

# Probabilistic Graphical Models

Parameter learning in Bayesian networks

Siamak Ravanbakhsh

Fall 2019

# Learning objectives

- likelihood function and MLE
- role of the sufficient statistics
- MLE for parameter learning in directed models
  - why is it easy?
- conjugate priors and Bayesian parameter learning

# Likelihood function through an example

a thumbtack with unknown prob. of heads & tails

Bernoulli dist.  $p(x; \theta) = \theta^x (1 - \theta)^{(1-x)}$



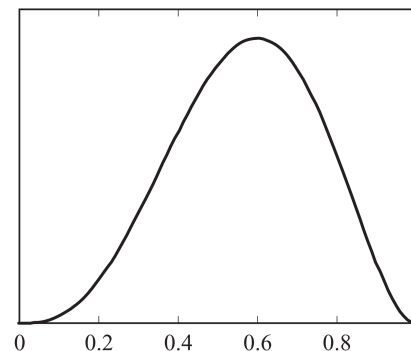
# Likelihood function through an example

a thumbtack with unknown prob. of heads & tails

Bernoulli dist.  $p(x; \theta) = \theta^x (1 - \theta)^{(1-x)}$

IID observations  $\mathcal{D} = \{1, 0, 0, 1, 1\}$

likelihood of  $\theta$  is  $L(\theta; \mathcal{D}) = \prod_{x \in \mathcal{D}} P(x; \theta) = \theta^3 (1 - \theta)^2$



likelihood function  $\theta$

not a pdf (it does not integrate to 1)

# Likelihood function through an example

a thumbtack with unknown prob. of heads & tails

Bernoulli dist.  $p(x; \theta) = \theta^x (1 - \theta)^{(1-x)}$

IID observations  $\mathcal{D} = \{1, 0, 0, 1, 1\}$

likelihood of  $\theta$  is  $L(\theta; \mathcal{D}) = \prod_{x \in \mathcal{D}} P(x; \theta) = \theta^3 (1 - \theta)^2$

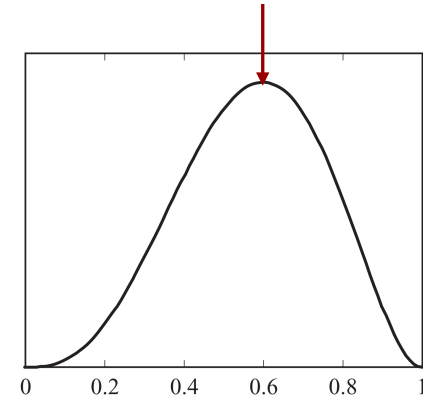
log-likelihood:  $\log L(\theta; \mathcal{D}) = 3 \log \theta + 2 \log(1 - \theta)$

maximizing the log-likelihood (M-projection of  $P_{\mathcal{D}}$ )

$$\frac{\partial}{\partial \theta} (3 \log \theta + 2 \log(1 - \theta)) = \frac{3}{\theta} - \frac{2}{1-\theta} = \frac{3-5\theta}{\theta(1-\theta)} = 0 \Rightarrow \hat{\theta} = \frac{3}{5}$$



**max-likelihood estimate (MLE)**



likelihood function  $\theta$   
not a pdf (it does not integrate to 1)

# Sufficient statistics through an example

IID observations  $\mathcal{D} = \{1, 0, 0, 1, 1\}$

likelihood of  $\theta$  is  $L(\theta, \mathcal{D}) = \prod_{x \in \mathcal{D}} P(x; \theta) = \theta^3(1 - \theta)^2$

heads  $\equiv 1$

tails  $\equiv 0$



all we needed to know about the data:

- number of heads and tails

given a distribution  $P(x; \theta)$

- its **sufficient statistics** is function  $\phi = [\phi_1, \dots, \phi_K]$  such that

$$\mathbb{E}_{\mathcal{D}}[\phi(x)] = \mathbb{E}_{\mathcal{D}'}[\phi(x')] \quad \Rightarrow \quad \frac{1}{|\mathcal{D}|} L(\theta, \mathcal{D}) = \frac{1}{|\mathcal{D}'|} L(\theta, \mathcal{D}') \quad \forall \mathcal{D}, \mathcal{D}', \theta$$

**sufficient statistics of the dataset** is all that matters about the data

# Revisiting exponential family

given a distribution  $P(x; \theta)$

- its **sufficient statistics** is function  $\phi = [\phi_1, \dots, \phi_K]$  such that

$$\mathbb{E}_{\mathcal{D}}[\phi(x)] = \mathbb{E}_{\mathcal{D}'}[\phi(x')] \quad \Rightarrow \quad \frac{1}{|\mathcal{D}|} L(\theta, \mathcal{D}) = \frac{1}{|\mathcal{D}'|} L(\theta, \mathcal{D}') \quad \forall \mathcal{D}, \mathcal{D}', \theta \quad L(\theta, \mathcal{D}) = \prod_{x \in \mathcal{D}} p(x; \theta)$$

the (linear) exponential family:  $p(x) \propto \exp(\langle \theta, \phi(x) \rangle)$

- **max-entropy distribution** subject to  $\mathbb{E}_p[\phi(x)] = \mu$

# Revisiting exponential family

given a distribution  $P(x; \theta)$

- its **sufficient statistics** is function  $\phi = [\phi_1, \dots, \phi_K]$  such that

$$\mathbb{E}_{\mathcal{D}}[\phi(x)] = \mathbb{E}_{\mathcal{D}'}[\phi(x')] \quad \Rightarrow \quad \frac{1}{|\mathcal{D}|} L(\theta, \mathcal{D}) = \frac{1}{|\mathcal{D}'|} L(\theta, \mathcal{D}') \quad \forall \mathcal{D}, \mathcal{D}', \theta \quad \boxed{L(\theta, \mathcal{D}) = \prod_{x \in \mathcal{D}} p(x; \theta)}$$

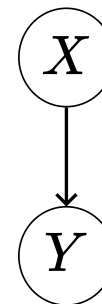
the (linear) exponential family:  $p(x) \propto \exp(\langle \theta, \phi(x) \rangle)$

- **max-entropy distribution** subject to  $\mathbb{E}_p[\phi(x)] = \mu$
- if  $\phi_1, \dots, \phi_k$  are linearly independent, then  $\theta \leftrightarrow \mu$



# MLE for Bayesian networks an example

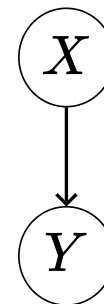
a simple network  $p(x, y; \theta) = p(x; \theta_X)p(y|x; \theta_{Y|X})$



# MLE for Bayesian networks an example

a simple network  $p(x, y; \theta) = p(x; \theta_X)p(y|x; \theta_{Y|X})$

likelihood 
$$L(\mathcal{D}; \theta) = \prod_{(x,y) \in \mathcal{D}} p(x; \theta_X)p(y|x; \theta_{Y|X})$$
$$= \underbrace{\left( \prod_{(x) \in \mathcal{D}} p(x; \theta_X) \right)}_{\text{likelihood of x}} \underbrace{\left( \prod_{(x,y) \in \mathcal{D}} p(y|x; \theta_{Y|X}) \right)}_{\text{cond. likelihood of y}}$$

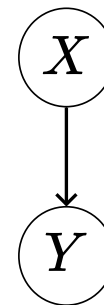


# MLE for Bayesian networks an example

a simple network  $p(x, y; \theta) = p(x; \theta_X)p(y|x; \theta_{Y|X})$

likelihood 
$$L(\mathcal{D}; \theta) = \prod_{(x,y) \in \mathcal{D}} p(x; \theta_X)p(y|x; \theta_{Y|X})$$

$$= \underbrace{\left( \prod_{(x) \in \mathcal{D}} p(x; \theta_X) \right)}_{\text{likelihood of x}} \underbrace{\left( \prod_{(x,y) \in \mathcal{D}} p(y|x; \theta_{Y|X}) \right)}_{\text{cond. likelihood of y}}$$

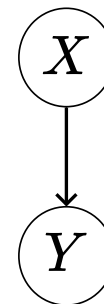


for discrete vars.

$$L(\mathcal{D}; \theta) = \left( \prod_{\ell \in \text{Val}(X)} \theta_{X,\ell}^{N(x=\ell)} \right) \left( \prod_{\ell, \ell' \in \text{Val}(X) \times \text{Val}(Y)} \theta_{Y|X,\ell,\ell'}^{N(x=\ell,y=\ell')} \right)$$

$\downarrow$   $p(X = \ell)$                        $\downarrow$   $p(X = \ell | Y = \ell')$

# MLE for Bayesian networks an example



a simple network  $p(x, y; \theta) = p(x; \theta_X)p(y|x; \theta_{Y|X})$

$$\begin{aligned}
 \text{likelihood } L(\mathcal{D}; \theta) &= \prod_{(x,y) \in \mathcal{D}} p(x; \theta_X)p(y|x; \theta_{Y|X}) \\
 &= \underbrace{\left( \prod_{(x) \in \mathcal{D}} p(x; \theta_X) \right)}_{\text{likelihood of x}} \underbrace{\left( \prod_{(x,y) \in \mathcal{D}} p(y|x; \theta_{Y|X}) \right)}_{\text{cond. likelihood of y}}
 \end{aligned}$$

for discrete vars.

$$\begin{aligned}
 L(\mathcal{D}; \theta) &= \left( \prod_{\ell \in \text{Val}(X)} \theta_{X,\ell}^{N(x=\ell)} \right) \left( \prod_{\ell, \ell' \in \text{Val}(X) \times \text{Val}(Y)} \theta_{Y|X,\ell,\ell'}^{N(x=\ell,y=\ell')} \right) \\
 &\quad \downarrow \qquad \qquad \qquad \downarrow \\
 &\quad p(X = \ell) \qquad \qquad \qquad p(X = \ell | Y = \ell')
 \end{aligned}$$

number of times  $x = \ell$  in the dataset      number of times  $x = \ell, y = \ell'$  in the dataset

MLE : maximize **local likelihood** terms individually

$$\theta_{X,\ell} = \frac{N(x=\ell)}{|\mathcal{D}|} \quad \theta_{Y|X,\ell,\ell'} = \frac{N(x=\ell,y=\ell')}{|\mathcal{D}|}$$

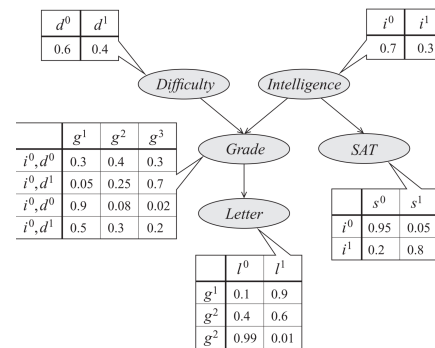
# MLE for Bayesian networks **general case**

Bayes-net  $p(x; \theta) = \prod_i p(x_i \mid Pa_{x_i}; \theta_{X_i \mid Pa_{X_i}})$

likelihood  $L(\mathcal{D}; \theta) = \prod_{x \in \mathcal{D}} \prod_i p(x_i \mid Pa_{x_i}; \theta_i \mid Pa_i)$

$$= \prod_i \prod_{(x_i, Pa_{x_i}) \in \mathcal{D}} p(x_i \mid Pa_{x_i}; \theta_i \mid Pa_i)$$

local likelihood terms



# MLE for Bayesian networks general case

Bayes-net  $p(\mathbf{x}; \theta) = \prod_i p(x_i \mid Pa_{x_i}; \theta_{X_i|Pa_{X_i}})$

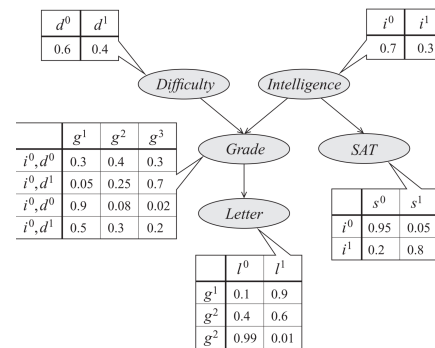
likelihood  $L(\mathcal{D}; \theta) = \prod_{\mathbf{x} \in \mathcal{D}} \prod_i p(x_i \mid Pa_{x_i}; \theta_i|Pa_i)$

$$= \prod_i \prod_{(x_i, Pa_{x_i}) \in \mathcal{D}} p(x_i \mid Pa_{x_i}; \theta_i|Pa_i)$$

local likelihood terms

maximizing the conditional likelihood for each node

- similar to solving individual prediction problems



# MLE for Bayesian networks general case

Bayes-net  $p(\mathbf{x}; \theta) = \prod_i p(x_i \mid Pa_{x_i}; \theta_{X_i|Pa_{X_i}})$

likelihood  $L(\mathcal{D}; \theta) = \prod_{\mathbf{x} \in \mathcal{D}} \prod_i p(x_i \mid Pa_{x_i}; \theta_i|Pa_i)$

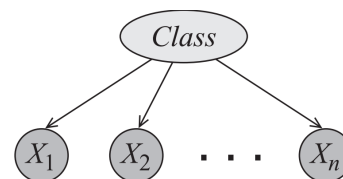
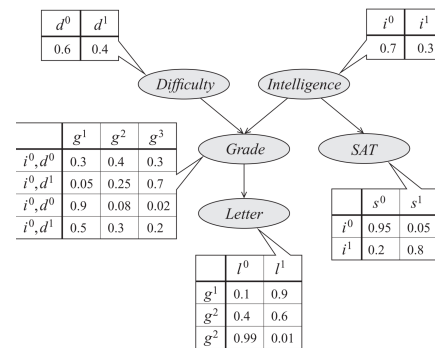
$$= \prod_i \prod_{(x_i, Pa_{x_i}) \in \mathcal{D}} p(x_i \mid Pa_{x_i}; \theta_i|Pa_i)$$

local likelihood terms

maximizing the conditional likelihood for each node

- similar to solving individual prediction problems

**Example** how to learn a naive Bayes?



# Bayesian parameter estimation

max-likelihood is the same  $\hat{\theta} = \frac{1}{3}$  for

- case 1.  $N(x = 1) = 1, N(x = 0) = 2$
- case 2.  $N(x = 1) = 100, N(x = 0) = 200$

heads  $\equiv 1$

tails  $\equiv 0$



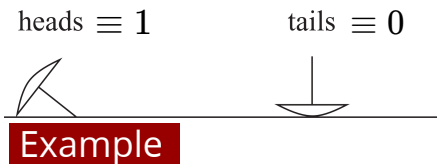
Example



# Bayesian parameter estimation

max-likelihood is the same  $\hat{\theta} = \frac{1}{3}$  for

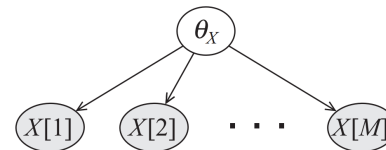
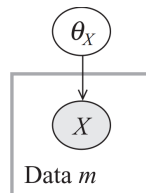
$$\begin{array}{l} \text{case 1.} \\ \text{case 2.} \end{array} \left| \begin{array}{l} N(x = 1) = 1, N(x = 0) = 2 \\ N(x = 1) = 100, N(x = 0) = 200 \end{array} \right.$$



need to model our uncertainty

Bayesian approach:

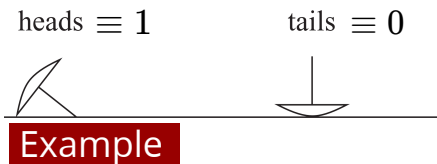
- assume a prior  $p(\theta)$
- estimate the posterior



# Bayesian parameter estimation

max-likelihood is the same  $\hat{\theta} = \frac{1}{3}$  for

$$\begin{array}{l} \text{case 1.} \\ \text{case 2.} \end{array} \left| \begin{array}{l} N(x=1) = 1, N(x=0) = 2 \\ N(x=1) = 100, N(x=0) = 200 \end{array} \right.$$



need to model our uncertainty

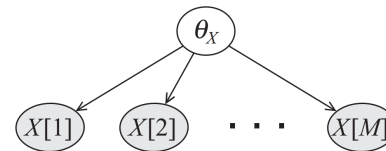
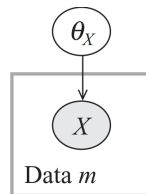
Bayesian approach:

- assume a prior  $p(\theta)$
- estimate the posterior

$$p(\theta | \mathcal{D}) = \frac{p(\theta)p(\mathcal{D}|\theta)}{p(\mathcal{D})} \propto p(\theta)p(\mathcal{D} | \theta)$$

**posterior**      **marginal likelihood**       $\prod_{x \in \mathcal{D}} p(x|\theta)$

**prior**      **likelihood**



# Bayesian parameter estimation

assuming a **uniform prior**  $p(\theta) = \begin{cases} 1 & 0 \leq \theta \leq 1 \\ 0 & o.w. \end{cases}$

posterior  $p(\theta | \mathcal{D}) \propto p(\theta)p(\mathcal{D} | \theta) \propto p(\mathcal{D} | \theta)$

heads  $\equiv 1$

tails  $\equiv 0$



# Bayesian parameter estimation

assuming a **uniform prior**  $p(\theta) = \begin{cases} 1 & 0 \leq \theta \leq 1 \\ 0 & \text{o.w.} \end{cases}$

heads  $\equiv 1$

tails  $\equiv 0$

posterior  $p(\theta | \mathcal{D}) \propto p(\theta)p(\mathcal{D} | \theta) \propto p(\mathcal{D} | \theta)$



**posterior predictive:** predicting heads/tails using the posterior

rather than a single MLE value

$$p(x | \mathcal{D}) = \int_0^1 p(\theta | \mathcal{D}) p(x | \theta) d\theta$$
$$\propto \theta^{N(1)} (1 - \theta)^{N(0)} \theta^x (1 - \theta)^{(1-x)}$$

# Bayesian parameter estimation

assuming a **uniform prior**  $p(\theta) = \begin{cases} 1 & 0 \leq \theta \leq 1 \\ 0 & o.w. \end{cases}$

heads  $\equiv 1$

tails  $\equiv 0$

posterior  $p(\theta | \mathcal{D}) \propto p(\theta)p(\mathcal{D} | \theta) \propto p(\mathcal{D} | \theta)$



**posterior predictive:** predicting heads/tails using the posterior

rather than a single MLE value

$$p(x | \mathcal{D}) = \int_0^1 p(\theta | \mathcal{D}) p(x | \theta) d\theta \\ \propto \theta^{N(1)} (1 - \theta)^{N(0)} \theta^x (1 - \theta)^{(1-x)}$$

if we do the integration above:  $p(x = 1 | \mathcal{D}) = \frac{N(1)+1}{N(0)+N(1)+2}$   
(and normalize)

Laplace correction

# Bayesian parameter estimation

assuming a **uniform prior**  $p(\theta) = \begin{cases} 1 & 0 \leq \theta \leq 1 \\ 0 & o.w. \end{cases}$

heads  $\equiv 1$

tails  $\equiv 0$

posterior  $p(\theta | \mathcal{D}) \propto p(\theta)p(\mathcal{D} | \theta) \propto p(\mathcal{D} | \theta)$



**posterior predictive:** predicting heads/tails using the posterior

rather than a single MLE value

$$p(x | \mathcal{D}) = \int_0^1 p(\theta | \mathcal{D}) p(x | \theta) d\theta \\ \propto \theta^{N(1)} (1 - \theta)^{N(0)} \theta^x (1 - \theta)^{(1-x)}$$

if we do the integration above:  $p(x = 1 | \mathcal{D}) = \frac{N(1)+1}{N(0)+N(1)+2}$   
(and normalize)

Laplace correction

compare with prediction using MLE  $p(x = 1 | \mathcal{D}) = \frac{N(1)}{N(0)+N(1)}$

# Conjugate priors

heads  $\equiv 1$

tails  $\equiv 0$



how about non-uniform priors? E.g., more likely to see heads

need an efficient way to get the posterior  $p(\theta | \mathcal{D}) \propto p(\theta)p(\mathcal{D} | \theta)$

# Conjugate priors



how about non-uniform priors? E.g., more likely to see heads

need an efficient way to get the posterior  $p(\theta | \mathcal{D}) \propto p(\theta)p(\mathcal{D} | \theta)$

**ideally** the prior  $p(\theta)$  & the posterior  $p(\theta|\mathcal{D})$  should have the same form  
 $p(\theta)$  is a **conjugate prior** to the likelihood  $p(\mathcal{D}|\theta)$



# Conjugate priors

heads  $\equiv 1$

tails  $\equiv 0$



how about non-uniform priors? E.g., more likely to see heads

need an efficient way to get the posterior  $p(\theta | \mathcal{D}) \propto p(\theta)p(\mathcal{D} | \theta)$

**ideally** the prior  $p(\theta)$  & the posterior  $p(\theta|\mathcal{D})$  should have the same form  
 $p(\theta)$  is a **conjugate prior** to the likelihood  $p(\mathcal{D}|\theta)$

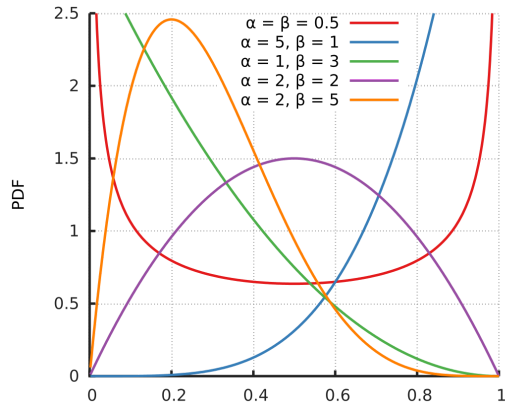
conjugate prior to the **Bernoulli likelihood** is the **Beta distribution**

$$p(\mathcal{D}|\theta) \propto \theta^{N(1)}(1 - \theta)^{N(0)}$$

$$p(\theta; \alpha, \beta) = \gamma \theta^{\alpha-1} (1 - \theta)^{\beta-1}$$

$$\gamma = \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)}$$

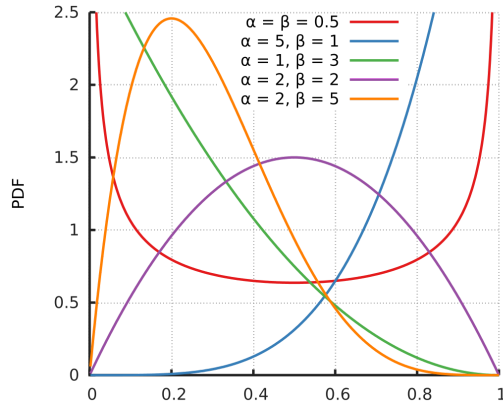
# Conjugate priors: **Beta-Bernoulli**



conjugate prior to the **Bernoulli likelihood** is the **Beta distribution**

$$p(\theta; \alpha, \beta) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \theta^{\alpha-1} (1 - \theta)^{\beta-1}$$

# Conjugate priors: **Beta-Bernoulli**

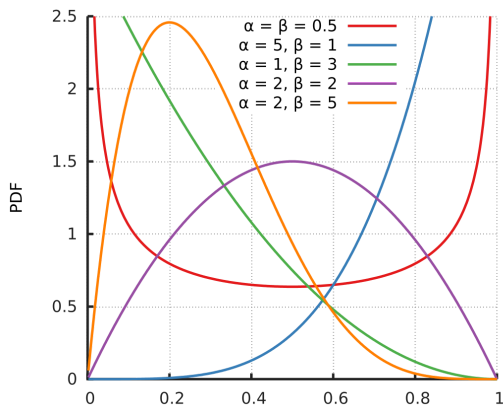


conjugate prior to the Bernoulli likelihood is the **Beta distribution**

$$p(\theta; \alpha, \beta) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \theta^{\alpha-1} (1 - \theta)^{\beta-1}$$

↓  
extension of factorial function  $\Gamma(n + 1) = n!$

# Conjugate priors: **Beta-Bernoulli**



conjugate prior to the Bernoulli likelihood is the **Beta distribution**

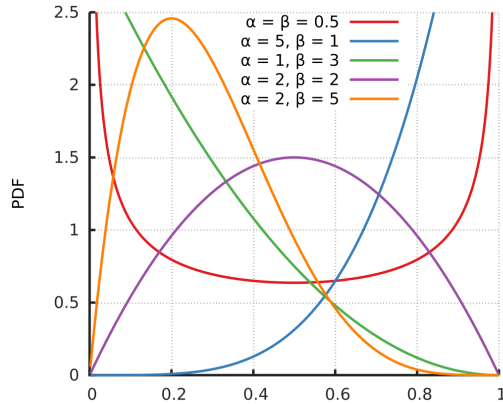
$$p(\theta; \alpha, \beta) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \theta^{\alpha-1} (1 - \theta)^{\beta-1}$$

extension of factorial function  $\Gamma(n + 1) = n!$

**hyper-parameters:** can be interpreted as # **imaginary** heads & tails

$$\text{prior predictive: } p(x = 1 \mid \mathcal{D} = \emptyset) = \int_{\theta} p(x = 1 \mid \theta) p(\theta; \alpha, \beta) d\theta = \frac{\alpha}{\alpha + \beta}$$

# Conjugate priors: Beta-Bernoulli



conjugate prior to the Bernoulli likelihood is the Beta distribution

$$p(\theta; \alpha, \beta) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \theta^{\alpha-1} (1 - \theta)^{\beta-1}$$

extension of factorial function  $\Gamma(n + 1) = n!$

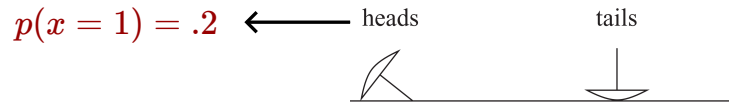
hyper-parameters: can be interpreted as # imaginary heads & tails

$$\text{prior predictive: } p(x = 1 | \mathcal{D} = \emptyset) = \int_{\theta} p(x = 1 | \theta) p(\theta; \alpha, \beta) d\theta = \frac{\alpha}{\alpha + \beta}$$

$$\text{posterior: } p(\theta | \mathcal{D}) \propto p(\theta) P(\mathcal{D} | \theta) \propto \theta^{\alpha-1} (1 - \theta)^{\beta-1} \theta^{N(1)} (1 - \theta)^{N(0)} = \theta^{\alpha-1+N(1)} (1 - \theta)^{\beta-1+N(0)}$$

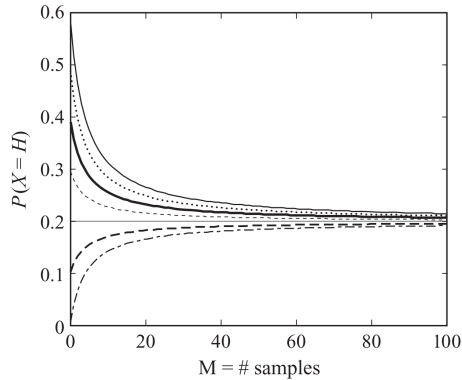
if the prior is  $p(\theta; \alpha, \beta)$ , the posterior is  $p(\theta; \alpha + N(1), \beta + N(0))$

# Beta-Bernoulli: Example

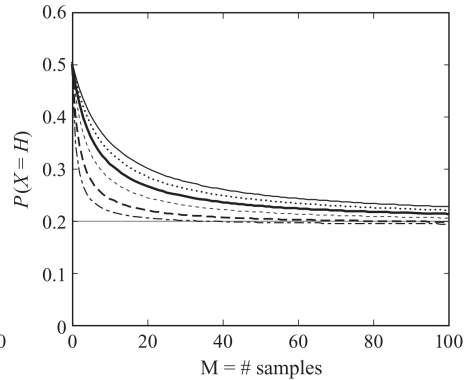


posterior for different priors and sample sizes

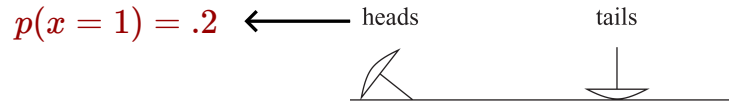
different prior means  $\frac{\alpha}{\alpha + \beta}$



different prior strength  $\alpha + \beta$

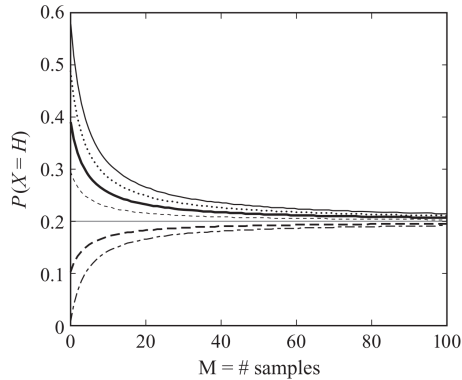


# Beta-Bernoulli: Example

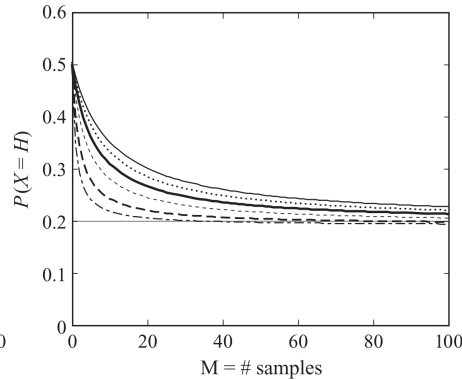


posterior for different priors and sample sizes

different prior means  $\frac{\alpha}{\alpha + \beta}$

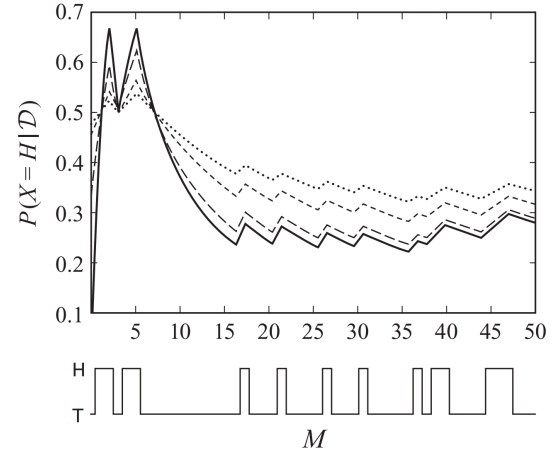


different prior strength  $\alpha + \beta$



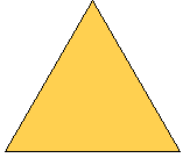
posterior predictive for online setting

- MLE
- - -  $\alpha = \beta = 1$
- · -  $\alpha = \beta = 5$

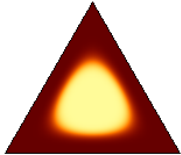


# Conjugate priors: Dirichlet-categorical

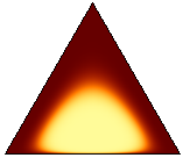
$\alpha = (1,1,1)$



$(7,7,7)$



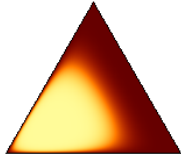
$(5,5,2)$



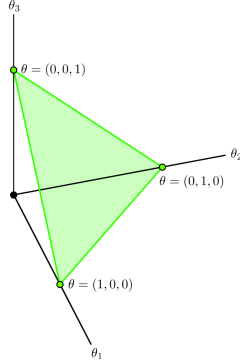
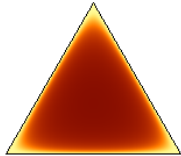
$(3,3,3)$



$(5,2,2)$



$(0.2,0.2,0.2)$



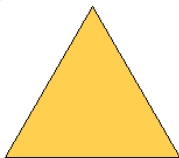
Bernoulli  $\Rightarrow$  Beta  
Categorical  $\Rightarrow$  Dirichlet

$$p(\theta; \alpha) = \frac{\Gamma(\sum_d \alpha_d)}{\prod_d \Gamma(\alpha_d)} \prod_d \theta_d^{\alpha_d - 1}$$



# Conjugate priors: Dirichlet-categorical

$$\alpha = (1, 1, 1)$$



$$(7, 7, 7)$$



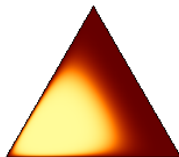
$$(5, 5, 2)$$



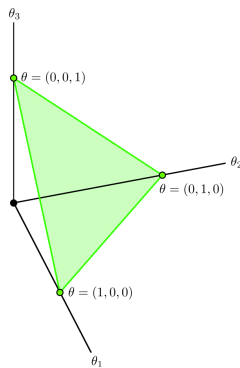
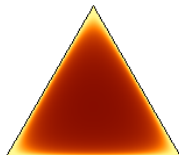
$$(3, 3, 3)$$



$$(5, 2, 2)$$



$$(0.2, 0.2, 0.2)$$



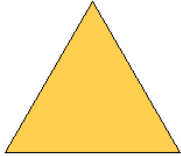
Bernoulli  $\Rightarrow$  Beta  
Categorical  $\Rightarrow$  Dirichlet

$$p(\theta; \alpha) = \frac{\Gamma(\sum_d \alpha_d)}{\prod_d \Gamma(\alpha_d)} \prod_d \theta_d^{\alpha_d - 1}$$

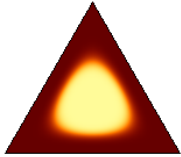
$\alpha \in (\mathbb{R}^+)^D$  pseudo-counts for different categories

# Conjugate priors: Dirichlet-categorical

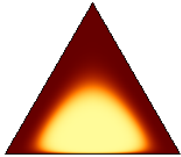
$$\alpha = (1, 1, 1)$$



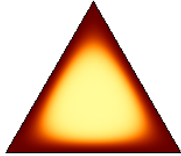
$$(7, 7, 7)$$



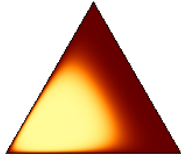
$$(5, 5, 2)$$



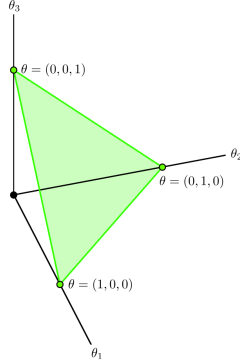
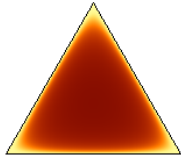
$$(3, 3, 3)$$



$$(5, 2, 2)$$



$$(0.2, 0.2, 0.2)$$



Bernoulli  $\Rightarrow$  Beta  
Categorical  $\Rightarrow$  Dirichlet

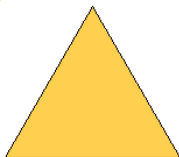
$$p(\theta; \alpha) = \frac{\Gamma(\sum_d \alpha_d)}{\prod_d \Gamma(\alpha_d)} \prod_d \theta_d^{\alpha_d - 1}$$

$\alpha \in (\mathbb{R}^+)^D$  pseudo-counts for different categories

prior:  $p(\theta; \alpha)$

# Conjugate priors: Dirichlet-categorical

$$\alpha = (1, 1, 1)$$



$$(7, 7, 7)$$



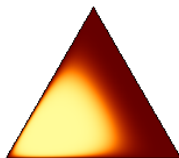
$$(5, 5, 2)$$



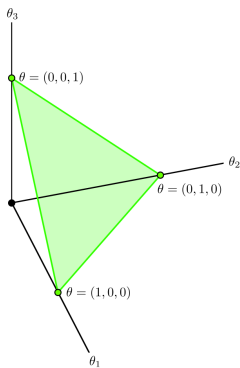
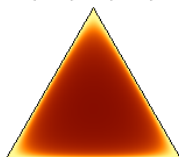
$$(3, 3, 3)$$



$$(5, 2, 2)$$



$$(0.2, 0.2, 0.2)$$



Bernoulli  $\Rightarrow$  Beta  
Categorical  $\Rightarrow$  Dirichlet

$$p(\theta; \alpha) = \frac{\Gamma(\sum_d \alpha_d)}{\prod_d \Gamma(\alpha_d)} \prod_d \theta_d^{\alpha_d - 1}$$

$\alpha \in (\mathbb{R}^+)^D$  pseudo-counts for different categories

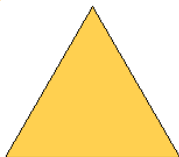
prior:  $p(\theta; \alpha)$

likelihood:  $p(\mathcal{D} | \theta) \propto \prod_{x \in \mathcal{D}} \prod_d \theta_d^{\mathbb{I}(x=d)} = \prod_d \theta_d^{N(d)}$

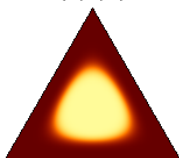
posterior:  $p(\theta | \mathcal{D}) \propto p(\theta) p(\mathcal{D} | \theta)$

# Conjugate priors: Dirichlet-categorical

$$\alpha = (1, 1, 1)$$



$$(7, 7, 7)$$



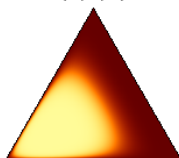
$$(5, 5, 2)$$



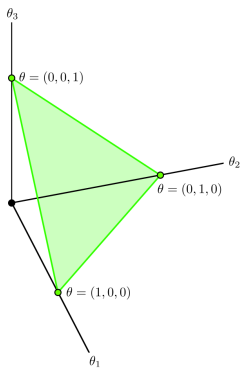
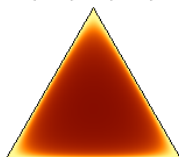
$$(3, 3, 3)$$



$$(5, 2, 2)$$



$$(0.2, 0.2, 0.2)$$



Bernoulli  $\Rightarrow$  Beta  
Categorical  $\Rightarrow$  Dirichlet

$$p(\theta; \alpha) = \frac{\Gamma(\sum_d \alpha_d)}{\prod_d \Gamma(\alpha_d)} \prod_d \theta_d^{\alpha_d - 1}$$

$\alpha \in (\mathbb{R}^+)^D$  pseudo-counts for different categories

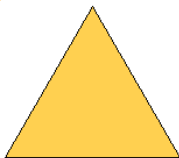
prior:  $p(\theta; \alpha)$

$$\text{likelihood: } p(\mathcal{D} | \theta) \propto \prod_{x \in \mathcal{D}} \prod_d \theta_d^{\mathbb{I}(x=d)} = \prod_d \theta_d^{N(d)}$$

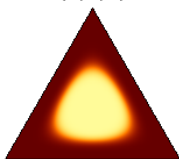
$$\text{posterior: } p(\theta | \mathcal{D}) \propto p(\theta) p(\mathcal{D} | \theta) \propto \prod_d \theta_d^{N(d)} \theta_d^{\alpha_d - 1} = \prod_d \theta_d^{\alpha_d + N(d) - 1}$$

# Conjugate priors: Dirichlet-categorical

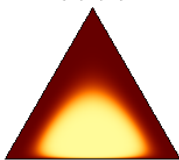
$$\alpha = (1, 1, 1)$$



$$(7, 7, 7)$$



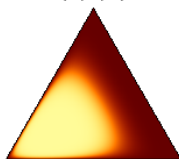
$$(5, 5, 2)$$



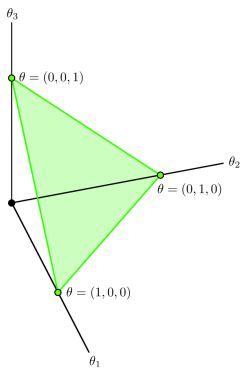
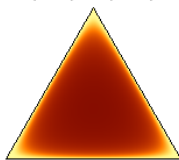
$$(3, 3, 3)$$



$$(5, 2, 2)$$



$$(0.2, 0.2, 0.2)$$



Bernoulli  $\Rightarrow$  Beta  
Categorical  $\Rightarrow$  Dirichlet

$$p(\theta; \alpha) = \frac{\Gamma(\sum_d \alpha_d)}{\prod_d \Gamma(\alpha_d)} \prod_d \theta_d^{\alpha_d - 1}$$

$\alpha \in (\mathbb{R}^+)^D$  pseudo-counts for different categories

prior:  $p(\theta; \alpha)$

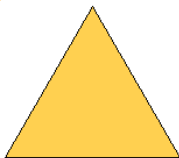
$$\text{likelihood: } p(\mathcal{D} | \theta) \propto \prod_{x \in \mathcal{D}} \prod_d \theta_d^{\mathbb{I}(x=d)} = \prod_d \theta_d^{N(d)}$$

$$\text{posterior: } p(\theta | \mathcal{D}) \propto p(\theta) p(\mathcal{D} | \theta) \propto \prod_d \theta_d^{N(d)} \theta_d^{\alpha_d - 1} = \prod_d \theta_d^{\alpha_d + N(d) - 1}$$

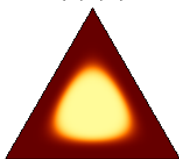
$$\text{posterior predictive: } p(x = \bar{x} | \mathcal{D}) = \int_{\theta} p(\theta | \mathcal{D}) p(x = \bar{x} | \theta) d\theta$$

# Conjugate priors: Dirichlet-categorical

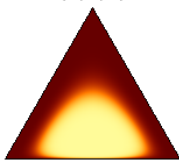
$$\alpha = (1, 1, 1)$$



$$(7, 7, 7)$$



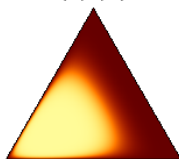
$$(5, 5, 2)$$



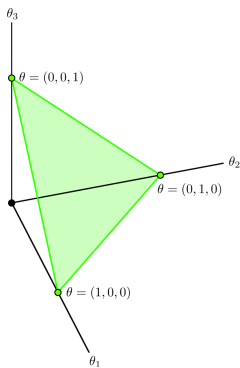
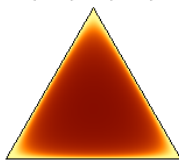
$$(3, 3, 3)$$



$$(5, 2, 2)$$



$$(0.2, 0.2, 0.2)$$



Bernoulli  $\Rightarrow$  Beta  
Categorical  $\Rightarrow$  Dirichlet

$$p(\theta; \alpha) = \frac{\Gamma(\sum_d \alpha_d)}{\prod_d \Gamma(\alpha_d)} \prod_d \theta_d^{\alpha_d - 1}$$

$\alpha \in (\mathbb{R}^+)^D$  pseudo-counts for different categories

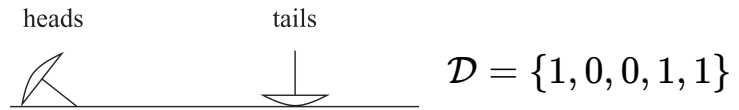
prior:  $p(\theta; \alpha)$

$$\text{likelihood: } p(\mathcal{D} | \theta) \propto \prod_{x \in \mathcal{D}} \prod_d \theta_d^{\mathbb{I}(x=d)} = \prod_d \theta_d^{N(d)}$$

$$\text{posterior: } p(\theta | \mathcal{D}) \propto p(\theta) p(\mathcal{D} | \theta) \propto \prod_d \theta_d^{N(d)} \theta_d^{\alpha_d - 1} = \prod_d \theta_d^{\alpha_d + N(d) - 1}$$

$$\text{posterior predictive: } p(x = \bar{x} | \mathcal{D}) = \int_{\theta} p(\theta | \mathcal{D}) p(x = \bar{x} | \theta) d\theta = \frac{\alpha_{\bar{x}} + N(\bar{x})}{N + \sum_d \alpha_d}$$

# Marginal likelihood vs. maximum likelihood



# Marginal likelihood vs. maximum likelihood

heads                      tails  
  $\mathcal{D} = \{1, 0, 0, 1, 1\}$

maximum likelihood value:  $P(\mathcal{D}|\hat{\theta}) = \left(\frac{3}{5}\right)^3 \left(\frac{2}{5}\right)^2 \approx .035$



# Marginal likelihood vs. maximum likelihood

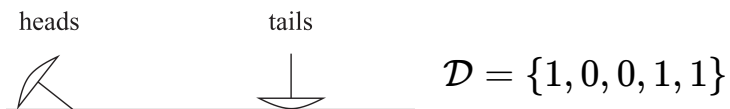


maximum likelihood value:  $P(\mathcal{D}|\hat{\theta}) = \left(\frac{3}{5}\right)^3 \left(\frac{2}{5}\right)^2 \approx .035$

marginal likelihood value:

**1**  $P(\mathcal{D}) = \int_{\theta \in [0,1]} P(\theta)P(\mathcal{D}|\theta)d\theta$

# Marginal likelihood vs. maximum likelihood



maximum likelihood value:  $P(\mathcal{D}|\hat{\theta}) = \left(\frac{3}{5}\right)^3 \left(\frac{2}{5}\right)^2 \approx .035$

marginal likelihood value:

1  $P(\mathcal{D}) = \int_{\theta \in [0,1]} P(\theta)P(\mathcal{D}|\theta)d\theta$

2 chain rule:  $P(\mathcal{D}) = \prod_{m=1}^M P(x^{(m)}|x^{(1)}, \dots, x^{(m-1)})$

# Marginal likelihood vs. maximum likelihood



maximum likelihood value:  $P(\mathcal{D}|\hat{\theta}) = \left(\frac{3}{5}\right)^3 \left(\frac{2}{5}\right)^2 \approx .035$

marginal likelihood value:

1  $P(\mathcal{D}) = \int_{\theta \in [0,1]} P(\theta) P(\mathcal{D}|\theta) d\theta$

2 chain rule:  $P(\mathcal{D}) = \prod_{m=1}^M P(x^{(m)} | x^{(1)}, \dots, x^{(m-1)})$

$$= \frac{\alpha_1}{\alpha} \cdot \frac{\alpha_0}{\alpha+1} \cdot \frac{\alpha_0+1}{\alpha+2} \cdot \frac{\alpha_1+1}{\alpha+3} \cdot \frac{\alpha_2+1}{\alpha+4}$$

# Marginal likelihood vs. maximum likelihood



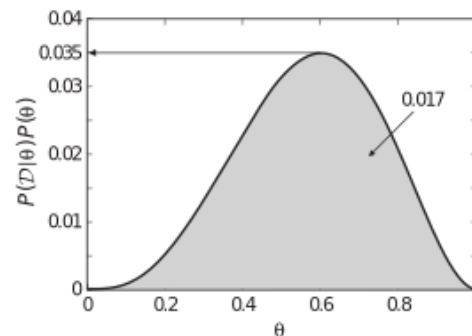
maximum likelihood value:  $P(\mathcal{D}|\hat{\theta}) = \left(\frac{3}{5}\right)^3 \left(\frac{2}{5}\right)^2 \approx .035$

marginal likelihood value:

1  $P(\mathcal{D}) = \int_{\theta \in [0,1]} P(\theta)P(\mathcal{D}|\theta)d\theta$

2 chain rule:  $P(\mathcal{D}) = \prod_{m=1}^M P(x^{(m)}|x^{(1)}, \dots, x^{(m-1)})$

$$= \frac{\alpha_1}{\alpha} \cdot \frac{\alpha_0}{\alpha+1} \cdot \frac{\alpha_0+1}{\alpha+2} \cdot \frac{\alpha_1+1}{\alpha+3} \cdot \frac{\alpha_2+1}{\alpha+4} = \frac{\Gamma(\alpha)}{\Gamma(\alpha+5)} \cdot \frac{\Gamma(\alpha_1+3)}{\Gamma(\alpha_1)} \cdot \frac{\Gamma(\alpha_0+2)}{\Gamma(\alpha_0)} \approx .017$$



using  $\Gamma(x+1) = x\Gamma(x)$

# Marginal likelihood vs. maximum likelihood



maximum likelihood value:  $P(\mathcal{D}|\hat{\theta}) = \left(\frac{3}{5}\right)^3 \left(\frac{2}{5}\right)^2 \approx .035$

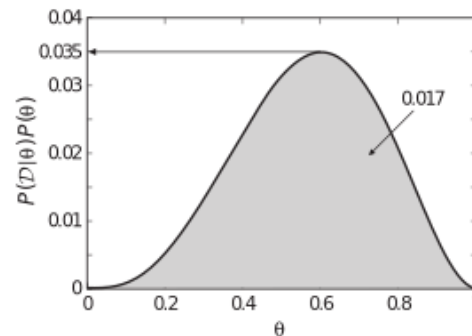
marginal likelihood value:

1  $P(\mathcal{D}) = \int_{\theta \in [0,1]} P(\theta)P(\mathcal{D}|\theta)d\theta$

2 chain rule:  $P(\mathcal{D}) = \prod_{m=1}^M P(x^{(m)}|x^{(1)}, \dots, x^{(m-1)})$

$$= \frac{\alpha_1}{\alpha} \cdot \frac{\alpha_0}{\alpha+1} \cdot \frac{\alpha_0+1}{\alpha+2} \cdot \frac{\alpha_1+1}{\alpha+3} \cdot \frac{\alpha_2+1}{\alpha+4} = \frac{\Gamma(\alpha)}{\Gamma(\alpha+5)} \cdot \frac{\Gamma(\alpha_1+3)}{\Gamma(\alpha_1)} \cdot \frac{\Gamma(\alpha_0+2)}{\Gamma(\alpha_0)} \approx .017$$

using  $\Gamma(x+1) = x\Gamma(x)$



marginal likelihood for Dirichlet  $P(\mathcal{D}) = \frac{\Gamma(\alpha)}{\Gamma(\alpha+|\mathcal{D}|)} \prod_i \frac{\Gamma(\alpha_i+|\mathcal{D}|p_{\mathcal{D}}(i))}{\Gamma(\alpha_i)}$

## Conjugate priors: **exponential family**

for the likelihood function:  $p(x | \theta) = \exp(\langle \phi(x), \theta \rangle - A(\theta))$

## Conjugate priors: exponential family

for the likelihood function:  $p(x | \theta) = \exp(\langle \phi(x), \theta \rangle - A(\theta))$

suppose we observe N instances  $p(\mathcal{D} | \theta) = \exp(\langle \sum_{x \in \mathcal{D}} \phi(x), \theta \rangle - NA(\theta))$

## Conjugate priors: exponential family

for the likelihood function:  $p(x | \theta) = \exp(\langle \phi(x), \theta \rangle - A(\theta))$

suppose we observe N instances  $p(\mathcal{D} | \theta) = \exp(\langle \sum_{x \in \mathcal{D}} \phi(x), \theta \rangle - NA(\theta))$

conjugate prior  $p(\theta; \eta, \nu) = \exp(\langle \nu \eta, \theta \rangle - \nu A(\theta))$



# Conjugate priors: exponential family

for the likelihood function:  $p(x | \theta) = \exp(\langle \phi(x), \theta \rangle - A(\theta))$

suppose we observe  $N$  instances  $p(\mathcal{D} | \theta) = \exp(\langle \sum_{x \in \mathcal{D}} \phi(x), \theta \rangle - NA(\theta))$

conjugate prior  $p(\theta; \eta, \nu) = \exp(\langle \nu \eta, \theta \rangle - \nu A(\theta))$

imaginary expected sufficient statistics

imaginary counts

posterior:  $p(\theta | \mathcal{D}; \eta, \nu) = \exp(\langle \nu \eta + \sum_{x \in \mathcal{D}} \phi(x), \theta \rangle - (\nu + N)A(\theta))$

# Bayesian learning for Bayes-nets

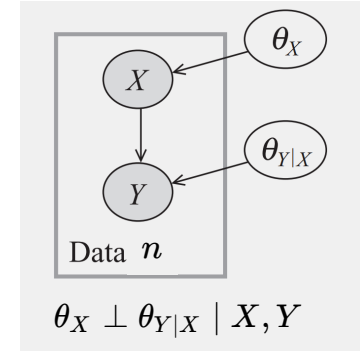
assumption

- global parameter independence: prior decomposes  $p(\theta) = \prod_i p(\theta_{X_i|Pa_{X_i}})$

conclusion

- posterior also decomposes  $p(\theta | \mathcal{D}) = \prod_i p(\theta_{X_i|Pa_{X_i}} | \mathcal{D})$

example



# Bayesian learning for Bayes-nets

assumption

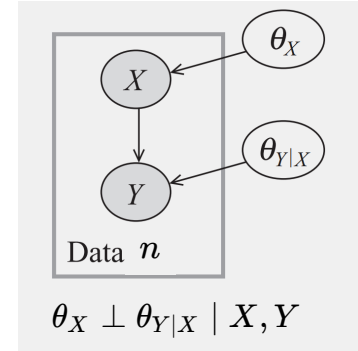
- global parameter independence: prior decomposes  $p(\theta) = \prod_i p(\theta_{X_i|Pa_{X_i}})$

conclusion

- posterior also decomposes  $p(\theta | \mathcal{D}) = \prod_i p(\theta_{X_i|Pa_{X_i}} | \mathcal{D})$

$$p(\theta | \mathcal{D}) = \underbrace{\prod_i p(\theta_{X_i|Pa_{X_i}})}_{\text{prior}} \underbrace{\prod_{x \in \mathcal{D}} \prod_i p(x_i | Pa_{x_i}; \theta_{X_i|Pa_{X_i}})}_{\text{likelihood}}$$

example



# Bayesian learning for Bayes-nets

assumption

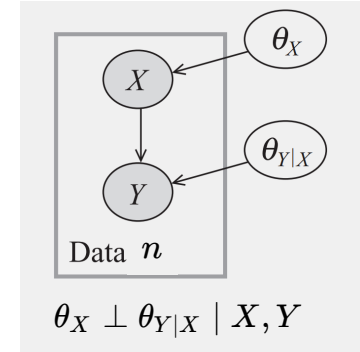
- global parameter independence: prior decomposes  $p(\theta) = \prod_i p(\theta_{X_i|Pa_{X_i}})$

conclusion

- posterior also decomposes  $p(\theta | \mathcal{D}) = \prod_i p(\theta_{X_i|Pa_{X_i}} | \mathcal{D})$

$$p(\theta | \mathcal{D}) = \underbrace{\prod_i p(\theta_{X_i|Pa_{X_i}})}_{\text{prior}} \underbrace{\prod_{x \in \mathcal{D}} \prod_i p(x_i | Pa_{x_i}; \theta_{X_i|Pa_{X_i}})}_{\text{likelihood}} = \underbrace{\prod_i p(\theta_{X_i|Pa_{X_i}} | \mathcal{D})}_{p(\theta_{X_i|Pa_{X_i}} | \mathcal{D})}$$

example



# Bayesian learning for Bayes-nets

assumption

- global parameter independence: prior decomposes  $p(\theta) = \prod_i p(\theta_{X_i|Pa_{X_i}})$

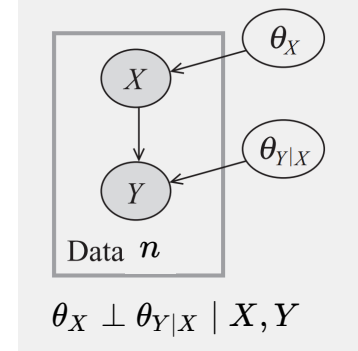
conclusion

- posterior also decomposes  $p(\theta | \mathcal{D}) = \prod_i p(\theta_{X_i|Pa_{X_i}} | \mathcal{D})$

$$p(\theta | \mathcal{D}) = \underbrace{\prod_i p(\theta_{X_i|Pa_{X_i}})}_{\text{prior}} \underbrace{\prod_{x \in \mathcal{D}} \prod_i p(x_i | Pa_{x_i}; \theta_{X_i|Pa_{X_i}})}_{\text{likelihood}} = \underbrace{\prod_i p(\theta_{X_i|Pa_{X_i}})}_{p(\theta_{X_i|Pa_{X_i}} | \mathcal{D})} \underbrace{\prod_{x \in \mathcal{D}} p(x_i | Pa_{x_i}; \theta_{X_i|Pa_{X_i}})}_{\text{individual posteriors}}$$

- we can apply Bayesian learning to individual conditional distributions

example



# Bayesian learning for Bayes-nets

assumption

- global parameter independence: prior decomposes  $p(\theta) = \prod_i p(\theta_{X_i|Pa_{X_i}})$

conclusion

- posterior also decomposes  $p(\theta | \mathcal{D}) = \prod_i p(\theta_{X_i|Pa_{X_i}} | \mathcal{D})$

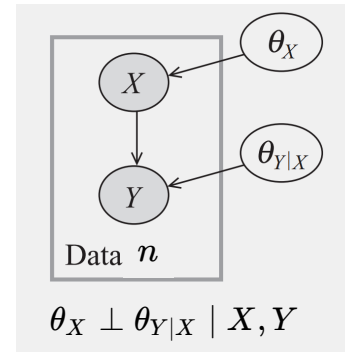
$$p(\theta | \mathcal{D}) = \underbrace{\prod_i p(\theta_{X_i|Pa_{X_i}})}_{\text{prior}} \underbrace{\prod_{x \in \mathcal{D}} \prod_i p(x_i | Pa_{x_i}; \theta_{X_i|Pa_{X_i}})}_{\text{likelihood}} = \underbrace{\prod_i p(\theta_{X_i|Pa_{X_i}})}_{p(\theta_{X_i|Pa_{X_i}} | \mathcal{D})} \underbrace{\prod_{x \in \mathcal{D}} p(x_i | Pa_{x_i}; \theta_{X_i|Pa_{X_i}})}_{\text{individual posteriors}}$$

- we can apply Bayesian learning to individual conditional distributions
- posterior predictive also decomposes:  $p(x' | \mathcal{D}) = \prod_i p(x'_i | \mathcal{D})$

$$\downarrow$$

$$\int_{\theta} p(\theta_{X_i|Pa_{X_i}} | \mathcal{D}) p(x'_i | Pa_{x'_i}; \theta_{X_i|Pa_{X_i}}) d\theta_{X_i|Pa_{X_i}}$$

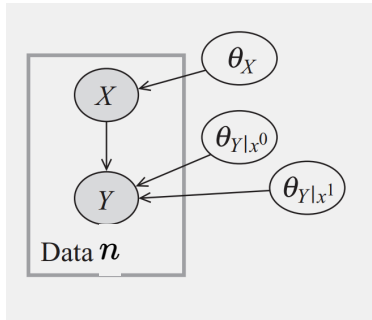
example



# Bayesian learning for Bayes-nets

discrete case: conditional probability tables (CPTs)

we can further decompose the prior & posterior



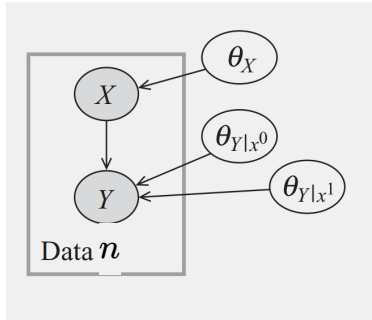
- assume a decomposed prior  $p(\theta_{Y|X}) = p(\theta_{Y|x^0})p(\theta_{Y|x^1})$  for binary  $X$ 
  - one for each assignment to the parent (e.g., cols of the table)

*local parameter independence*

# Bayesian learning for Bayes-nets

discrete case: conditional probability tables (CPTs)

we can further decompose the prior & posterior



- assume a decomposed prior  $p(\theta_{Y|X}) = p(\theta_{Y|x^0})p(\theta_{Y|x^1})$  for binary X
  - one for each assignment to the parent (e.g., cols of the table)

- posterior is also decomposed  $p(\theta_{Y|X} | \mathcal{D}) = p(\theta_{Y|x^0} | \mathcal{D})p(\theta_{Y|x^1} | \mathcal{D})$

$$\downarrow$$
$$p(\theta_{Y|x^0}) \prod_{(x^0, y) \in \mathcal{D}} p(y|x^0; \theta_{Y|x^0})$$



# Bayesian learning for Bayes-nets

**discrete case:** conditional probability tables (CPTs)

In practice this means:

- keep a vector of *pseudo-counts* for each node
- after observing N samples:
  - update these based on the frequency of different (x,y) values

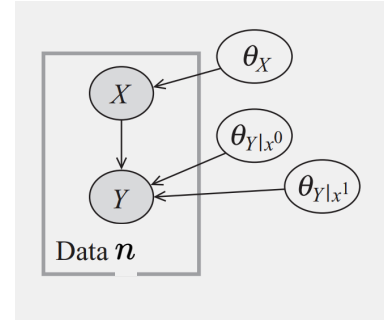
**K2 prior**  $\alpha_{Y|x^0} = \alpha_{Y|x^1} = [1, \dots, 1]$

similar to Laplace smoothing

**BDe prior** use a second Bayes-net to

- keep frequencies  $P'(x_i, pa_{X_i})$
- keep a total pseudo-count  $\alpha$

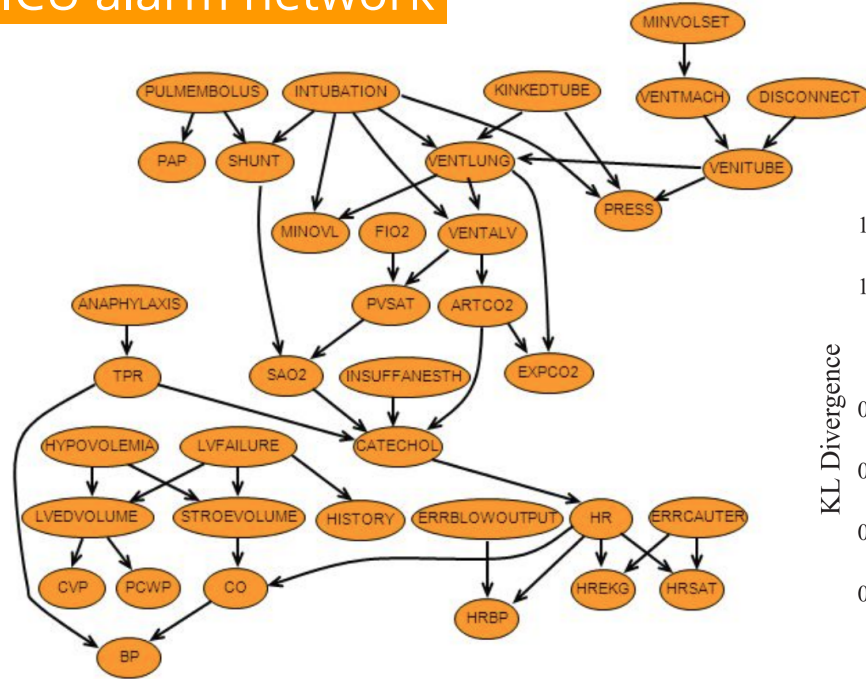
then  $\alpha_{x_i|pa_{X_i}} = \alpha P'(x_i, pa_{X_i})$



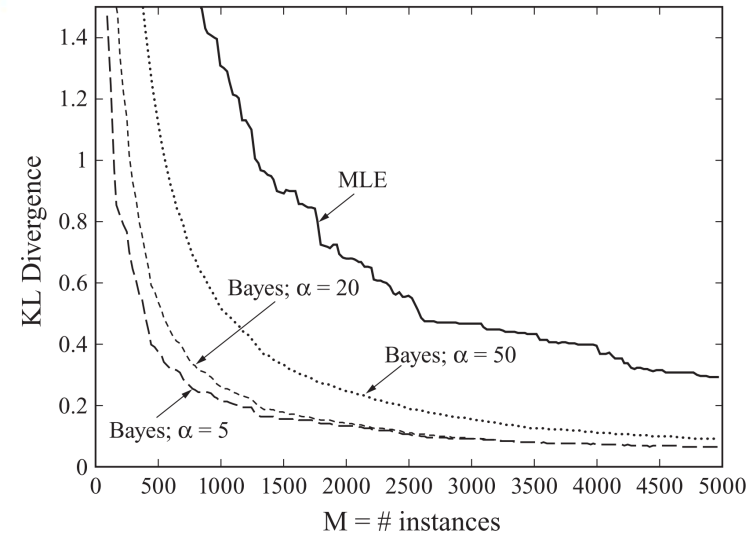
# Bayesian learning for Bayes-nets

example

ICU alarm network



## Bayesian learning vs MLE



# Summary

learn the parameter by **maximizing the likelihood**

it does not reflect uncertainty:

- maintain a distribution over the parameters
- for conjugate pairs (prior-likelihood), this maintenance is easy

# Summary

learn the parameter by **maximizing the likelihood**

it does not reflect uncertainty:

- maintain a distribution over the parameters
- for conjugate pairs (prior-likelihood), this maintenance is easy

In Bayes-nets:

- both **MLE and Bayesian learning is easy**
  - they have a **decomposed form**