mathcal{P} = \{P_{\boldsymbol{\theta}} \mid \boldsymbol{\theta} \in \Theta\}$$

- $\boldsymbol{\tau}(\mathbf{x}) = \langle \mathbf{x} == \text{heads}, \mathbf{x} == \text{tails} \rangle$
- $\Theta = \mathbb{R}$
- $Z(\theta) = 1 + \exp(\theta)$

# Natural Parameter Space for the Multinomial

$$P_{\boldsymbol{\theta}}(\mathbf{x}) = \frac{1}{Z(\boldsymbol{\theta})} A(\mathbf{x}) \exp(\mathbf{t}(\boldsymbol{\theta})^\top \boldsymbol{\tau}(\mathbf{x}))$$

$$Z(\boldsymbol{\theta}) = \sum_{\mathbf{x} \in \text{Val}(\mathbf{X})} A(\mathbf{x}) \exp(\mathbf{t}(\boldsymbol{\theta})^\top \boldsymbol{\tau}(\mathbf{x}))$$

$$\mathcal{P} = \{P_{\boldsymbol{\theta}} \mid \boldsymbol{\theta} \in \Theta\}$$

- $\boldsymbol{\tau}(\mathbf{x}) = \langle x_1, x_2, x_3, \dots, x_K \rangle$
- $\Theta = \mathbb{R}^{K-1}$
- What is  $\mathbf{t}$ ?

# Log-Linear Models

$$P_{\boldsymbol{\theta}}(\mathbf{x}) = \frac{1}{Z(\boldsymbol{\theta})} A(\mathbf{x}) \exp(\mathbf{t}(\boldsymbol{\theta})^{\top} \boldsymbol{\tau}(\mathbf{x}))$$
$$Z(\boldsymbol{\theta}) = \sum_{\mathbf{x} \in \text{Val}(\mathbf{X})} A(\mathbf{x}) \exp(\mathbf{t}(\boldsymbol{\theta})^{\top} \boldsymbol{\tau}(\mathbf{x}))$$
$$\mathcal{P} = \{P_{\boldsymbol{\theta}} \mid \boldsymbol{\theta} \in \Theta\}$$

- $\boldsymbol{\tau}(\mathbf{x}) = \langle f_1(\mathbf{d}_1), \dots, f_K(\mathbf{d}_K) \rangle$
- By extension, discrete Markov networks are exponential families.

Putting together exponential  
families ...

# Products

- It's easy to show that a product of exponential factor families' *unnormalized* probabilities gives ... another exponential family.
- Same thing for linear exponential factor families!
- Markov net as a Gibbs distribution

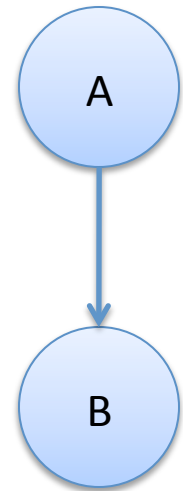
$$1/Z(\theta) \prod_{c \in \text{cliques}} \psi(c) = 1/Z(\theta) \prod_{c \in \text{cliques}} \exp(\eta^T f(c)) = 1/Z(\theta) \exp(\eta^T \sum_{c \in \text{cliques}} f(c))$$

# Bayesian Networks

- A Bayesian network with exponential CPDs is an exponential family.
  - Multinomials
  - Linear Gaussians
- However: local normalization is important here!

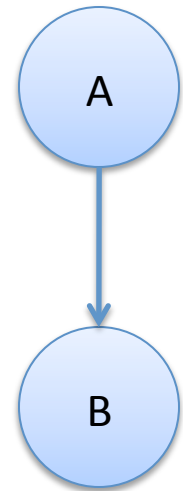
# Bayesian Network as an Exponential Family

- $\mathbf{t}_{P(A)} = \langle \ln P(A = 0), \ln P(A = 1) \rangle$
- $\mathbf{t}_{P(B | A)} = \langle \ln P(B = 0 | A = 0),$   
 $\ln P(B = 1 | A = 0),$   
 $\ln P(B = 0 | A = 1),$   
 $\ln P(B = 1 | A = 1) \rangle$
- 6 parameters, but only 3 degrees of freedom.



# Bayesian Network as an Exponential Family

- $\mathbf{t}_{P(A)} = \langle \ln P(A = 1) / \ln P(A = 0) \rangle$
- $\mathbf{t}_{P(B | A)} = \langle \ln P(B = 1 | A = 0) / \ln P(B = 0 | A = 0), \ln P(B = 1 | A = 1) / \ln P(B = 0 | A = 1) \rangle$
- Consider  $P_{\theta}(A = 1, B = 1)$  and  $P_{\theta}(A = 0, B = 0)$  ... what is  $Z(\theta)$ ?



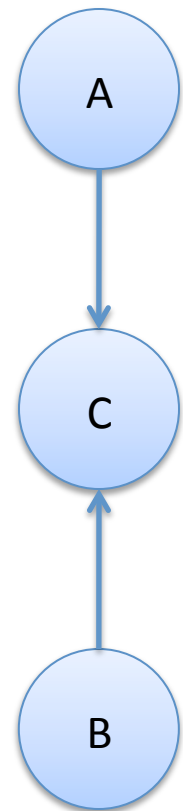


# Understanding the Failure

- The global normalizer  $Z(\theta)$  cannot do the work of local normalization within each conditional probability distribution.
- Need to make sure each CPD is normalized.

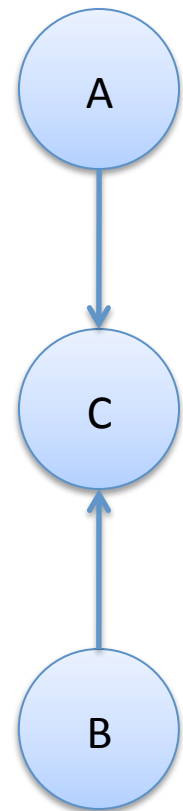
# Another Example

- Allowing for general CPDs, need the sufficient statistics to modify any of these without modifying the others:
  - $P(A = 0, B = 0, C = 0)$
  - $P(A = 0, B = 1, C = 0)$
  - $P(A = 1, B = 0, C = 0)$
  - $P(A = 1, B = 1, C = 0)$



# Another Example

- The  $\tau$  coordinates for these four outcomes must be linearly independent.
  - Without loss of generality, use indicator functions for those four outcomes for  $\tau$ .
- Resulting linear family is a Markov network over the clique  $\{A, B, C\}$ .
- Marginal independence between A and B may not hold!



# Are Bayesian Networks *Linear* Exponential Families?

- Not in general; any network with immoralities does not induce a *linear* exponential family.

# Some Consequences of Using an Exponential Family

- Entropy
- KL Divergence
- Projections

# Entropy

- Number of bits or nats needed on average to encode an outcome.

$$H(P(\mathbf{X})) = - \sum_{\mathbf{x} \in \text{Val}(\mathbf{X})} P(\mathbf{X} = \mathbf{x}) \ln P(\mathbf{X} = \mathbf{x})$$

- For an exponential family:

$$H(P_{\boldsymbol{\theta}}(\mathbf{X})) = \ln Z(\boldsymbol{\theta}) - \mathbf{t}(\boldsymbol{\theta})^{\top} \mathbb{E}_{P_{\boldsymbol{\theta}}}[\boldsymbol{\tau}(\mathbf{X})]$$

# Entropy of a Markov Network

$$H(P(\mathbf{X})) = \ln Z + \sum_k \mathbb{E}_P[-\log_2 \phi_k(\mathbf{D}_k)]$$

- (Are these quantities easy to compute?)

# Entropy of a Bayesian Network

$$H(P(\mathbf{X})) = \sum_i H(P(X_i | \text{Parents}(X_i)))$$

- (Are these quantities easy to compute?)



# Entropy of a Bayesian Network

$$H(P(\mathbf{X})) \leq \sum_i \max_{\mathbf{u} \in \text{Val}(\text{Parents}(X_i))} H(P(X_i | \text{Parents}(X_i) = \mathbf{u}))$$

$$H(P(\mathbf{X})) \geq \sum_i \min_{\mathbf{u} \in \text{Val}(\text{Parents}(X_i))} H(P(X_i | \text{Parents}(X_i) = \mathbf{u}))$$

- Bounds.

# KL Divergence (Relative Entropy)

$$D(P\|Q) = \sum_{\mathbf{x} \in \text{Val}(\mathbf{X})} P(\mathbf{X} = \mathbf{x}) \log \frac{P(\mathbf{X} = \mathbf{x})}{Q(\mathbf{X} = \mathbf{x})}$$

- A measurement of “distance” between two distributions.
  - Not symmetric.
- For exponential family P and any Q:

$$D(Q\|P_{\boldsymbol{\theta}}) = -H(Q(\mathbf{X})) - \mathbf{t}(\boldsymbol{\theta})^{\top} \mathbb{E}_Q[\boldsymbol{\tau}(\mathbf{X})] + \ln Z(\boldsymbol{\theta})$$

# KL Divergence

- Two distributions in the same exponential family:

$$D(P_{\boldsymbol{\theta}_1} \| P_{\boldsymbol{\theta}_2}) = (\mathbf{t}(\boldsymbol{\theta}_1) - \mathbf{t}(\boldsymbol{\theta}_2))^\top \mathbb{E}_{P_{\boldsymbol{\theta}_1}}[\boldsymbol{\tau}(\mathbf{X})] - \ln \frac{Z(\boldsymbol{\theta}_1)}{Z(\boldsymbol{\theta}_2)}$$

# Projections

- Given a distribution  $P$  and an exponential family  $\mathcal{Q}$ , find the distribution from  $\mathcal{Q}$  that is closest to  $P$ .
  - I-projection (information projection):

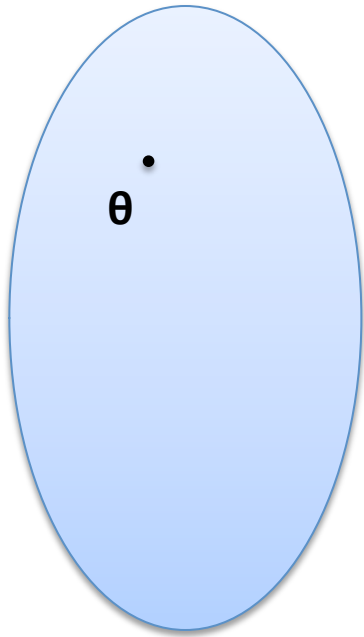
$$\arg \min_{Q \in \mathcal{Q}} D(Q \| P)$$

- M-projection (moment projection):

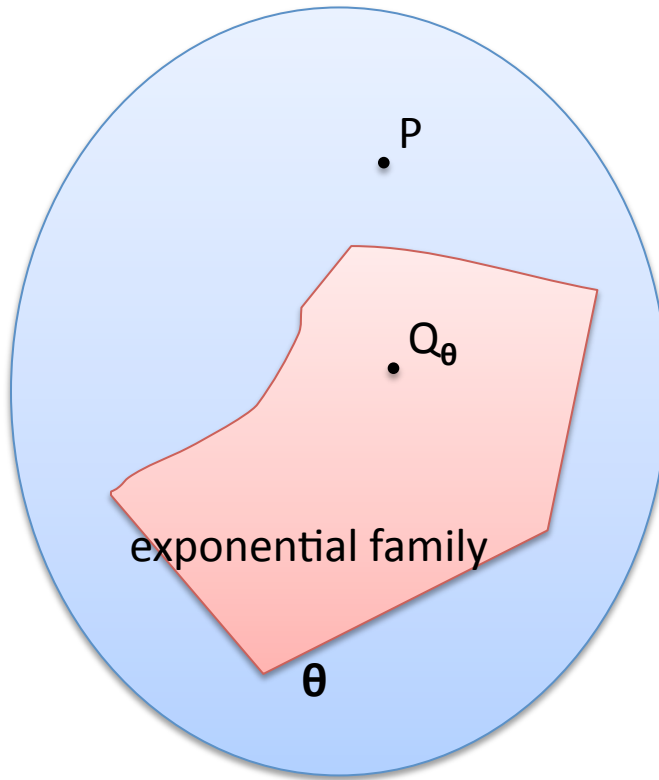
$$\arg \min_{Q \in \mathcal{Q}} D(P \| Q)$$

# M-Projection

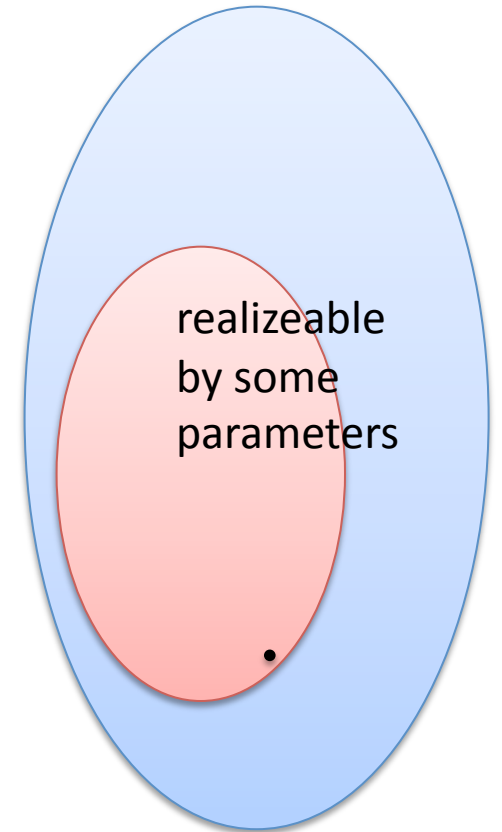
- Example: let  $\mathcal{Q}$  be the graph with no edges.
  - M-projection is the product of marginals from  $P$ .
- General result: if there is a  $\theta$  such that the expected sufficient statistics  $E_{Q_\theta}[\tau(\mathbf{X})]$  match  $E_P[\tau(\mathbf{X})]$ , then  $Q_\theta$  is the M-projection of  $P$ .



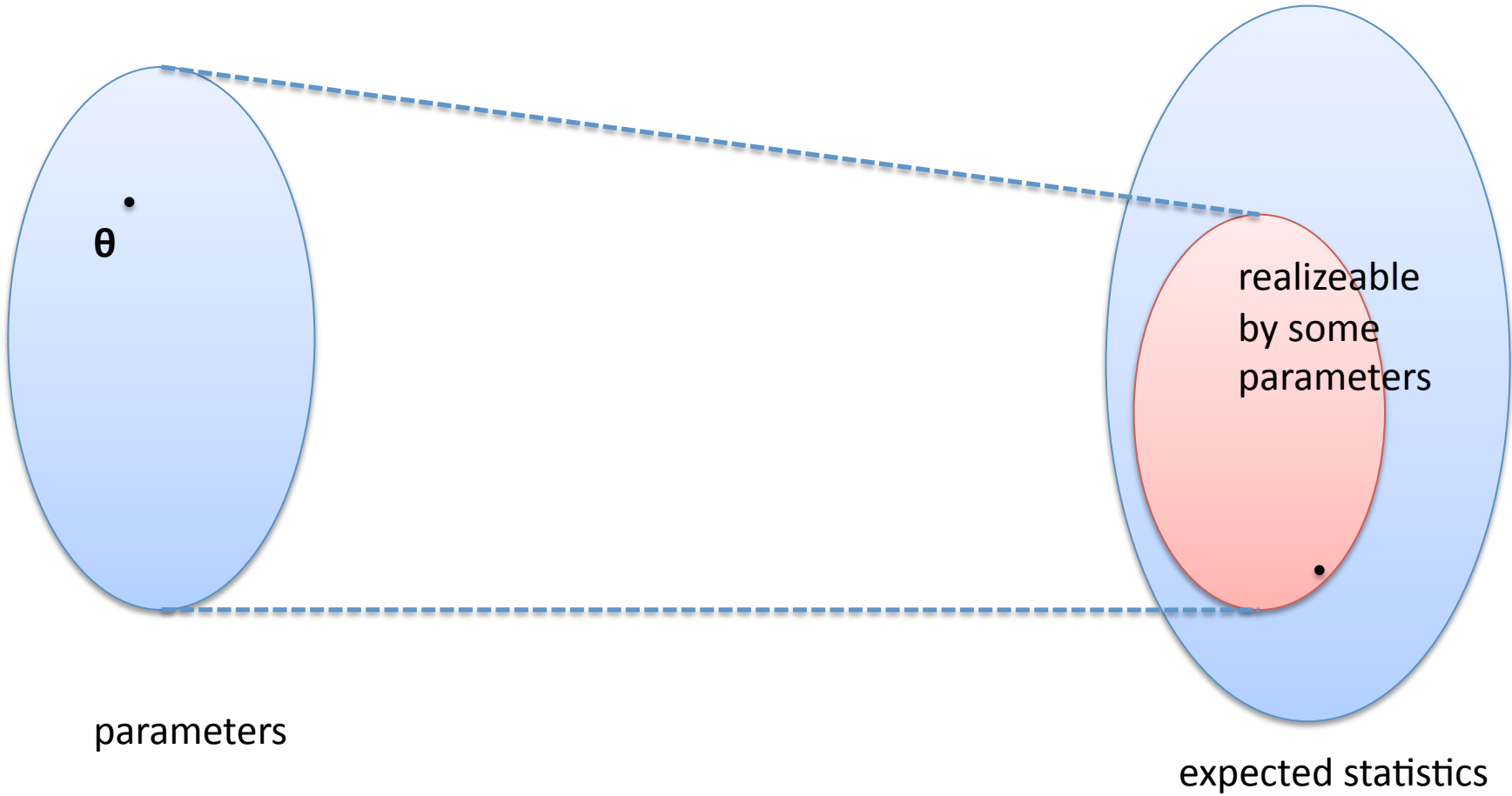
parameters

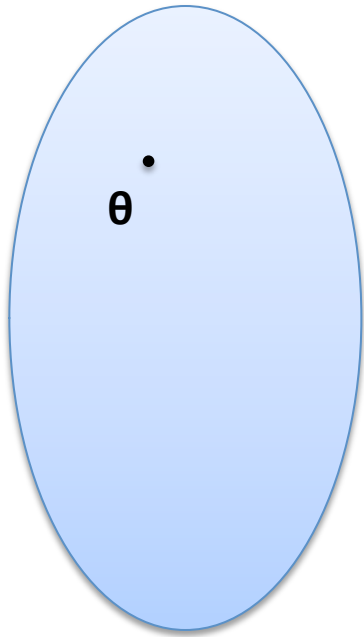


distributions

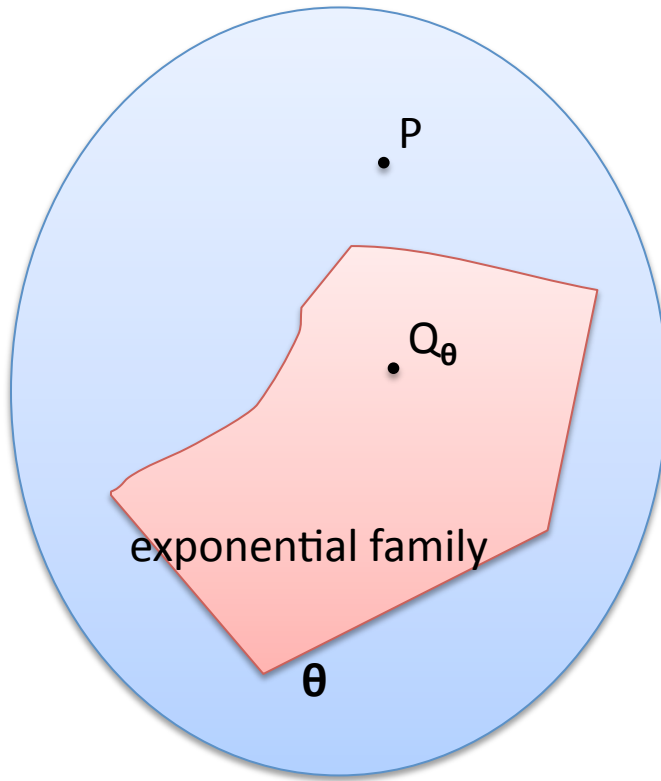


expected statistics

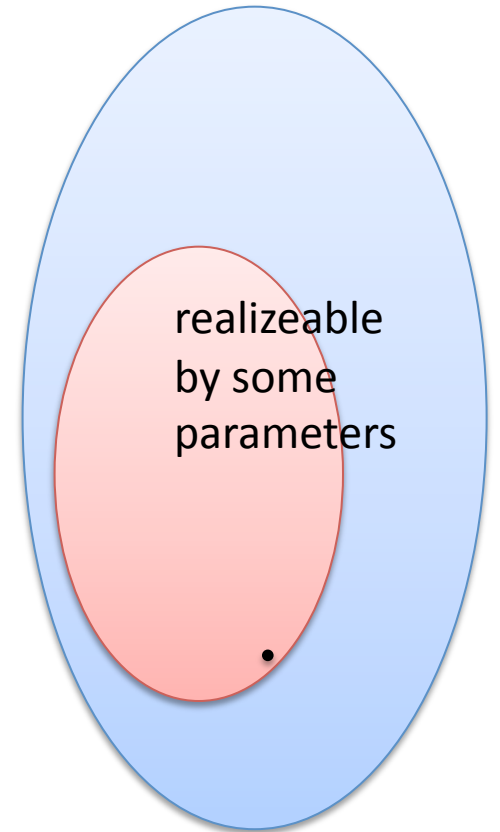




parameters

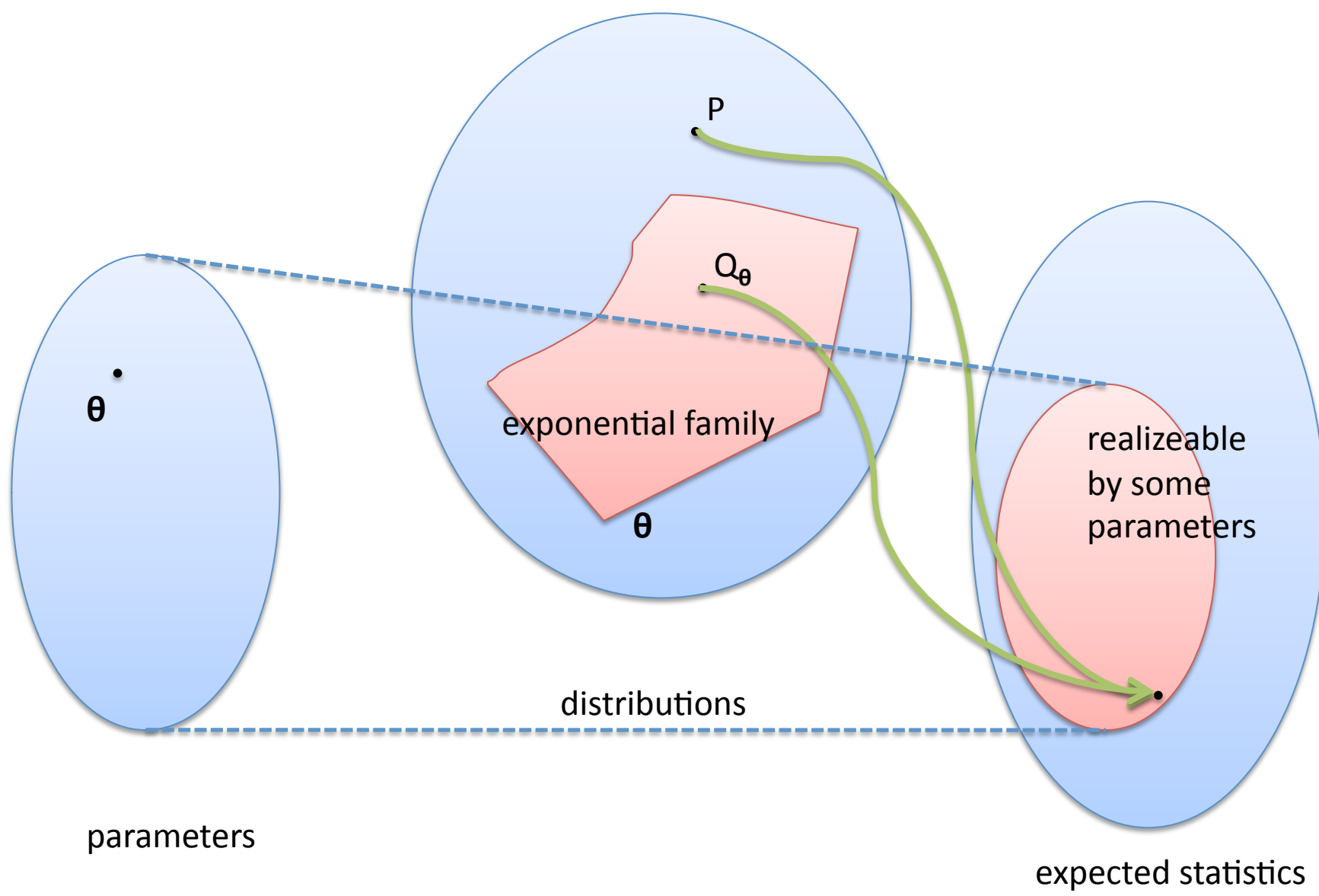


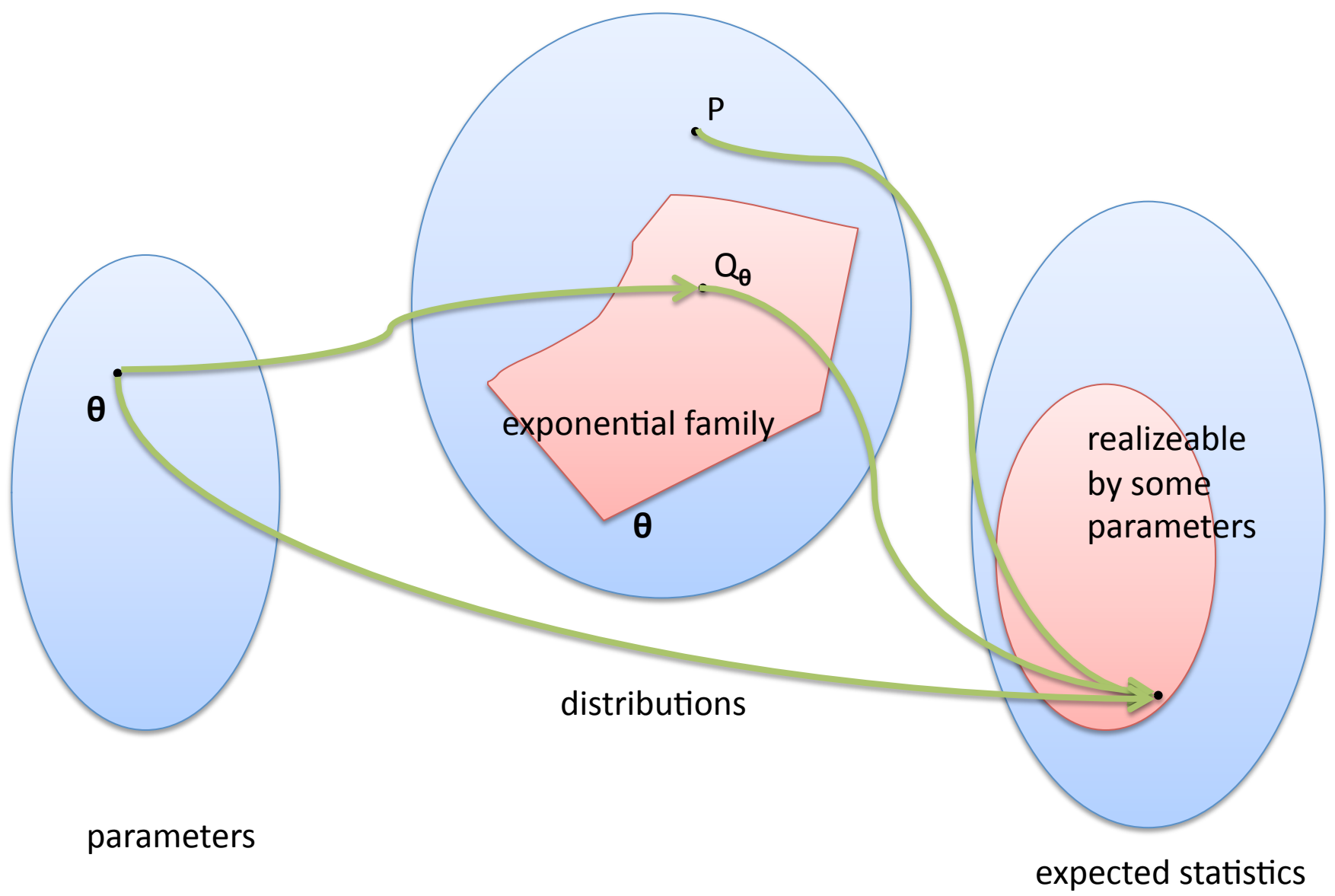
distributions



expected statistics







# Projections

- I-Projection is a bit harder:

$$D(Q_{\theta} \| P) = -H(Q_{\theta}(\mathbf{X})) - \mathbb{E}_{Q_{\theta}}[\ln P(\mathbf{X})]$$

- The exponential form of Q isn't very helpful about the second term. If P has some structure, this could help.

# Why Projections?

- M-projection is a foundation for learning.
- I-projection is used when answering difficult probabilistic queries.
- Note: I am not yet talking about algorithms for solving either!

# What You Now Know

- What is an exponential family?
- What is a linear exponential family?
- How these relate to Bayesian and Markov networks.
- Information theoretic quantities.