

# Parameter Estimation for Bayesian Networks

Sargur Srihari

[srihari@cedar.buffalo.edu](mailto:srihari@cedar.buffalo.edu)

# Topics

- Problem Statement
- Maximum Likelihood Approach
  - Thumb-tack (Bernoulli), Multinomial, Gaussian
  - Application to Bayesian Networks
- Bayesian Approach
  - Thumbtack vs Coin Toss
    - Uniform Prior vs Beta Prior
  - Multinomial
    - Dirichlet Prior

# Problem Statement

- Network structure is fixed
- Data set consists of fully observed instances of the network variables

$$\mathcal{D} = \{\xi[1], \dots, \xi[M]\}$$

- Arises commonly in practice
  - Since hard to elicit from human experts
- Forms building block for structure learning and learning from incomplete data

# Two Main Approaches

- Maximum Likelihood Estimation
- Bayesian
- Start with the simplest context
  - Bayesian network with a single variable
  - Generalize to arbitrary network structures

# Thumbtack Example

- Classical Statistics Problem
  - Thumbtack is tossed to land head/tail



- Tossed several times obtaining either heads or tails
  - Based on data set wish to estimate probability that next flip will land heads/tails
- Probability of heads is  $\theta$  which is the unknown parameter

# Results of Tosses

- Toss 100 times
  - We get 35 heads
- What is our estimate for  $\theta$ ?
  - Intuition suggests 0.35
  - If it was 0.1 chances of 35/100 would be lower
- Law of larger number says
  - as no of tosses grows it is increasingly unlikely that fraction will be far from  $\theta$
  - For large M fraction of heads is a good estimate with high probability
- How to formalize this intuition?

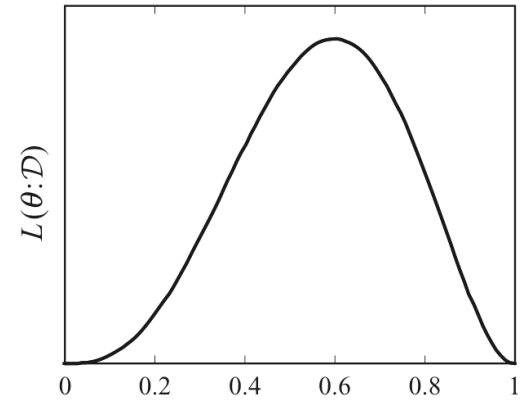
# Maximum Likelihood

- Since tosses are independent, probability of sequence is

$$P(H,T,T,H,H:\theta) = \theta (1-\theta)(1-\theta) \theta\theta = \theta^3 (1-\theta)^2$$

- Probability depends on value of  $\theta$
- Define the likelihood function to be

$$L(\theta:\langle H,T,T,H,H \rangle) = P(H,T,T,H,H:\theta) = \theta^3 (1-\theta)^2$$



$\theta=0.6$

Maximizes likelihood  
of sequence HTTHH

- Parameter values with higher likelihood are more likely to generate the observed sequences
- Thus likelihood is measure of parameter quality

# MLE for General Case

- If there are  $M[1]$  heads and  $M[0]$  tails

- Likelihood function is

$$L(\theta: \mathcal{D}) = \theta^{M[1]} (1-\theta)^{M[0]}$$

- Log-likelihood is

$$l(\theta: \mathcal{D}) = M[1] \log \theta + M[0] \log(1-\theta)$$

- Differentiating and setting equal to zero

$$\hat{\theta} = \frac{M[1]}{M[1] + M[0]}$$

# Maximum Likelihood Principle

- Observe random variables  $\chi$  from unknown distribution  $P^*(\chi)$
- Training samples are  $\mathcal{D} = \{\xi[1], \dots, \xi[M]\}$
- We are given a parametric model defined by  $P(\xi; \theta)$ , we wish to estimate parameters  $\theta$
- For each choice of parameter  $\theta$ ,  $P(\xi; \theta)$  is a legal distribution

$$\sum_{\xi} P(\xi; \theta) = 1$$

- We need to define parameter space  $\Theta$  which is the set of allowable parameters

# Parameter Space

<ul style="list-style-type: none"> <li>• <u>Thumb-tack</u></li> <li>• <math>\Theta_{thumbtack} = [0, 1]</math></li> <li>• <math>P(x:\theta) = \begin{cases} \theta &amp; \text{if } x=H \\ 1-\theta &amp; \text{if } x=T \end{cases}</math></li> <li>• <math>\theta \in [0, 1]</math></li> </ul>	<ul style="list-style-type: none"> <li>• <u>Multinomial</u></li> <li>• <math>X</math> takes one of <math>K</math> values <math>x^1 \dots x^K</math></li> <li>• <math>\theta \in [0, 1]^K</math> <ul style="list-style-type: none"> <li>• <math>K</math> real values in interval</li> </ul> </li> <li>• <math>P(x:\theta) = \theta_k</math> if <math>x = x^k</math></li> <li>• <math>\Theta_{multinomial} = \{\theta \in [0, 1]^K : \sum \theta_i = 1\}</math></li> </ul>	<ul style="list-style-type: none"> <li>• <u>Gaussian</u></li> <li>• <math>P(x:\mu, \sigma) = \frac{1}{\sqrt{2\pi\sigma}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}</math></li> <li>• <math>\Theta_{Gaussian} = R \times R^+</math></li> <li>• Any real value of <math>\mu</math></li> <li>• Any positive value of <math>\sigma</math></li> </ul>
--	--	---

Likelihood Function

$$L(\theta : D) = \prod_m P(\xi[m] : \theta)$$

# Sufficient Statistics Examples

- In thumbtack
  - Likelihood is in terms of  $M[0]$  and  $M[1]$  which are counts of 0s and 1s,  $L(\theta : \mathcal{D}) = \theta^{M[1]} (1-\theta)^{M[0]}$ 
    - No need for other aspects of training data e.g., toss order
    - They are summaries of data for computing likelihood
- For multinomial
  - Tuple of counts  $M[1], \dots, M[K]$  which are counts of  $x^k$ 
    - Number of times  $x^k$  appears in training data
  - Sum of instance level statistics
    - obtained using  $\tau(x^k) = 0, \dots, 0, 1, 0, \dots, 0$
  - Likelihood function is

$$\left\{ \begin{array}{l} \text{1-of-}K \text{ representation} \\ \sum_{i=1}^N X_i \text{ yields } K\text{-dim vector} \end{array} \right.$$

$$L(\theta : D) = \prod_k \theta_k^{M[k]}$$

# Sufficient Statistics: Gaussian

$$P_{Gaussian}(x : \mu, \sigma) = e^{-x^2 \frac{1}{2\sigma^2} + x \frac{\mu}{\sigma^2} - \frac{\mu^2}{2\sigma^2} - \frac{1}{2} \log(2\pi) - \log \sigma}$$

Likelihood function will involve  $\sum x_i^2$ ,  $\sum x_i$  and  $\sum 1$

Sufficient statistic is to apply the following function to each sample

$$S_{Gaussian}(x) = \langle 1, x, x^2 \rangle$$

# Maximum Likelihood Estimates

- Given data set  $\mathcal{D}$  choose parameters that satisfy  $L(\hat{\theta} : D) = \max_{\theta \in \Theta} L(\theta : D)$

- Multinomial

$$\hat{\theta}_k = \frac{M[k]}{M}$$

- Gaussian

$$\hat{\mu} = \frac{1}{M} \sum_m x[m]$$

$$\hat{\sigma} = \sqrt{\frac{1}{M} \sum_m (x[m] - \hat{\mu})^2}$$

# MLE for Bayesian Networks

- Structure of Bayesian network allows us to reduce parameter estimation problem into a set of unrelated problems

# Simple Bayesian Network

- Network with two binary variables  $X \rightarrow Y$
- Parameter vector  $\theta$  defines CPD parameters
- Each training sample is a tuple  $\langle x[m], y[m] \rangle$

$$L(\theta : D) = \prod_{m=1}^M P(x[m], y[m] : \theta)$$

- **Network model specifies**  $P(X, Y) = P(X)P(Y|X)$

$$L(\theta : D) = \prod_{m=1}^M P(x[m] : \theta) P(y[m] | x[m] : \theta)$$

changing order of multiplication

$$L(\theta : D) = \prod_{m=1}^M P(x[m] : \theta) \prod_{m=1}^M P(y[m] | x[m] : \theta)$$

Likelihood decomposes into two separate terms, each a local likelihood function  
Measures how well variable is predicted given its parents

# Bayesian Parameter Estimation

- Thumbtack example
  - Toss tack and get 3 heads out of 10
    - Conclude that parameter  $\theta$  is set to 0.3
- Coin toss example
  - Get 3 heads out of 10, Can we conclude  $\theta = 0.3$ ?
    - No. We have lot more experience and have prior knowledge about their behavior
    - If we observe 1 million tosses and 300,000 came out heads we conclude a trick coin with  $\theta$  is 0.3
- MLE cannot distinguish between
  1. Thumbtack and coin
  2. Between 10 tosses and  $10^6$  tosses

# Joint Probabilistic Model

- Encode prior knowledge about  $\theta$  with a probability distribution
- Represents how likely we are about different choices of parameters
- Create a joint distribution  $\theta$  involving parameter and samples  $X[1], \dots, X[M]$
- Joint distribution captures our assumptions about experiment

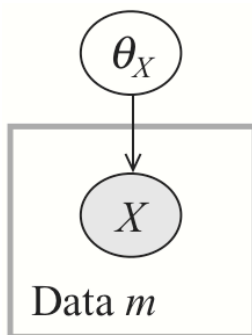
# Toss Experiment

- In MLE Tosses were assumed independent
  - Assumption made when  $\theta$  was fixed
- If we do not know  $\theta$  then one toss tells us something about the parameter  $\theta$ 
  - Thus about next toss
  - Thus tosses are not marginally independent
- Once we know  $\theta$  we cannot learn probability of one toss from others
  - So tosses are conditionally independent give  $\theta$

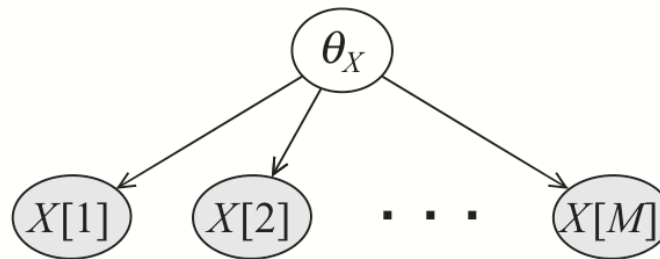
# Probabilistic model of Experiments

Meta network of independent identically distributed samples

Plate Model



Ground Bayesian network



Joint distribution needs probability sample given parameter and probability of parameter

# Local probability models

## 1. Probability distribution of sample given parameter $\theta$

$$P(x[m] | \theta) = \begin{cases} \theta & \text{if } x[m] = x^1 \\ 1 - \theta & \text{if } x[m] = x^0 \end{cases}$$

It is equal to  $\theta$  if sample is heads and  $1-\theta$  if sample is tails

## 2. Probability distribution of parameter $\theta$

Prior distribution is a continuous density over interval  $[0,1]$

## Substituting into joint probability distribution

$P(x[1], \dots, x[M], \theta) = P(x[1], \dots, x[M] | \theta)P(\theta)$  by Bayes rule

$$\begin{aligned} &= P(\theta) \prod_{m=1}^M P(x[m] | \theta) \quad \text{Given } \theta \text{ samples independent} \\ &= P(\theta) \theta^{M^{[1]}} (1-\theta)^{M^{[0]}} \end{aligned}$$

# Using the Probability Model

## 1. Determine posterior distribution

- From observed data set  $\mathcal{D}$  instantiate values  $x[1], \dots, x[M]$  then compute posterior over  $\theta$

$$P(\theta | x[1], \dots, x[M]) = \frac{P(x[1], \dots, x[M] | \theta) P(\theta)}{P(x[1], \dots, x[M])} \quad P(\theta | \mathcal{D}) = \frac{P(\mathcal{D} | \theta) P(\theta)}{P(\mathcal{D})}$$

- Posterior proportional to: likelihood  $\times$  prior
- If  $P(\theta)$  is uniform, i.e.,  $P(\theta) = 1$  for all  $\theta \in [0, 1]$  posterior is just the normalized likelihood function

## 2. Predict probability over next toss

# Prediction

- For predicting probability of heads in next toss
  - In Bayesian approach instead of single value for  $\theta$  we use its distribution
- Parameter is unknown, so we consider all possible values

$$P(x[M+1] | D) = \int P(x[M+1] | \theta, D) P(\theta | D) d\theta \quad \text{By Bayes rule \& Sum rule}$$

$$= \int P(x[M+1] | \theta) P(\theta | D) d\theta \quad \text{since samples are independent given } \theta$$

- Integrating posterior over  $\theta$  to predict prob heads in next toss

- Thumbtack with Uniform Prior

- Assume uniform prior over  $[0, 1]$ , i.e.,  $P(\theta) = 1$

- Substituting  $P(\theta | \mathcal{D}) = \frac{P(\mathcal{D} | \theta)P(\theta)}{P(\mathcal{D})}$  and  $P(\mathcal{D} | \theta) = \theta^{M[1]}(1-\theta)^{M[0]}$

- We get

$$P(X[M+1] = x^1 | D) = \frac{M[1]+1}{M[1]+M[0]+2}$$

Similar to MLE estimate except adds an imaginary sample to each count  
Called “Laplace’s Correction”

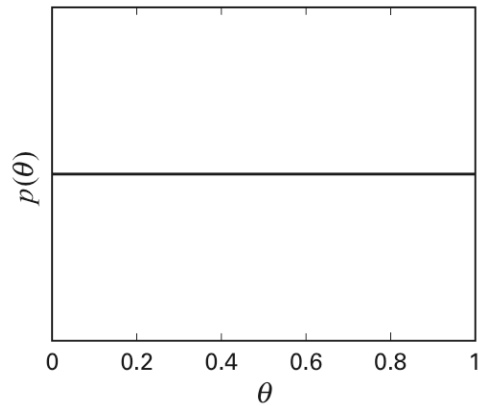
$$\hat{\theta} = \frac{M[1]}{M[1]+M[0]}$$

# Non-uniform Prior

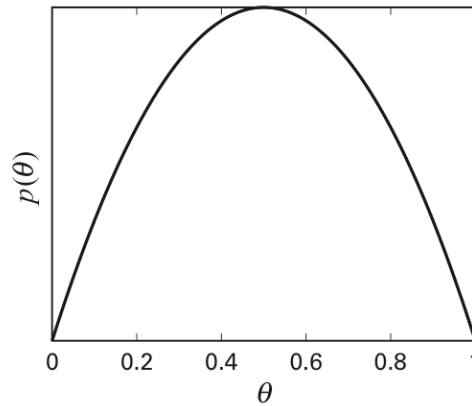
- Pick a distribution that can be represented
  1. Compactly
  2. Updated efficiently as we get new data
- Beta Distribution
  - Parameterized by two hyper-parameters
    - $\alpha_0$  and  $\alpha_1$  which are positive reals
    - They are imaginary numbers of heads and tails we have seen before starting the experiment

# Examples of Beta Distribution

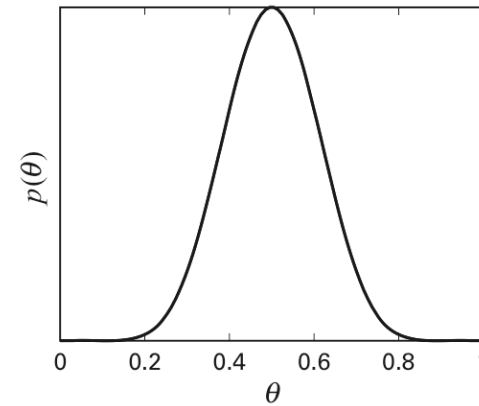
## For different hyperparameters



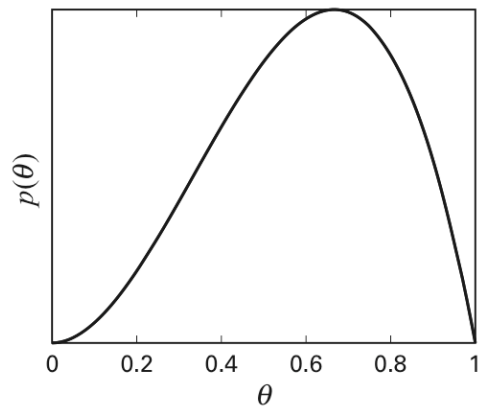
*Beta(1,1)*



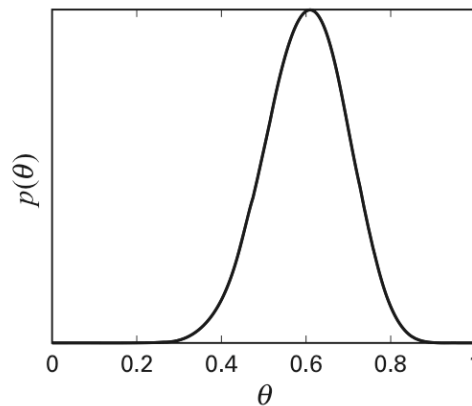
*Beta(2,2)*



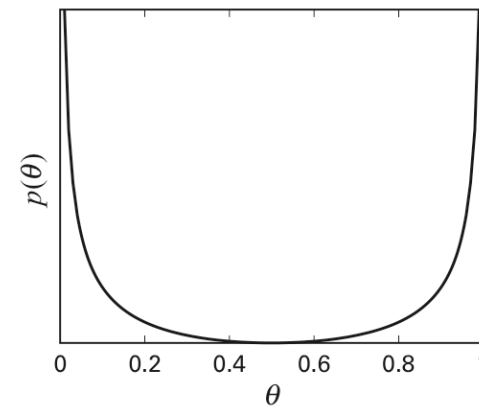
*Beta(10,10)*



*Beta(3,2)*



*Beta(15,10)*



*Beta(0.5,0.5)*

# Details of Beta Distribution

- Defined as

$$\theta \sim \text{Beta}(\alpha_1, \alpha_0) \text{ if } p(\theta) = \gamma \theta^{\alpha_1 - 1} (1 - \theta)^{\alpha_0 - 1}$$

- $\gamma$  is a normalizing constant, defined as

$$\gamma = \frac{\Gamma(\alpha_1 + \alpha_0)}{\Gamma(\alpha_1)\Gamma(\alpha_0)} \text{ where } \Gamma(x) = \int_0^{\infty} t^{x-1} e^{-t} dt \text{ is the Gamma function}$$

- Gamma function is a continuous generalization of factorials

$$\Gamma(1) = 1 \text{ and } \Gamma(x + 1) = x\Gamma(x)$$

Thus  $\Gamma(n + 1) = n!$  when  $n$  is large

# Thumbtack with Beta Prior

- Marginal probability of  $X$  based on  $P(\theta)$

$$\begin{aligned}
 P(X[1] = x^1) &= \int_0^1 P(X[1] = x^1 | \theta) \cdot P(\theta) d\theta \\
 &= \int_0^1 \theta \cdot P(\theta) d\theta = \frac{\alpha_1}{\alpha_1 + \alpha_0} \quad \text{Using standard integration techniques}
 \end{aligned}$$

$$\begin{aligned}
 P(\theta | \mathcal{D}) &\propto P(\mathcal{D} | \theta) P(\theta) \\
 &= \theta^{\alpha_1 + M[1] - 1} (1 - \theta)^{\alpha_0 + M[0] - 1}
 \end{aligned}$$

Which is precisely  $Beta(\alpha_1 + M[1], \alpha_0 + M[0])$

- Key property of Beta distribution
  - Prior conditioned on evidence is also a Beta distribution
    - Beta is conjugate with Bernoulli
  - Immediate consequence is computing probability over next toss:

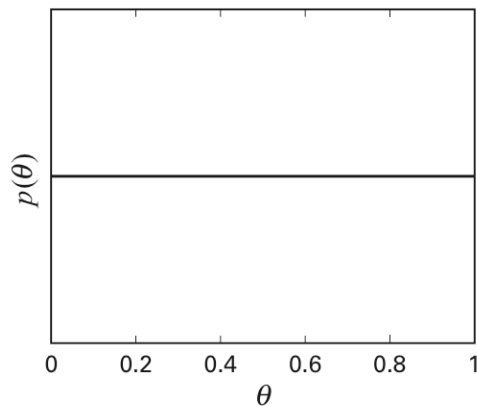
$$P(X[M+1] = x^1 | D) = \frac{M[1] + \alpha_1}{M[1] + M[0] + \alpha_1 + \alpha_0}$$

Compare with uniform prior  $P(X[M+1] = x^1 | D) = \frac{M[1] + 1}{M[1] + M[0] + 2}$

# Bayesian Approach overcomes both shortcomings of MLE

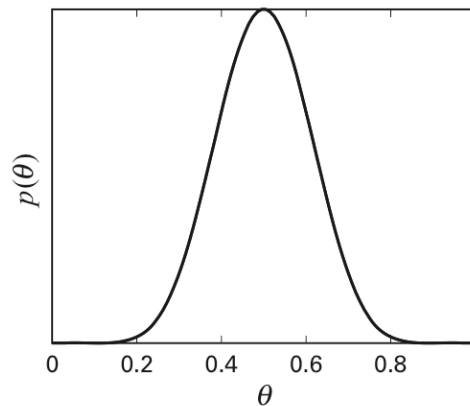
- Strength of prior captured by choice of hyper-parameters

Thumb-tack:  $\alpha_1 = \alpha_0 = 1$



$Beta(1,1)$

Coin:  $\alpha_1 = \alpha_0 = 100$



$Beta(10,10)$

- Distinction between few samples and many is captured by peakedness of posterior

# Multinomial Case: Bayesian Approach

- Multinom. Distribn

- $X$  takes one of  $K$  values  $x^1 \dots x^K$
- $\boldsymbol{\theta} \in [0, 1]^K$  i.e,  $K$  reals in  $[0, 1]$
- $P(x: \boldsymbol{\theta}) = \theta_k$  if  $x = x^k$

- Likelihood function

(probability of iid samples)

$$L(\boldsymbol{\theta}: \mathcal{D}) = \prod_k \theta_k^{M[k]}$$

- Since posterior is prior x likelihood
  - Require prior to have form similar to likelihood
  - Dirichlet has that form

$$\boldsymbol{\theta} \sim \text{Dirichlet}(\alpha_1, \dots, \alpha_k) \text{ if } P(\boldsymbol{\theta}) \propto \prod_k \theta_k^{\alpha_k - 1}$$

# Dirichlet Distribution

- Generalizes the Beta Distribution
- Specified by a set of hyper-parameters
- If  $P(\boldsymbol{\theta})$  is *Dirichlet*( $\alpha_1, \dots, \alpha_K$ ) then
  1.  $P(\boldsymbol{\theta} | \mathcal{D})$  is *Dirichlet*( $\alpha_1 + M[1], \dots, \alpha_K + M[K]$ )
  2.  $E[\theta_k] = \alpha_k / \alpha$