



Computer  
Science

# CSC580: Probabilistic Graphical Models

## Probabilistic Graphical Models

Jason Pacheco

# Administrivia

- Homework submission
  - Make sure questions are answered in PDF
  - Match pages to questions
  - Put code in PDF (relevant parts of code at least)
  - Doublecheck your submission
- Midterm Exam
  - Thursday 10/12
  - No coding
  - Probably closed-book

- Probability Refresher
- Probabilistic Graphical Models
- Naïve Bayes

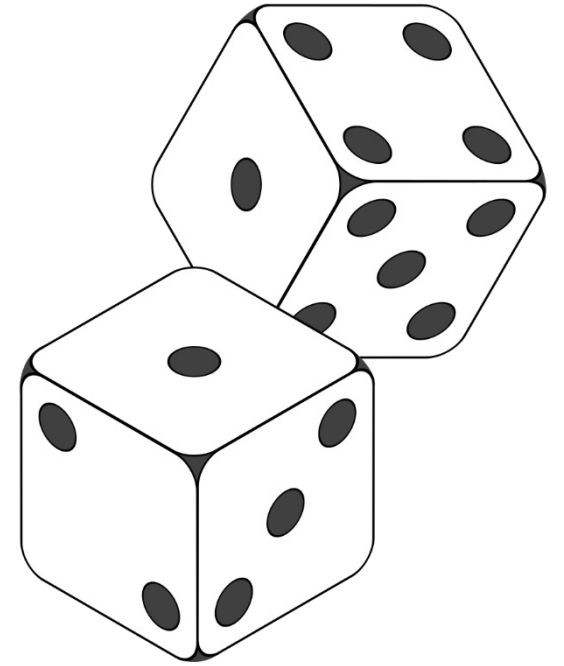
- Probability Refresher
- Probabilistic Graphical Models
- Naïve Bayes

Before we learn about probabilistic graphical models, we need to review probability...

# Random Events and Probability

***Suppose we roll two fair dice...***

- What are the possible outcomes?
- What is the *probability* of rolling **even** numbers?
- What is the *probability* of rolling **odd** numbers?



***...probability theory gives a mathematical formalism to addressing such questions...***

**Definition** An **experiment** or **trial** is any process that can be repeated with well-defined outcomes. It is *random* if more than one outcome is possible.

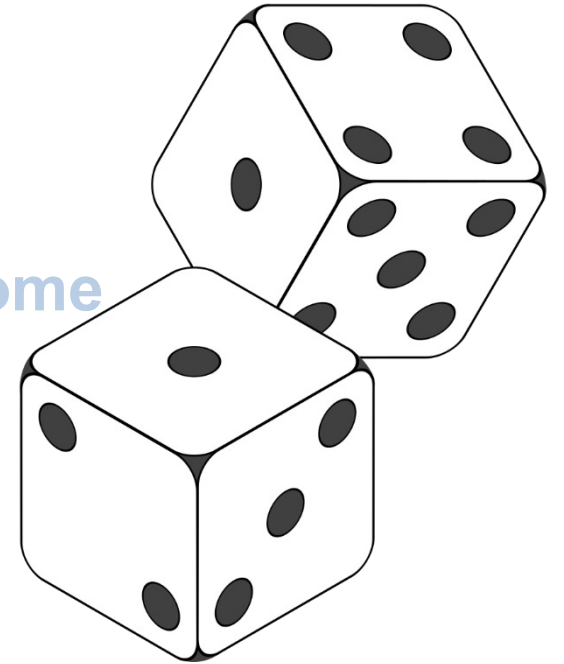
# Random Events and Probability

**Definition** An **outcome** is a possible result of an experiment or trial, and the collection of all possible outcomes is the **sample space** of the experiment,

**Example**  $(1,1), (1,2), \dots, (6,1), (6,2), \dots, (6,6)$

Sample Space

Outcome



**Definition** An **event** is a *set* of outcomes (a subset of the sample space),

**Example Event** Roll at least a single 1

$\{(1,1), (1,2), (1,3), \dots, (1,6), \dots, (6,1)\}$

# Random Variables

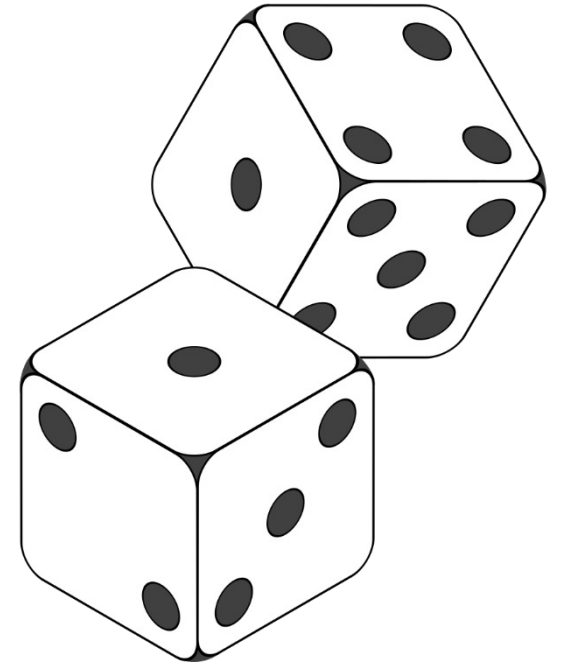
*(Informally) A random variable is an unknown quantity that maps events to numeric values.*

**Example**  $X$  is the *sum of two dice* with values,

$$X \in \{2, 3, 4, \dots, 12\}$$

**Example** Flip a coin and let random variable  $Y$  represent the outcome,

$$Y \in \{\text{Heads}, \text{Tails}\}$$



# Random Variables and Probability

Capitol letters represent  
random variables

Lowercase letters are  
realized *values*

$$X = x$$

$X = x$  is the **event** that  $X$  takes the value  $x$

**Example** Let  $X$  be the random variable (RV) representing the sum of two dice with values,

$$X \in \{2, 3, 4, \dots, 12\}$$

$X=5$  is the *event* that the dice sum to 5.



# Probability Mass Function

A function  $p(X)$  is a **probability mass function (PMF)** of a discrete random variable if the following conditions hold:

(a) It is nonnegative for all values in the support,

$$p(X = x) \geq 0$$

(b) The sum over all values in the support is 1,

$$\sum_x p(X = x) = 1$$

**Intuition** Probability mass is conserved, just as in physical mass. Reducing probability mass of one event must increase probability mass of other events so that the definition holds...

# Probability Mass Function

**Example** Let  $X$  be the outcome of a single fair die. It has the PMF,

$$p(X = x) = \frac{1}{6} \quad \text{for } x = 1, \dots, 6 \quad \text{Uniform Distribution}$$

**Example** We can often represent the PMF as a vector. Let  $S$  be an RV that is the *sum of two fair dice*. The PMF is then,

**Observe that  $S$  does not follow a uniform distribution**

$$p(S) = \begin{pmatrix} p(S = 2) \\ p(S = 3) \\ p(S = 4) \\ \vdots \\ p(S = 12) \end{pmatrix} = \begin{pmatrix} 1/36 \\ 1/18 \\ 1/2 \\ \vdots \\ 1/36 \end{pmatrix}$$

# PMF Notation

- We use  $p(X)$  to refer to the probability mass *function* (i.e. a function of the RV  $X$ )
- We use  $p(X=x)$  to refer to the probability of the *outcome*  $X=x$  (also called an “event”)
- We will often use  $p(x)$  as shorthand for  $p(X=x)$

# Joint Probability

**Definition** Two (discrete) RVs  $X$  and  $Y$  have a *joint PMF* denoted by  $p(X, Y)$  and the probability of the event  $X=x$  and  $Y=y$  denoted by  $p(X = x, Y = y)$  where,

(a) It is nonnegative for all values in the support,

$$p(X = x, Y = y) \geq 0$$

(b) The sum over all values in the support is 1,

$$\sum_x \sum_y p(X = x, Y = y) = 1$$

# Joint Probability

Let  $X$  and  $Y$  be *binary RVs*. We can represent the joint PMF  $p(X, Y)$  as a 2x2 array (table):

		Y	
		0	1
X	0	0.04	0.36
	1	0.30	0.30

**All values are nonnegative**

# Joint Probability

Let  $X$  and  $Y$  be *binary RVs*. We can represent the joint PMF  $p(X,Y)$  as a 2x2 array (table):

		Y	
		0	1
X	0	0.04	0.36
	1	0.30	0.30

**The sum over all values is 1:  
 $0.04 + 0.36 + 0.30 + 0.30 = 1$**

# Joint Probability

Let  $X$  and  $Y$  be *binary RVs*. We can represent the joint PMF  $p(X, Y)$  as a 2x2 array (table):

		Y	
		0	1
X	0	0.04	0.36
	1	0.30	0.30

$$P(X=1, Y=0) = 0.30$$

# Fundamental Rules of Probability

Given two RVs  $X$  and  $Y$  the **conditional distribution** is:

$$p(X | Y) = \frac{p(X, Y)}{p(Y)} = \frac{p(X, Y)}{\sum_x p(X=x, Y)}$$

Multiply both sides by  $p(Y)$  to obtain the **probability chain rule**:

$$p(X, Y) = p(Y)p(X | Y)$$

The probability chain rule extends to  $N$  RVs  $X_1, X_2, \dots, X_N$ :

$$p(X_1, X_2, \dots, X_N) = p(X_1)p(X_2 | X_1) \dots p(X_N | X_{N-1}, \dots, X_1)$$

Chain rule valid  
for any ordering

$$= p(X_1) \prod_{i=2}^N p(X_i | X_{i-1}, \dots, X_1)$$



# Fundamental Rules of Probability

## Law of total probability

$$p(Y) = \sum_x p(Y, X = x)$$

- $P(y)$  is a **marginal** distribution
- This is called **marginalization**

**Proof**

$$\begin{aligned} \sum_x p(Y, X = x) &= \sum_x p(Y) p(X = x | Y) && \text{( chain rule )} \\ &= p(Y) \sum_x p(X = x | Y) && \text{( distributive property )} \\ &= p(Y) && \text{( PMF sums to 1 )} \end{aligned}$$

*Generalization for conditionals:*

$$p(Y | Z) = \sum_x p(Y, X = x | Z)$$

# Tabular Method

Let  $X, Y$  be binary RVs with the joint probability table

For Binomial use K-by-K probability table.

		Y	
		$y_1$	$y_2$
X	$x_1$	0.04	0.36
	$x_2$	0.30	0.30

0.4  $P(x_1)$

0.6  $P(x_2)$

$P(x)$

$P(y_1) = P(x_1, y_1) + P(x_2, y_1)$   
 $P(y_2) = P(x_1, y_2) + P(x_2, y_2)$   
[i.e., sum down columns]

0.34  $P(y_1)$

0.66  $P(y_2)$

$P(y)$

$P(x_1) = P(x_1, y_1) + P(x_1, y_2)$   
 $P(x_2) = P(x_2, y_1) + P(x_2, y_2)$   
[i.e., sum across rows]

# Tabular Method

We don't care about event  $Y=y_2$

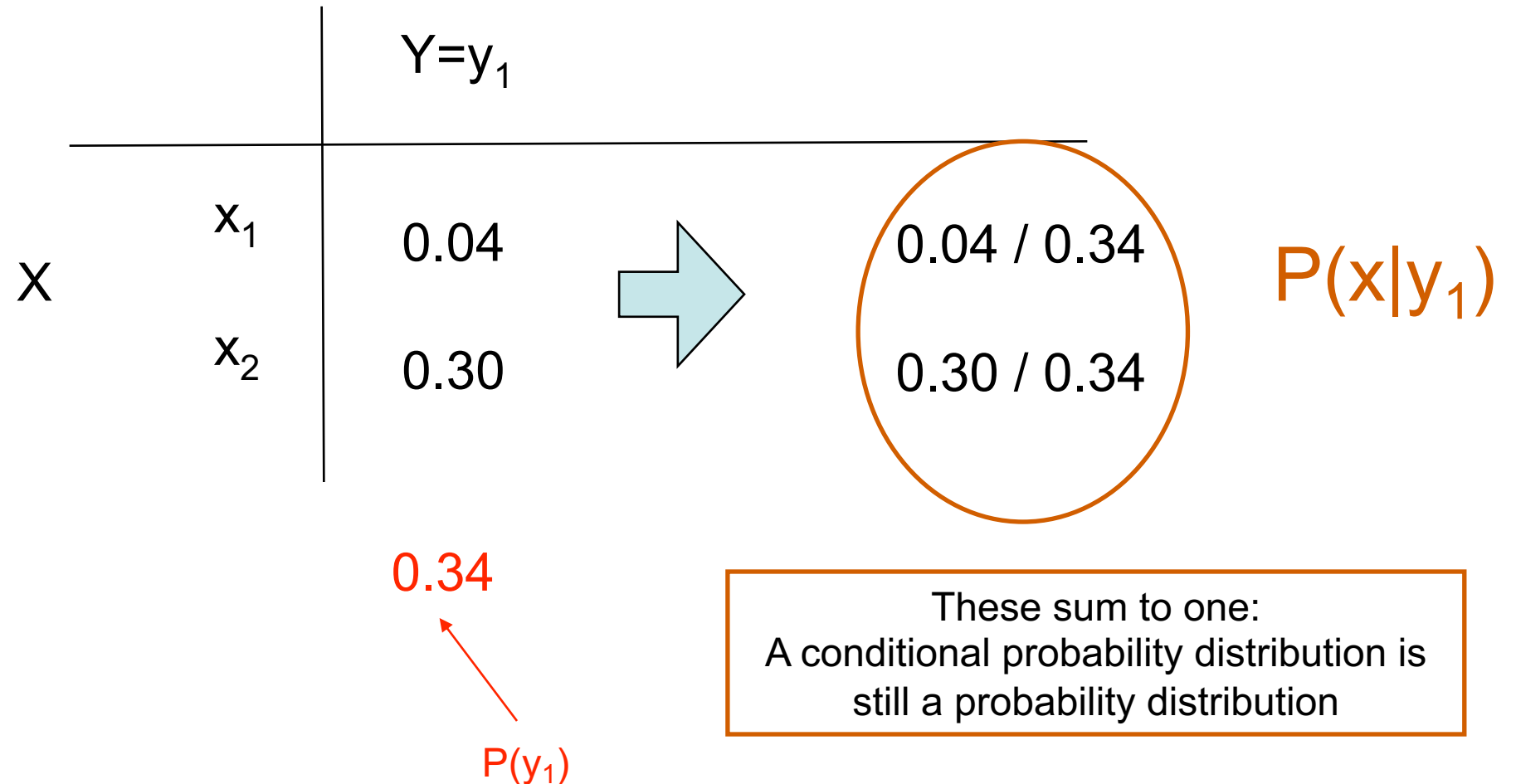
		Y	
		$y_1$	$y_2$
X	$x_1$	0.04	Censored!
	$x_2$	0.30	

$P(x|y_1)=?$

0.34

$P(y_1)$

# Tabular Method



# Intuition Check

Question: Roll two dice and let their outcomes be  $X_1, X_2 \in \{1, \dots, 6\}$  for die 1 and die 2, respectively. Recall the definition of conditional probability,

$$p(X_1 | X_2) = \frac{p(X_1, X_2)}{p(X_2)}$$

Which of the following are true?

a)  $p(X_1 = 1 | X_2 = 1) > p(X_1 = 1)$

b)  $p(X_1 = 1 | X_2 = 1) = p(X_1 = 1)$

Outcome of die 2 doesn't *affect* die 1

c)  $p(X_1 = 1 | X_2 = 1) < p(X_1 = 1)$

# Intuition Check

Question: Let  $X_1 \in \{1, \dots, 6\}$  be outcome of die 1, as before. Now let  $X_3 \in \{2, 3, \dots, 12\}$  be the sum of both dice. Which of the following are true?

a)  $p(X_1 = 1 | X_3 = 3) > p(X_1 = 1)$

b)  $p(X_1 = 1 | X_3 = 3) = p(X_1 = 1)$

c)  $p(X_1 = 1 | X_3 = 3) < p(X_1 = 1)$

Only 2 ways to get  $X_3 = 3$ , each with equal probability:

$$(X_1 = 1, X_2 = 2) \quad \text{or} \quad (X_1 = 2, X_2 = 1)$$

so

$$p(X_1 = 1 | X_3 = 3) = \frac{1}{2} > \frac{1}{6} = p(X_1 = 1)$$

# Dependence of RVs

Intuition...

Consider  $P(B|A)$  where you want to bet on  $B$

Should you pay to know  $A$ ?

In general you would pay something for  $A$  if it changed your belief about  $B$ . In other words if,

$$P(B|A) \neq P(B)$$

# Independence of RVs

**Definition** Two random variables  $X$  and  $Y$  are independent if and only if,

$$p(X = x, Y = y) = p(X = x)p(Y = y)$$

for all values  $x$  and  $y$ , and we say  $X \perp Y$ .

**Definition** RVs  $X_1, X_2, \dots, X_N$  are mutually independent if and only if,

$$p(X_1 = x_1, \dots, X_N = x_N) = \prod_{i=1}^N p(X_i = x_i)$$

- Independence is *symmetric*:  $X \perp Y \Leftrightarrow Y \perp X$
- Equivalent definition of independence:  $p(X | Y) = p(X)$



# Independence of RVs

**Definition** Two random variables  $X$  and  $Y$  are conditionally independent given  $Z$  if and only if,

$$p(X = x, Y = y \mid Z = z) = p(X = x \mid Z = z)p(Y = y \mid Z = z)$$

for all values  $x$ ,  $y$ , and  $z$ , and we say that  $X \perp Y \mid Z$ .

➤  $N$  RVs conditionally independent, given  $Z$ , if and only if:

$$p(X_1, \dots, X_N \mid Z) = \prod_{i=1}^N p(X_i \mid Z)$$

Shorthand notation  
Implies for all  $x, y, z$

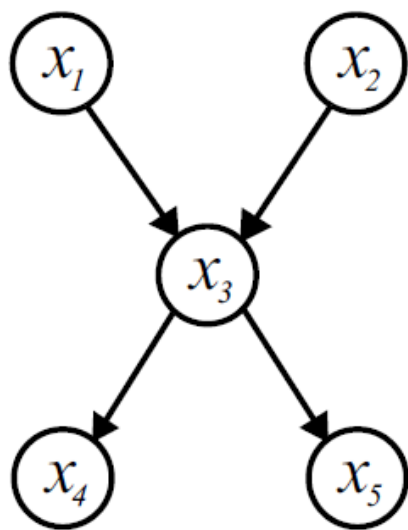
➤ Equivalent def'n of conditional independence:  $p(X \mid Y, Z) = p(X \mid Z)$

➤ Symmetric:  $X \perp Y \mid Z \Leftrightarrow Y \perp X \mid Z$

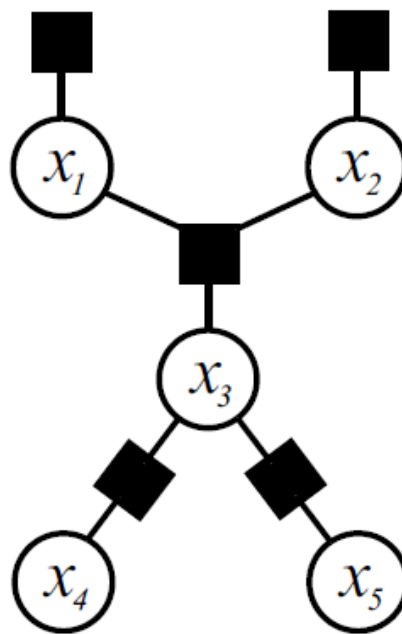
- Probability Refresher
- **Probabilistic Graphical Models**
- Naïve Bayes

# Graphical Models

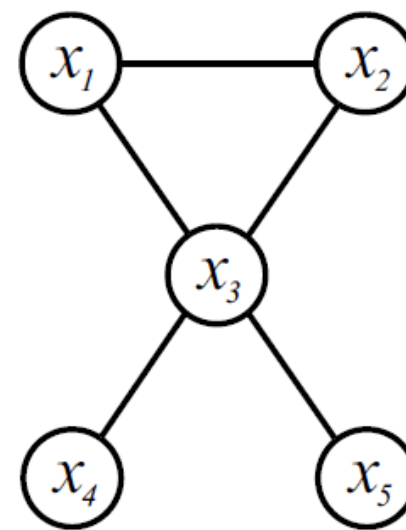
*A variety of graphical models can represent the same probability distribution*



**Bayes Network**



**Factor Graph**



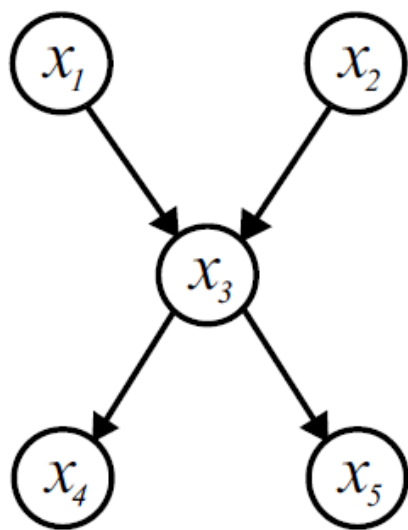
**Markov Random Field**

**Directed Models**

**Undirected Models**

# Graphical Models

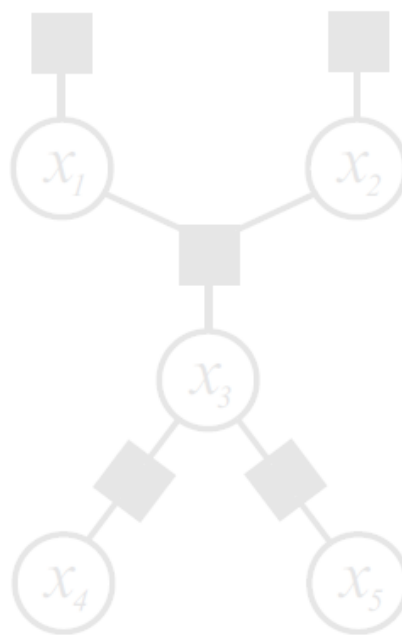
*A variety of graphical models can represent the same probability distribution*



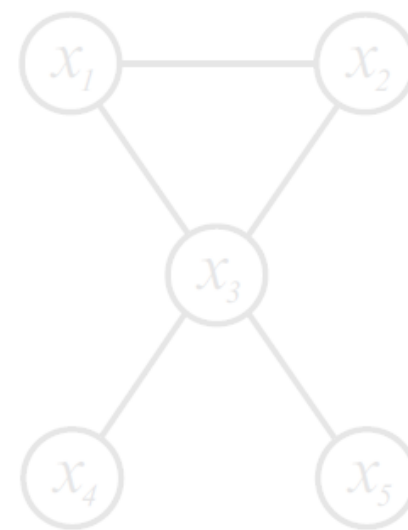
**Bayes Network**



**Directed Models**



**Factor Graph**



**Markov Random Field**

**Undirected Models**

# From Probabilities to Pictures

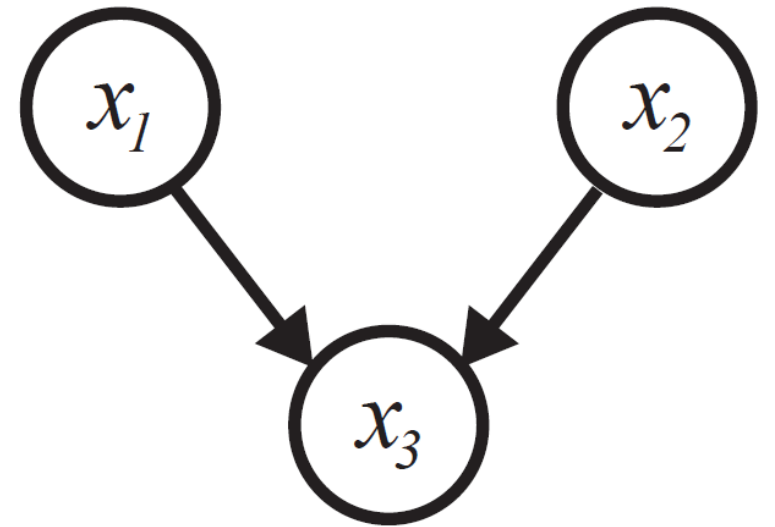
*A probabilistic graphical model allows us to pictorially represent a probability distribution*

**Probability Model:**

$$p(x_1, x_2, x_3) = p(x_1)p(x_2)p(x_3 \mid x_1, x_2)$$



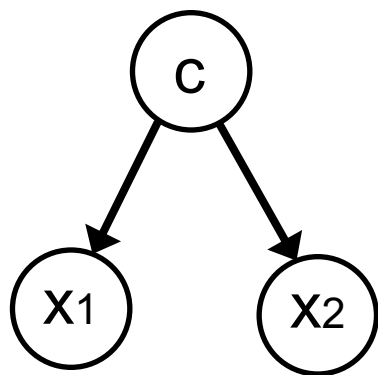
**Graphical Model:**



Conditional distribution on each RV is dependent on its parent nodes in the graph

# Directed Graphical Models

*Directed models are generative models...*



$$p(C, X_1, X_2) = p(C)p(X_1 | C)p(X_2 | C)$$

The graph and the formula say exactly the same thing.  
(The graph has very specific semantics.)

...tells how data are generated (called ***ancestral sampling***)

**Step 1** Sample root node (prior):  $c \sim p(C)$

**Step 2** Sample children, given sample of parent (likelihood):

$$x_1 \sim p(X_1 | C = c) \qquad x_2 \sim p(X_2 | C = c)$$

# Probability Chain Rule

Recall the **probability chain rule** says that we can decompose any joint distribution as a product of conditionals....

$$p(x_1, x_2, x_3, x_4) = p(x_1)p(x_2 | x_1)p(x_3 | x_1, x_2)p(x_4 | x_1, x_2, x_3)$$

Valid for *any ordering* of the random variables...

$$p(x_1, x_2, x_3, x_4) = p(x_3)p(x_1 | x_3)p(x_4 | x_1, x_3)p(x_2 | x_1, x_3, x_4)$$

For a collection of N RVs and any permutation  $\rho$  :

$$p(x_1, \dots, x_N) = p(x_{\rho(1)}) \prod_{i=2}^N p(x_{\rho(i)} | x_{\rho(i-1)}, \dots, x_{\rho(1)})$$

# Conditional Independence

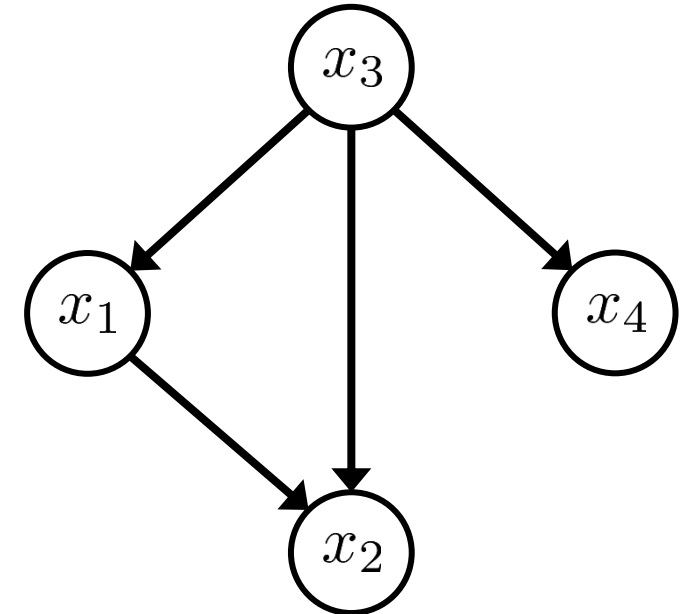
Recall two RVs  $X$  and  $Y$  are **conditionally independent** given  $Z$  (or  $X \perp Y \mid Z$ ) iff:

$$p(X \mid Y, Z) = p(X \mid Z)$$

**Idea** Apply *chain rule* with ordering that exploits conditional independencies to simplify the terms

**Ex.** Suppose  $x_4 \perp x_1 \mid x_3$  and  $x_2 \perp x_4 \mid x_1$  then:

$$\begin{aligned} p(x) &= p(x_3)p(x_1 \mid x_3)p(x_4 \mid x_1, x_3)p(x_2 \mid x_1, x_3, x_4) \\ &= p(x_3)p(x_1 \mid x_3)p(x_4 \mid x_3)p(x_2 \mid x_1, x_3) \end{aligned}$$

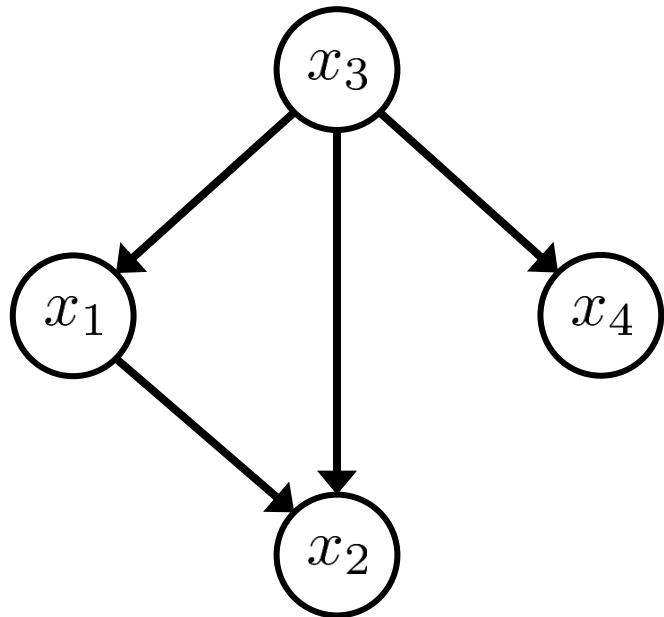


Can visualize conditional dependencies using **directed acyclic graph (DAG)**



# General Directed Graphs

**Def.** A directed graph is a graph with edges  $(s, t) \in \mathcal{E}$  (arcs) connecting parent vertex  $s \in \mathcal{V}$  to a child vertex  $t \in \mathcal{V}$



**Def.** Parents of vertex  $t \in \mathcal{V}$  are given by the set of nodes with arcs pointing to  $t$ ,

$$\text{Pa}(t) = \{s : (s, t) \in \mathcal{E}\}$$

Children of  $t \in \mathcal{V}$  are given by the set,

$$\text{Ch}(t) = \{t : (t, k) \in \mathcal{E}\}$$

Ancestors are parents-of-parents.

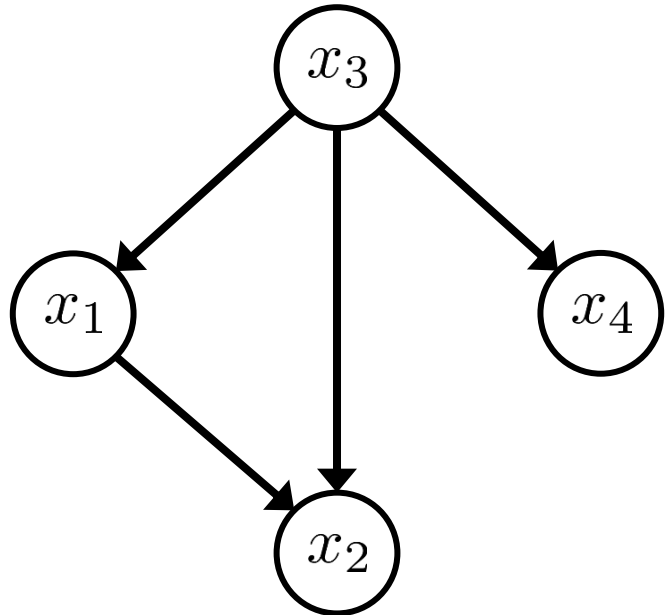
Descendants are children-of-children.

# Directed PGM = Bayes Network

Model factors are normalized conditional distributions:

$$p(x) = \prod_{s \in \mathcal{V}} p(x_s \mid x_{\text{Pa}(s)})$$

 Parents of node  $s$

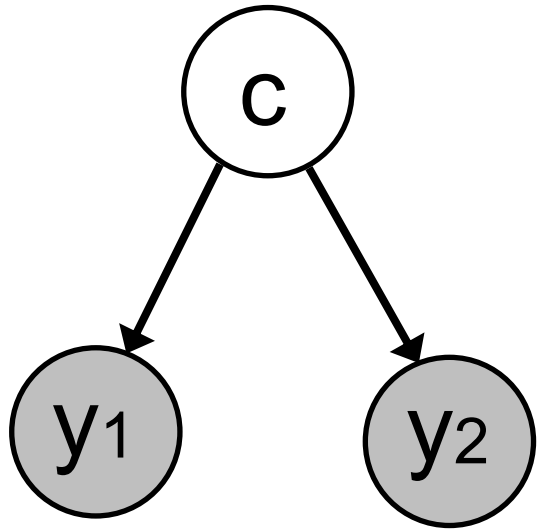


**Directed acyclic graph (DAG)** specifies factorized form of joint probability:

$$p(x) = p(x_3)p(x_1 \mid x_3)p(x_4 \mid x_3)p(x_2 \mid x_1, x_3)$$

*Locally normalized factors yield globally normalized joint probability*

# Inference



Denote observed data with shaded nodes,

$$Y_1 = y_1 \quad Y_2 = y_2$$

Infer *latent* variable C via Bayes' rule:

$$p(c | y_1, y_2) = \frac{p(c)p(y_1 | c)p(y_2 | c)}{p(y_1, y_2)}$$

- This is (obviously) a simple example
- Models and inference task can get really complicated
- But the fundamental concepts and approach are the same

# Bayes' Rule

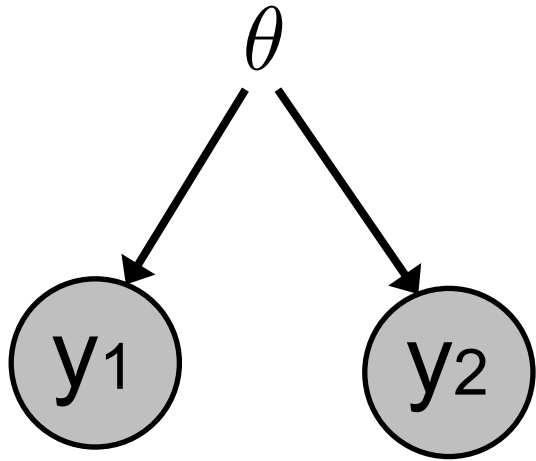
*Posterior represents all uncertainty after observing data...*

The diagram shows the Bayes' Rule equation with four labels and arrows pointing to the corresponding parts of the equation:

- prior** probability: points to  $p(c)$  in the numerator.
- likelihood** function for the parameters: points to  $p(y | c)$  in the numerator.
- posterior** probability: points to  $p(c | y)$  on the left side of the equation.
- marginal likelihood** or: **evidence** or: **partition function** or: **normalizer**: points to  $p(y)$  in the denominator.

$$p(c | y) = \frac{p(c)p(y | c)}{p(y)}$$

# Learning / Training



Model random data with hyperparameters  $\theta$  :

$$y \sim p(y | \theta)$$

Sometimes we use:

$$p(y; \theta)$$

Given training data:

$$\{y_i\}_{i=1}^n \stackrel{\text{i.i.d.}}{\sim} p(y | \theta)$$

Learn parameters, e.g. via *maximum likelihood estimation*:

$$\hat{\theta}^{\text{MLE}} = \arg \max_{\theta} \log p(y_1, \dots, y_n | \theta)$$

We will talk more  
about MLE in  
coming weeks

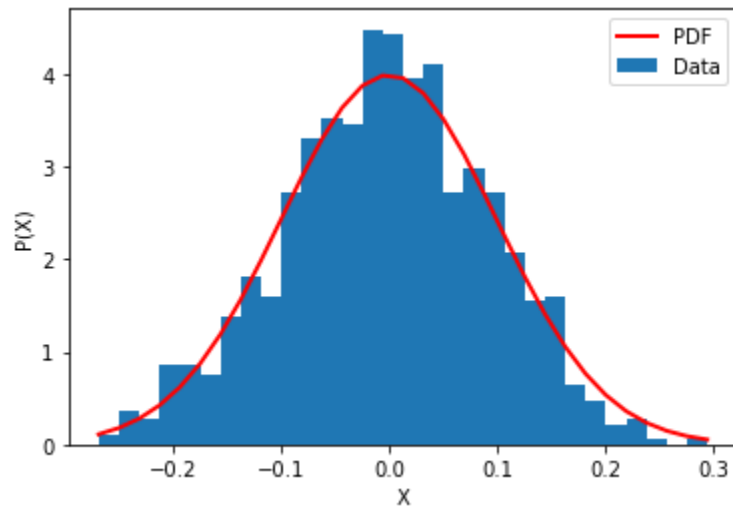
Other estimators are possible:

- *Maximum a posteriori (MAP)*
- *Minimum mean squared error (MMSE)*
- *Etc.*

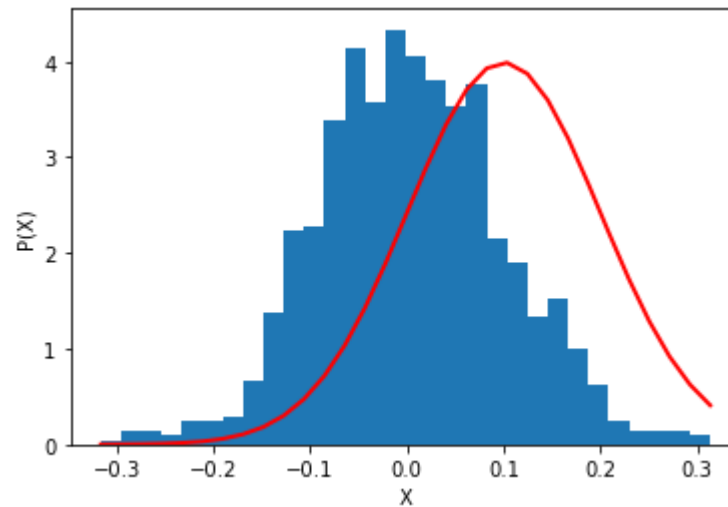
# Likelihood (Intuitively)

*Suppose we observe  $N$  data points from a Gaussian model and wish to estimate model parameters...*

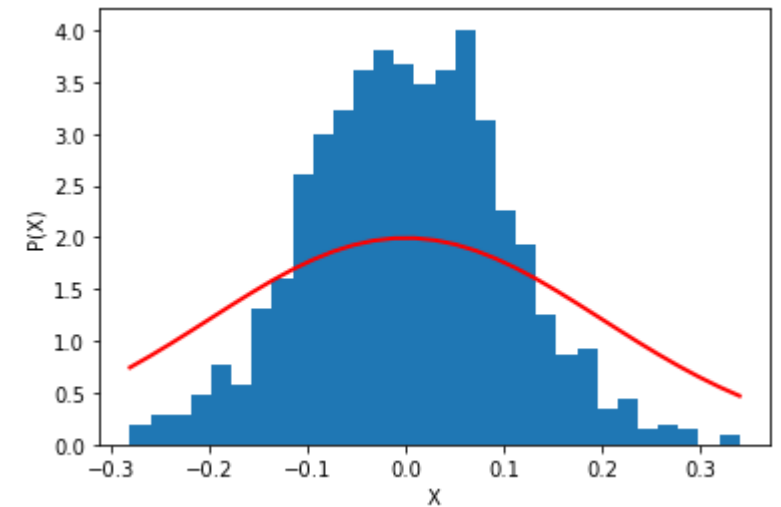
**High Likelihood**



**Low Likelihood (mean)**



**Low Likelihood (variance)**



**Likelihood Principle** *Given a statistical model, the likelihood function describes all evidence of a parameter that is contained in the data.*

# Likelihood Function

Suppose  $x_i \sim p(x; \theta)$ , then what is the **joint probability** over  $N$  *independent identically distributed* (iid) observations  $x_1, \dots, x_N$ ?

$$p(x_1, \dots, x_N; \theta) = \prod_{i=1}^N p(x_i; \theta)$$

- We call this the **likelihood function**, often denoted  $\mathcal{L}_N(\theta)$
- It is a function of the parameter  $\theta$ , the data are fixed
- Measures how well parameter  $\theta$  describes data (*goodness of fit*)

*How could we use this to estimate a parameter  $\theta$ ?*

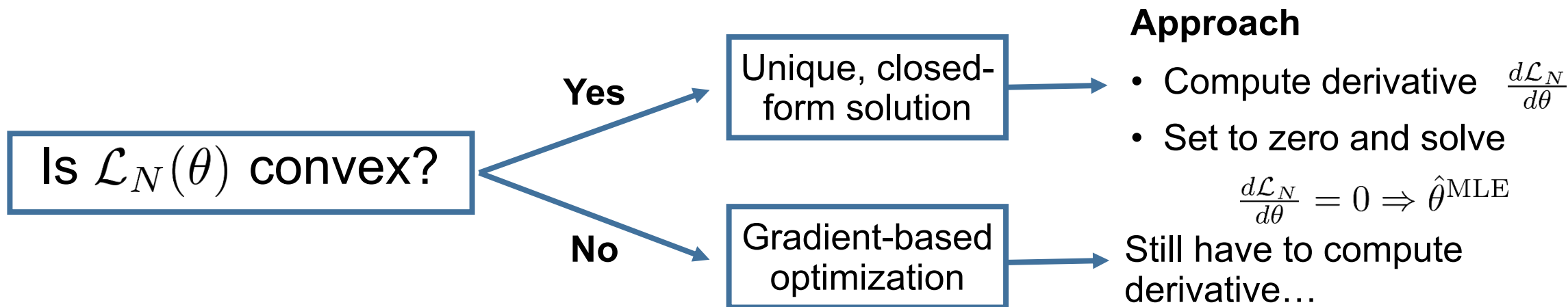
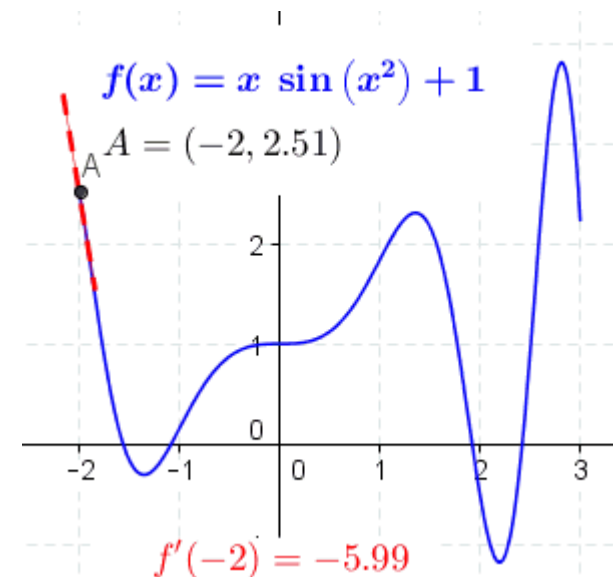
# Maximum Likelihood

**Maximum Likelihood Estimator (MLE)** as the name suggests, maximizes the likelihood function.

$$\hat{\theta}^{\text{MLE}} = \arg \max_{\theta} \mathcal{L}_N(\theta) = \prod_{i=1}^N p(x_i; \theta)$$

**Question** How do we find the MLE?

**Answer** Remember calculus...





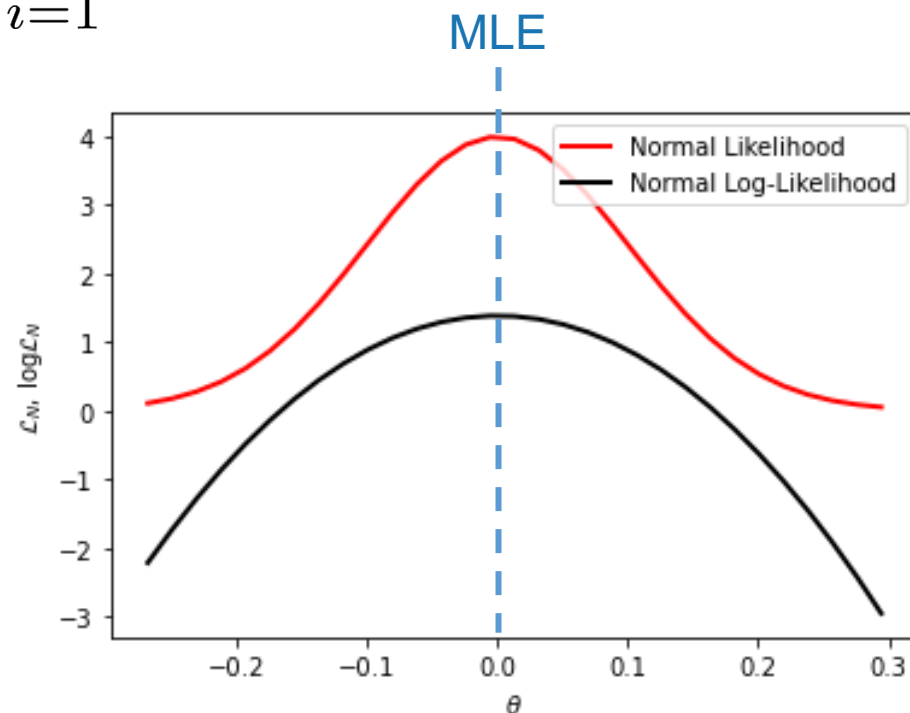
# Maximum Likelihood

Maximizing log-likelihood makes the math easier (as we will see) and doesn't change the answer (logarithm is an increasing function)

$$\hat{\theta}^{\text{MLE}} = \arg \max_{\theta} \log \mathcal{L}_N(\theta) = \sum_{i=1}^N \log p(x_i; \theta)$$

Derivative is a linear operator so,

$$\frac{d}{d\theta} \log \mathcal{L}_N(\theta) = \underbrace{\sum_{i=1}^N \frac{d}{d\theta} \log p(x_i; \theta)}_{\substack{\text{One term per data point} \\ \text{Can be computed in parallel} \\ \text{(big data)}}$$



# Maximum Likelihood

**Example** Suppose we have  $N$  coin tosses with  $X_1, \dots, X_n \sim \text{Bernoulli}(p)$  but we don't know the coin bias  $p$ . The likelihood function is,

$$\mathcal{L}_n(p) = \prod_{i=1}^n p^{x_i} (1-p)^{1-x_i} = p^S (1-p)^{n-S}$$

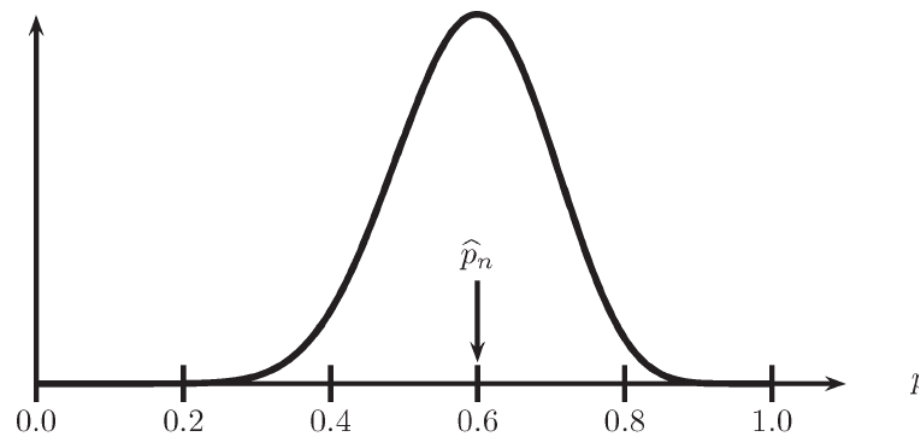
where  $S = \sum_i x_i$ . The log-likelihood is,

$$\log \mathcal{L}_n(p) = S \log p + (n - S) \log(1 - p)$$

Set the derivative of  $\log \mathcal{L}_n(p)$  to zero and solve,

$$\hat{p}^{\text{MLE}} = S/n = \frac{1}{n} \sum_{i=1}^n x_i$$

[ Source: Wasserman, L. 2004 ]



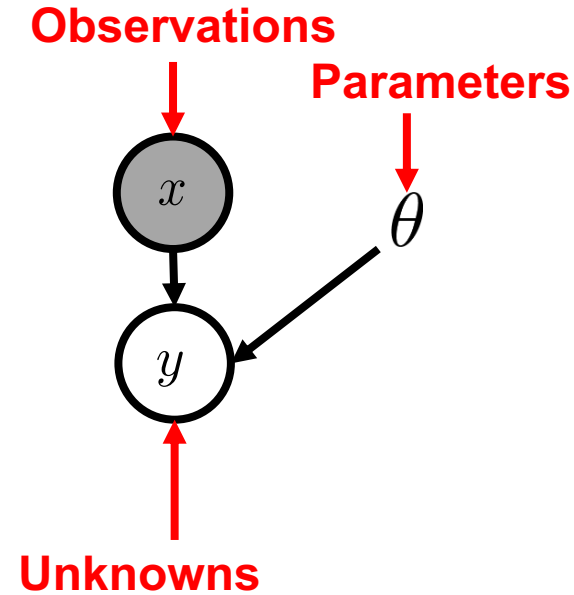
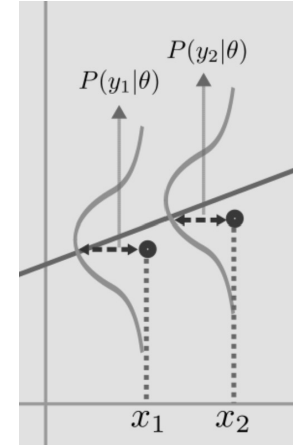
*Likelihood function for Bernoulli with  $n=20$  and  $\sum_i x_i = 12$  heads*

Maximum likelihood is equivalent to sample mean in Bernoulli

# Discriminative vs Generative modeling

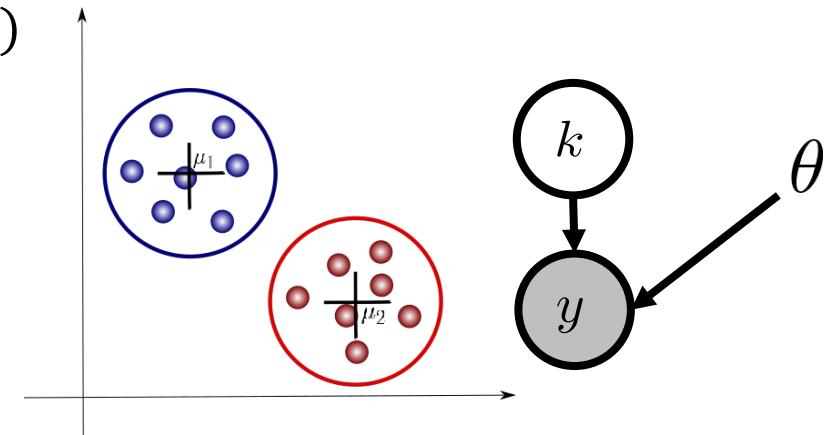
## Discriminative model:

- Only models  $P(y | x, \theta)$  -- i.e. *doesn't model data  $x$*
- Recall linear regression:  $y | x; \theta \sim N(x^\top \theta, \sigma^2)$
- Logistic regression:  $y | x; \theta \sim \text{Bernoulli}(\sigma(x^\top \theta))$



## Generative model:

- Models everything including data:  $P(k, y) = P(k)P(y | k, \theta)$
- e.g., Gaussian mixture model (GMM)
  - $\theta = (\pi_k, \mu_k, \Sigma_k)_{k=1}^K$
  - $k \sim \text{Categorical}(\pi)$  (*hidden*), i.e.  $P(k = l) = \pi_l$
  - $y | k \sim N(\mu_k, \Sigma_k)$



# Barbershop Example

Suppose you go to a barbershop at every last Friday of the month. You want to be able to predict the waiting time. You have collected 12 data points (i.e., how long it took to be served) from the last year:  $S = \{x_1, \dots, x_{12}\}$

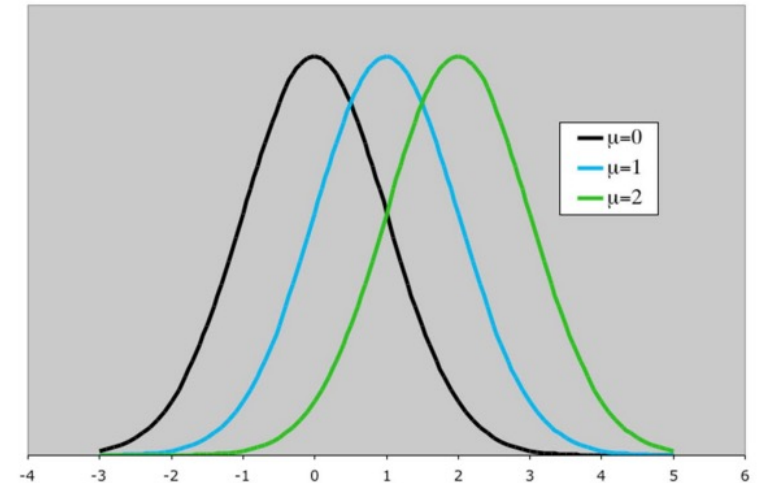
- 1. Modeling assumption:  $x_i \sim$  Gaussian distribution  $N(\mu, 1)$

- $p(x; \mu) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{(x-\mu)^2}{2}\right)$
- Observation: this distribution has mean  $\mu$

- 2. Find the MLE  $\hat{\mu}$  from data S

- (2.1) write down the neg. log likelihood of the sample

$$L_n(\mu) = -\ln P(x_1, \dots, x_n; \mu) = 12 \ln \sqrt{2\pi} + \frac{1}{2} \sum_{i=1}^{12} (x_i - \mu)^2$$



*Is this a generative or discriminative model?*

# Generative model: basic example I (cont'd)

## 2. Find the MLE $\hat{\mu}$ from data $S$

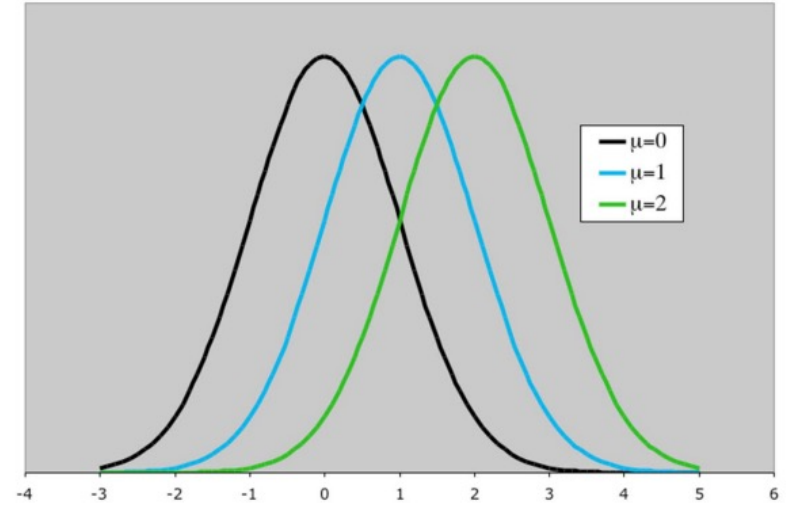
- (2.2) compute the first derivative, set it to 0, solve for  $\lambda$  (be sure to check convexity)

$$L'_n(\mu) = \sum_{i=1}^{12} (x_i - \mu) = 0 \Rightarrow \mu = \frac{x_1 + \dots + x_{12}}{12}$$

Sample Mean

## 3. The learned model $N(\hat{\mu}, 1)$ is yours!

- Simple prediction: e.g., predict the next wait time by  $\mathbb{E}_{X \sim N(\hat{\mu}, 1)}[X]$
- which is  $\hat{\mu} = \frac{x_1 + \dots + x_{12}}{12}$



## 4. (Optional: Model Checking) Generate some data... Does it look realistic?

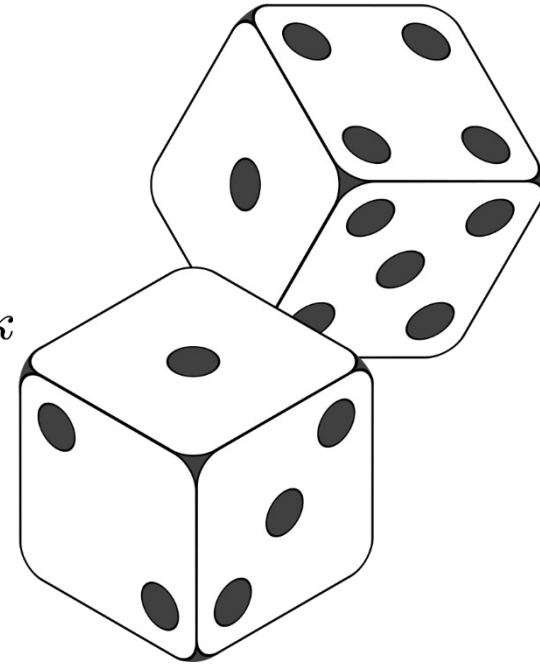
# (Aside) Categorical Distribution

*Distribution on integer-valued RV  $X \in \{1, \dots, K\}$*

$$p(X) = \prod_{k=1}^K \pi_k^{\mathbf{I}(X=k)} \quad \text{or} \quad p(X) = \sum_{k=1}^K \mathbf{I}(X = k) \cdot \pi_k$$

*with parameter  $p(X = k) = \pi_k$  and Kroenecker delta:*

$$\mathbf{I}(X = k) = \begin{cases} 1, & \text{If } X = k \\ 0, & \text{Otherwise} \end{cases}$$



Can also represent  $X$  as *one-hot* binary vector,

$$X \in \{0, 1\}^K \quad \text{where} \quad \sum_{k=1}^K X_k = 1 \quad \text{then} \quad p(X) = \prod_{k=1}^K \pi_k^{X_k}$$

# Basic Example II



**Data**  $S = \{y_i\}_{i=1}^n$ , where  $y_i \in \{1, \dots, C\}$

## Generative Story

$y \sim \text{Categorical}(\pi)$ , where  $\pi = (\pi_1, \dots, \pi_C) \in \Delta^{C-1}$  ( $\pi_c \geq 0$  and  $\pi_1 + \dots + \pi_C = 1$ )

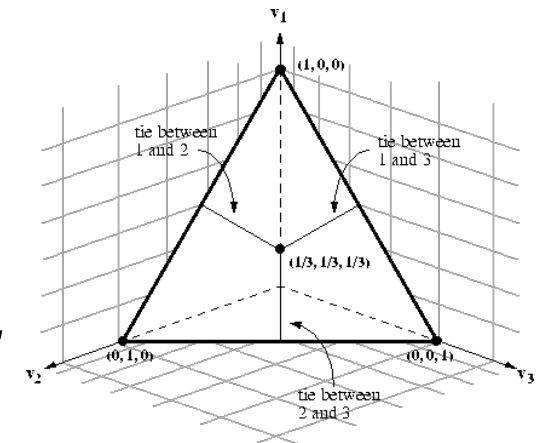
e.g.  $y_i$  = the color of  $i$ -th ball drawn randomly from a bin (with replacement)

$$p(y; \pi) = \pi_y \left( = \prod_{c=1}^C \pi_c^{I(y=c)} \right)$$

## Training

$$(2.1) L_n(\pi) = -\ln P(y_1, \dots, y_n; \pi) = \sum_{i=1}^n -\ln \pi_{y_i} = -\sum_{c=1}^C n_c \ln \pi_c,$$

where  $n_c = \#\{i: y_i = c\} = \sum_{i=1}^n I(y_i = c)$



# Basic Example II (Cont'd)



## Training

$$(2.2) \text{ minimize}_{\pi \in \Delta^{C-1}} L_n(\pi) := - \sum_{c=1}^C n_c \ln \pi_c$$

Constrained maximization problem; solve by Lagrange multipliers

$$\frac{\partial}{\partial \pi} \left( - \sum_{c=1}^C n_c \ln \pi_c - \lambda \left( \sum_{c=1}^C \pi_c - 1 \right) \right) = - \frac{n_c}{\pi_c} - \lambda = 0 \Rightarrow \pi_c = - \frac{n_c}{\lambda}$$

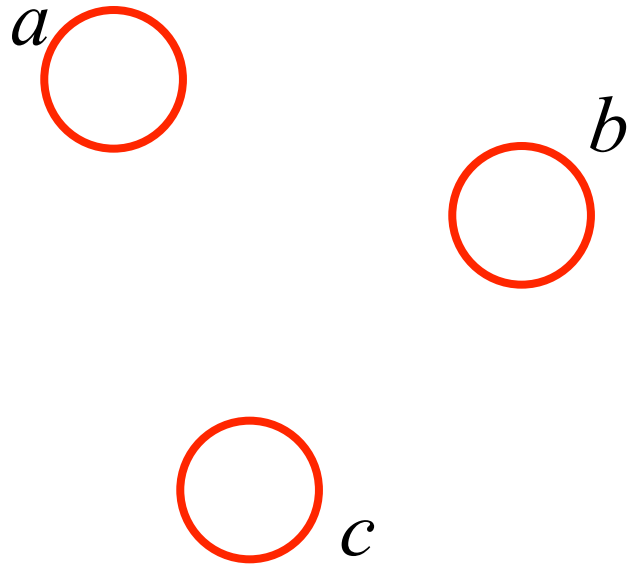
Combined with the constraint that  $\pi_1 + \dots + \pi_C = 1 \Rightarrow \hat{\pi}_c = \frac{n_c}{n}$ , for all  $c$

**Test** predict label  $\operatorname{argmax}_c P(y = c; \hat{\pi}) = \operatorname{argmax}_c \hat{\pi}_c$

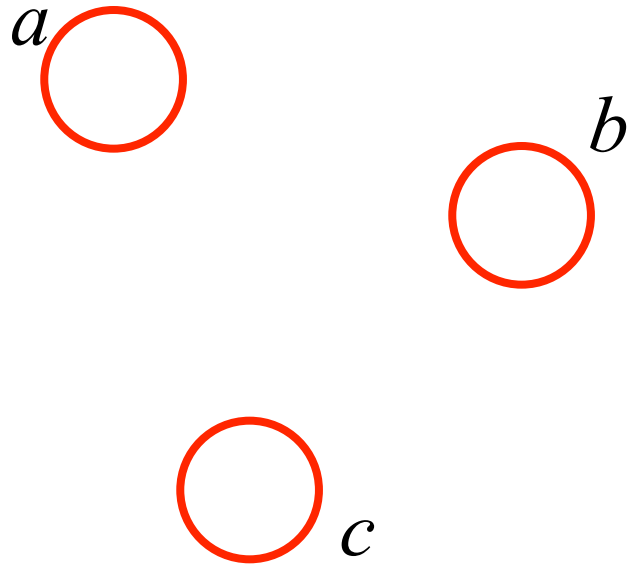


- Probability Refresher
- Probabilistic Graphical Models
- **Naïve Bayes**

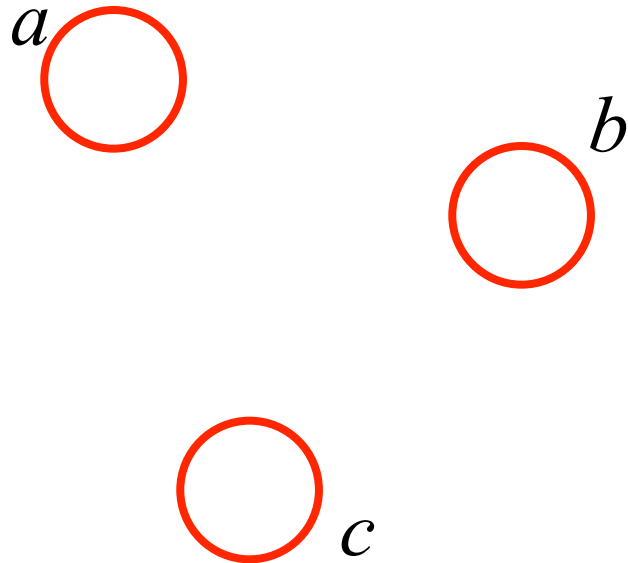
*What is the joint factorization?*



$$p(\mathbf{a}, \mathbf{b}, \mathbf{c}) = p(\mathbf{a})p(\mathbf{b})p(\mathbf{c})$$

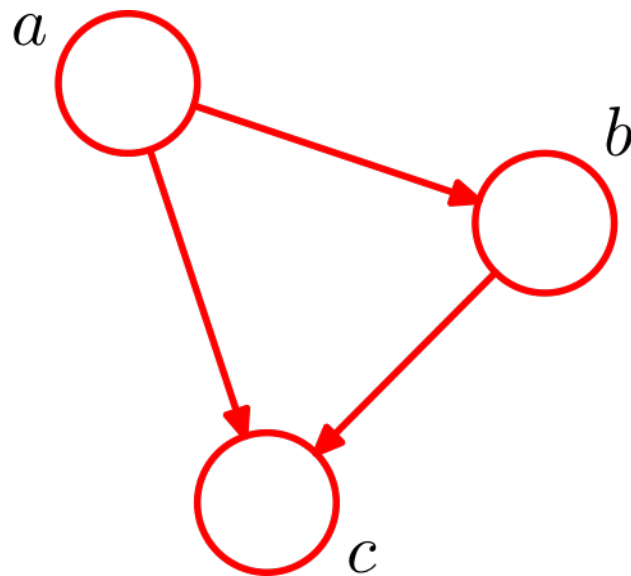


*Are  $a$  and  $b$  independent ( $a \perp b$ )?*



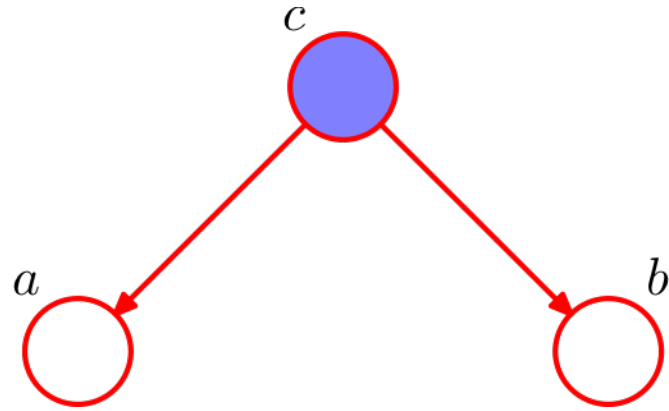
$$\mathbf{p(a,b,c) = p(a)p(b)p(c)}$$

$$p(a,b,c) = p(a)p(b|a)p(c|a,b)$$



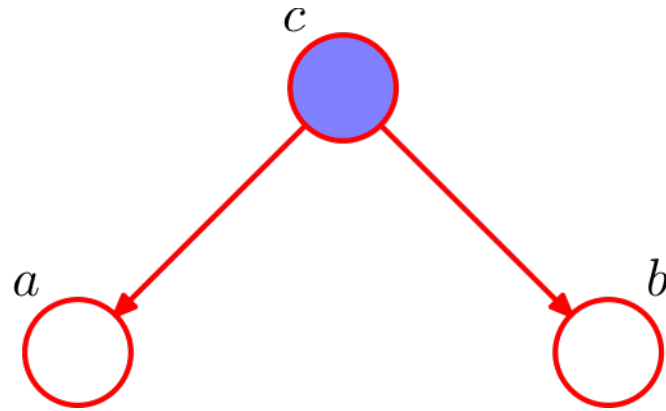
Note there are **no conditional independencies**

# Case one where c is observed



Is  $a \perp b \mid c$  ?

# Case one where c is observed



$$a \perp b \mid c$$

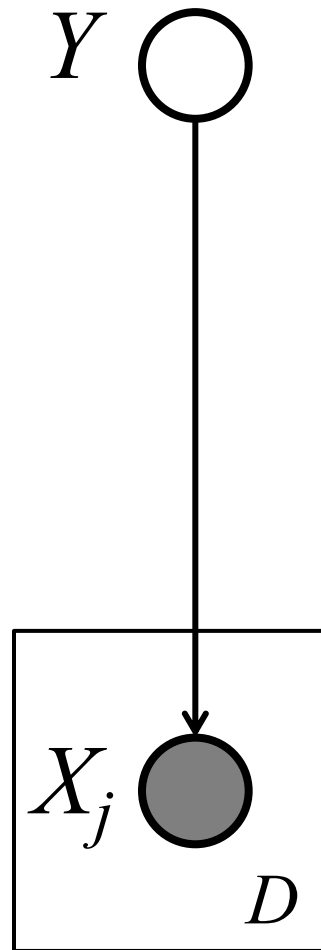
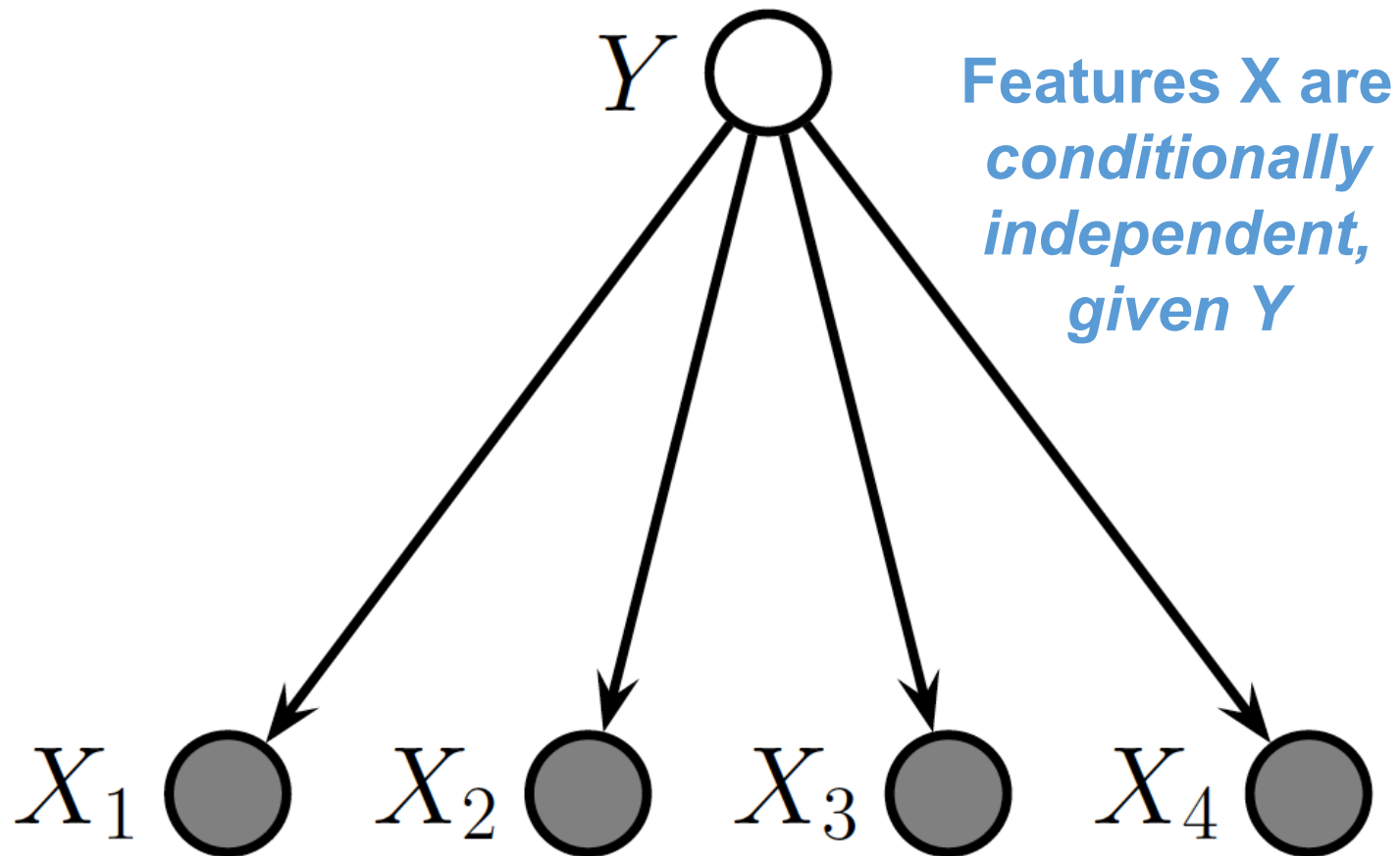
$$p(a, b, c) = p(c)p(a|c)p(b|c) \quad (\text{what the graph represents in general})$$

$$p(a, b|c) = p(a|c)p(b|c) \quad (\text{with } c \text{ observed})$$

This is the definition of  $a \perp b \mid c$

# Shading & Plate Notation

*Convention: Shaded nodes are observed, open nodes are latent/hidden/unobserved*



*Plates denote replication of random variables*

$$p(y, \mathbf{x}) = p(y) \prod_{j=1}^D p(x_j | y)$$



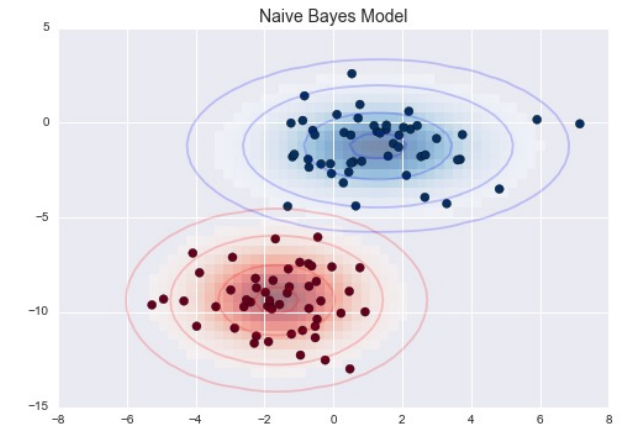
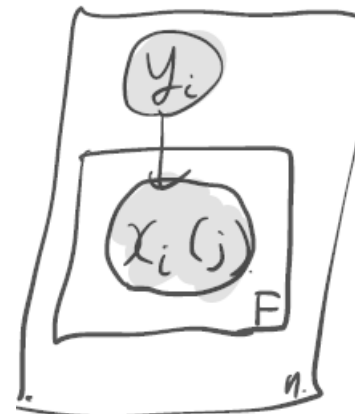
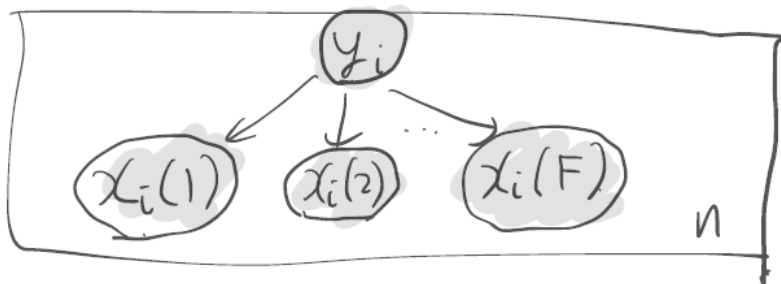
# Naïve Bayes for supervised learning

- Motivation: supervised learning for classification
- high-dimensional  $x = (x(1), \dots, x(F))$ , modeling  $P(x | y)$  can be tricky
- In general,  $P(x | y) = P(x(1) | y) \cdot P(x(2) | x(1), y) \cdot \dots \cdot P(x(F) | x(1), \dots, x(F - 1), y)$
- A modeling assumption:  $x(1), \dots, x(F)$  are conditionally independent given  $y$   
i.e. for all  $i$

$$x(i) \perp\!\!\!\perp (x(1), \dots, x(i - 1), x(i + 1), \dots, x(F)) | y$$

(Conditional independence notation:  $A \perp\!\!\!\perp B | C$ )

- Equivalently  $P(x | y) = P(x(1) | y) \cdot \dots \cdot P(x(F) | y)$



# Recall : Class Preference Prediction

Define the labeled training dataset  $S = \{(x_i, y_i)\}_{i=1}^m$

To make this a binary classification we set  
“Liked” =  $\{+2, +1, 0\}$   
“Nah” =  $\{-1, -2\}$

Features	Rating	Easy?	AI?	Sys?	Thy?	Morning?
	+2	y	y	n	y	n
Feature Values	+2	y	y	n	y	n
	+2	n	y	n	n	n
	+2	n	n	n	y	n
	+2	n	y	y	n	y
	+1	y	y	n	n	n
	+1	y	y	n	y	n
	+1	n	y	n	y	n
Labels	0	n	n	n	n	y
	0	y	n	n	y	y
	0	n	y	n	y	n
	0	y	y	y	y	y
	-1	y	y	y	n	y
	-1	n	n	y	y	n
	-1	n	n	y	n	y
	-1	y	n	y	n	y
	-2	n	n	y	y	n
	-2	n	y	y	n	y
Data Point	-2	y	n	y	n	n
	-2	y	n	y	n	y

# Naïve Bayes: binary-valued features

**Training Data**  $S = \{(x_i, y_i)\}_{i=1}^n$ ,

$$x_i \in \{0,1\}^F$$

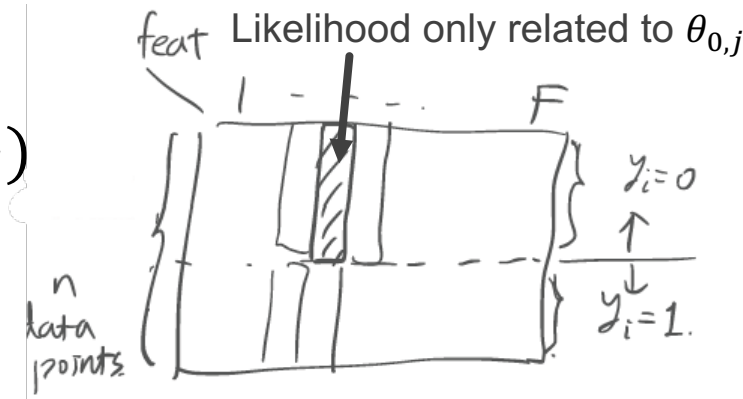
$$y_i \in \{0,1\}$$

## Generative Story

$y \sim \text{Bernoulli}(\pi)$ ; for all  $j \in [F]$ ,  $x(j) \mid y = c \sim \text{Bernoulli}(\theta_{c,j})$

#parameters =  $1 + 2F$

**Training** (denote by  $\theta = \{\theta_{c,j}\}$ )



$$\begin{aligned} \max_{\pi, \theta} \sum_{i=1}^n \ln P(x_i, y_i; \pi, \theta) &= \sum_{i=1}^n \ln P(y_i; \pi) + \sum_{i=1}^n \ln P(x_i \mid y_i; \theta) \\ &= \max_{\pi} \sum_{i=1}^n \ln P(y_i; \pi) + \max_{\{\theta_{0,j}\}} \sum_{i:y_i=0} \ln P(x_i \mid y_i; \theta) + \max_{\{\theta_{1,j}\}} \sum_{i:y_i=1} \ln P(x_i \mid y_i; \theta) \end{aligned}$$

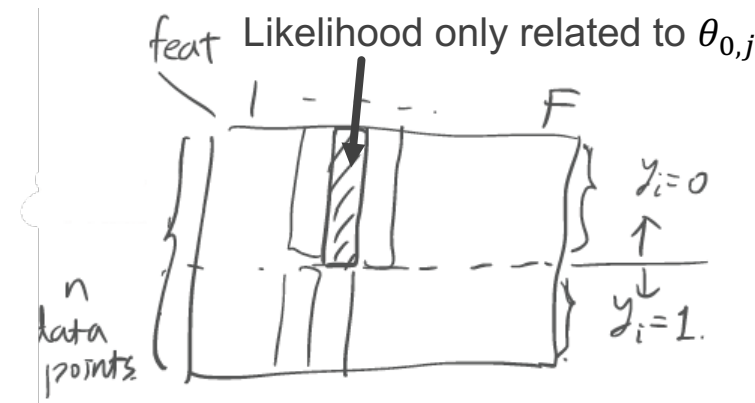
Key observation: optimal  $\pi$ , optimal  $\{\theta_{0,j}\}$ , optimal  $\{\theta_{1,j}\}$  can be found separately

$$\text{Optimal } \pi: \max_{\pi} \sum_{i=1}^n \ln P(y_i; \pi) = \max_{\pi} n_0 \ln(1 - \pi) + n_1 \ln(\pi) \Rightarrow \hat{\pi} = \frac{n_1}{n}$$

# Naïve Bayes: binary-valued features (cont'd)

By the Naïve Bayes modeling assumption,

$$\begin{aligned}\max_{\{\theta_{0,j}\}} \sum_{i:y_i=0} \ln P(x_i | y_i; \theta) &= \max_{\{\theta_{0,j}\}} \sum_{j=1}^F \sum_{i:y_i=0} \ln P(x_i(j) | y_i; \theta_{0,j}) \\ &= \sum_{j=1}^F \max_{\theta_{0,j}} \sum_{i:y_i=0} \ln P(x_i(j) | y_i; \theta_{0,j})\end{aligned}$$



Again, can optimize each  $\theta_{0,j}$  separately,

- Optimal  $\theta_{0,j}$ :  $\max_{\theta_{0,j}} \sum_{i:y_i=0, x_i(j)=1} \ln \theta_{0,j} + \sum_{i:y_i=0, x_i(j)=0} \ln (1 - \theta_{0,j})$

$$\hat{\theta}_{0,j} = \frac{\#\{i: y_i=0, x_i(j)=1\}}{\#\{i: y_i=0\}}; \quad \text{similarly,} \quad \hat{\theta}_{1,j} = \frac{\#\{i: y_i=1, x_i(j)=1\}}{\#\{i: y_i=1\}}$$

# Naïve Bayes: binary-valued features (cont'd)

**Test** Given  $\hat{\pi}$ ,  $\{\hat{\theta}_{c,j}\}$ , Bayes optimal classifier

$$\hat{f}_{BO}(x) = \operatorname{argmax}_y P(x, y; \hat{\pi}, \{\hat{\theta}_{c,j}\}) = \operatorname{argmax}_y \log P(x, y; \hat{\pi}, \{\hat{\theta}_{c,j}\})$$

- $\log P(x, y = 0; \pi, \{\theta_{c,j}\}) = \ln(1 - \pi) + \sum_{j=1}^F \ln P(x(j) | y; \theta_{0,j})$   
 $= \ln(1 - \pi) + \sum_{j=1}^F \ln(1 - \theta_{0,j}) I(x(j) = 0) + \ln(\theta_{0,j}) I(x(j) = 1)$   
 $= \ln(1 - \pi) + \sum_{j=1}^F \ln(1 - \theta_{0,j}) + \sum_{j=1}^F x(j) \ln \frac{\theta_{0,j}}{1 - \theta_{0,j}}$
- Similarly,  $\log P(x, y = 1; \pi, \{\theta_{c,j}\}) = \ln(\pi) + \sum_{j=1}^F \ln(1 - \theta_{1,j}) + \sum_{j=1}^F x(j) \ln \frac{\theta_{1,j}}{1 - \theta_{1,j}}$
- Therefore,  $\hat{f}_{BO}(x) = 1 \Leftrightarrow \underbrace{\ln\left(\frac{\pi}{1 - \pi}\right) + \sum_{j=1}^F \ln\left(\frac{1 - \theta_{1,j}}{1 - \theta_{0,j}}\right)}_b + \sum_{j=1}^F x(j) \underbrace{\left(\ln \frac{\theta_{1,j}}{1 - \theta_{1,j}} - \ln \frac{\theta_{0,j}}{1 - \theta_{0,j}}\right)}_{w(j)} \geq 0$
- I.e. Bayes classifier is *linear*

# Naïve Bayes: Discrete (Categorical-valued) features

Data  $S = \{(x_i, y_i)\}_{i=1}^n$ ,

$x_i \in [W]^F$

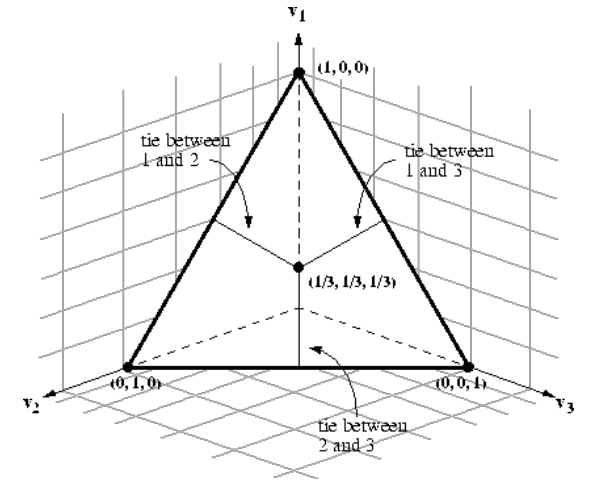
$y_i \in \{0,1\}$

## Generative story

$y \sim \text{Bernoulli}(\pi)$ ; for all  $j \in [F]$ ,  $x(j) \mid y = c \sim \text{Categorical}(\theta_c)$  ( $\theta_c \in \Delta^{W-1}$ )

#parameters =  $1 + 2W$

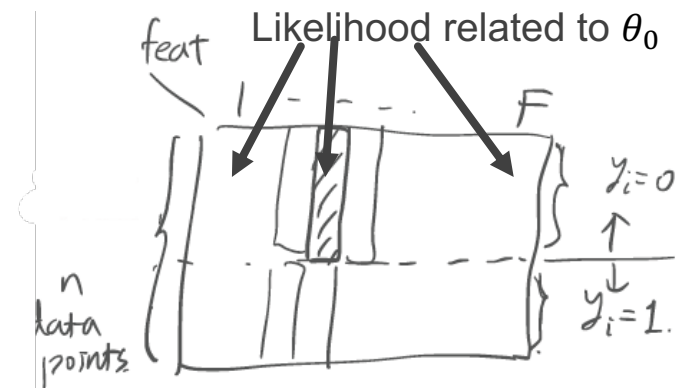
Note: in this example,  $\theta_c$  shared across all features!



## Training

Similar to previous example, optimal  $\pi$ , optimal  $\theta_0$ , optimal  $\theta_1$  can be found separately, by maximizing the respective part of the likelihood function (exercise)

Optimal  $\pi$  same as previous example



# Naïve Bayes: Discrete features (cont'd)

## Training

Optimal  $\theta_c$ :

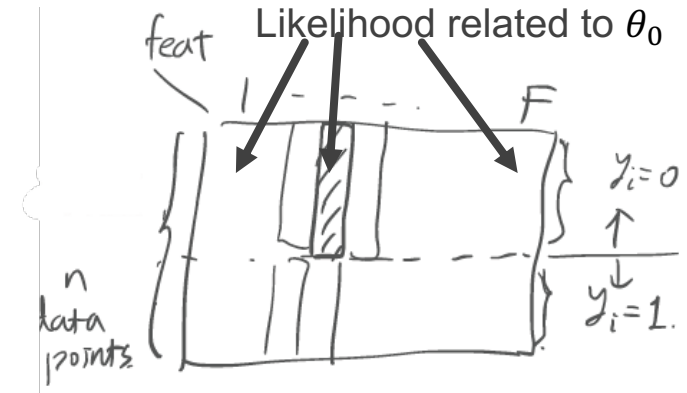
$$\begin{aligned}\max_{\theta_0} \sum_{i: y_i=0} \ln P(x_i | y_i; \theta_0) &= \max_{\theta_0} \sum_{j=1}^F \sum_{i: y_i=0} \ln P(x_i(j) | y_i; \theta_0) \\ &= \max_{\theta_0} \sum_{w=1}^W \sum_{j=1}^F \sum_{i: y_i=0} I(x_i(j) = w) \ln \theta_{0,w} \\ &= \max_{\theta_0} \sum_{w=1}^W \ln \theta_{0,w} \#\{(i, j): y_i = 0, x_i(j) = w\}\end{aligned}$$

$$\Rightarrow \hat{\theta}_{c,w} = \frac{\#\{(i, j): y_i=c, x_i(j)=w\}}{\#\{i: y_i=c\} \times F}$$

Exercise: how to extend this to variable-length  $x_i$ 's (e.g. for text classification)?

## Test

Bayes optimal classification rule with  $(\hat{\pi}, \hat{\theta}_0, \hat{\theta}_1)$  (exercise)



# Summary

## Fundamental rules of Probability:

- Law of total probability:  $p(Y) = \sum_x p(Y, X = x)$
- Probability chain rule:  $p(X | Y) = \frac{p(X, Y)}{p(Y)}$
- Conditional probability:  $p(X, Y) = p(Y)p(X | Y)$

## Independence of Random Variables:

- Two RVs are independent if:  $p(X = x, Y = y) = p(X = x)p(Y = y)$
- Or:  $p(X | Y) = p(X)$
- They are *conditionally independent* if:

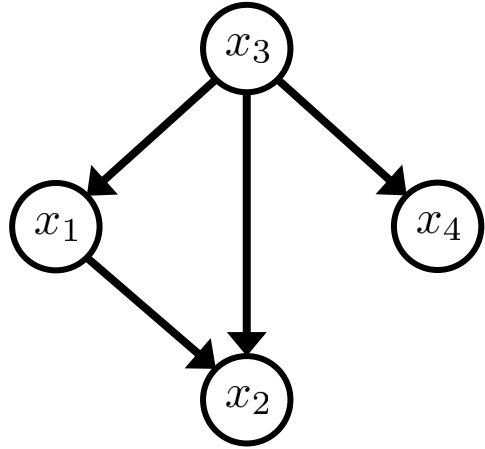
$$p(X = x, Y = y | Z = z) = p(X = x | Z = z)p(Y = y | Z = z)$$

- Or:  $p(X | Y, Z) = p(X | Z)$



# Summary

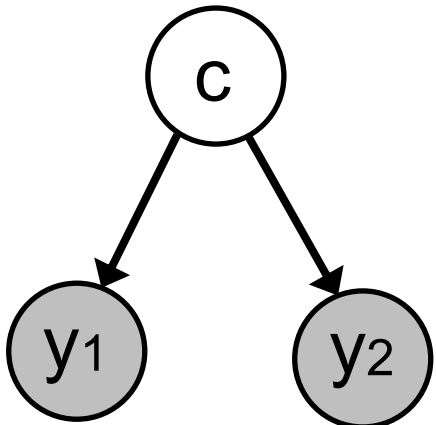
A Bayes Network expresses a unique probability factorization:



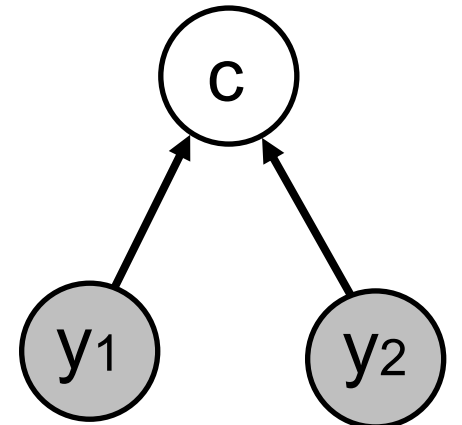
$$p(x) = \prod_{s \in \mathcal{V}} p(x_s \mid x_{\text{Pa}(s)})$$

 Parents of node  $s$

Inference is performed by Bayes' rule (posterior distribution):

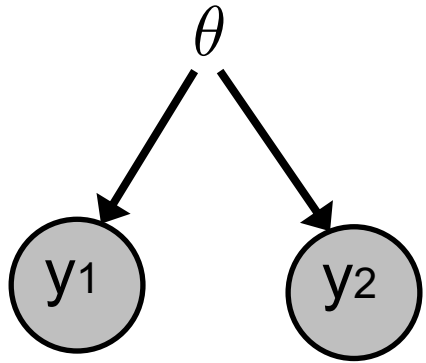


$$p(c \mid y_1, y_2) = \frac{p(c)p(y_1 \mid c)p(y_2 \mid c)}{p(y_1, y_2)}$$



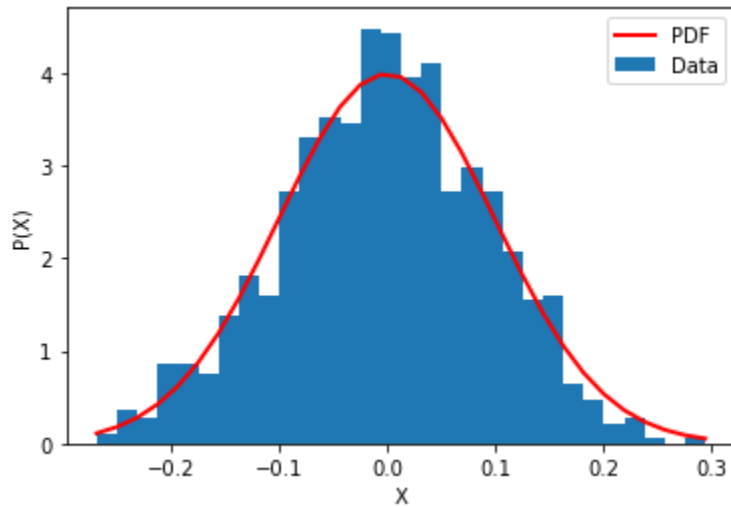
# Summary

Hyperparameters must be estimated (e.g. Maximum Likelihood):

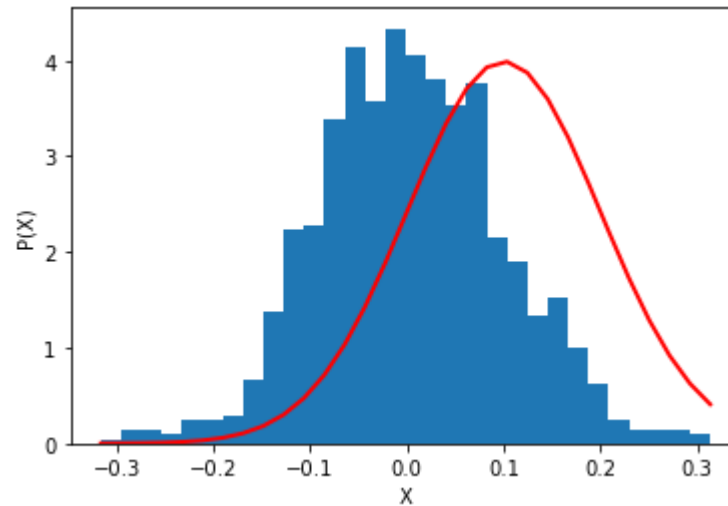


$$\hat{\theta}^{\text{MLE}} = \arg \max_{\theta} \log p(y_1, \dots, y_n \mid \theta)$$

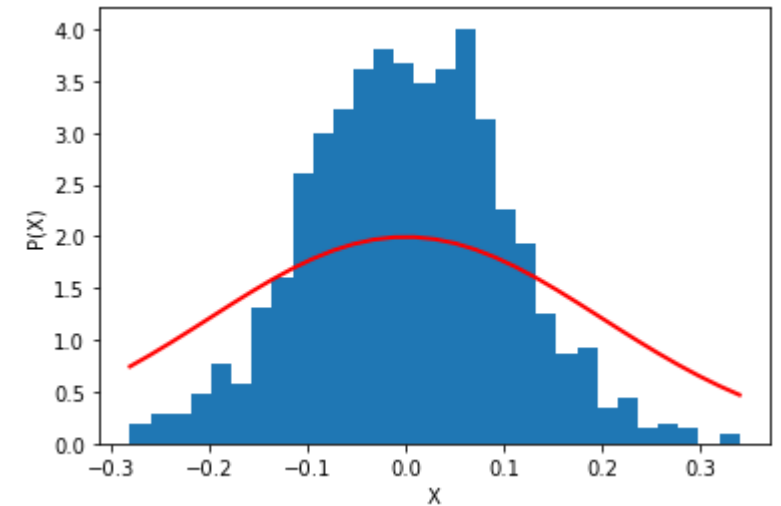
**High Likelihood**



**Low Likelihood (mean)**



**Low Likelihood (variance)**



# Summary

Naïve Bayes classifier assumes features are *conditionally independent* given class  $Y$ :

$$x(j) \perp\!\!\!\perp (x(1), \dots, x(j-1), x(j+1), \dots, x(D)) \mid y$$

Joint distribution factorizes as:

$$p(x, y) = p(y) \prod_{\{j=1\}}^D p(x(j) \mid y)$$

Allows easier fitting of hyperparameters for *class conditional distributions* (they can be fit independently of each other)

