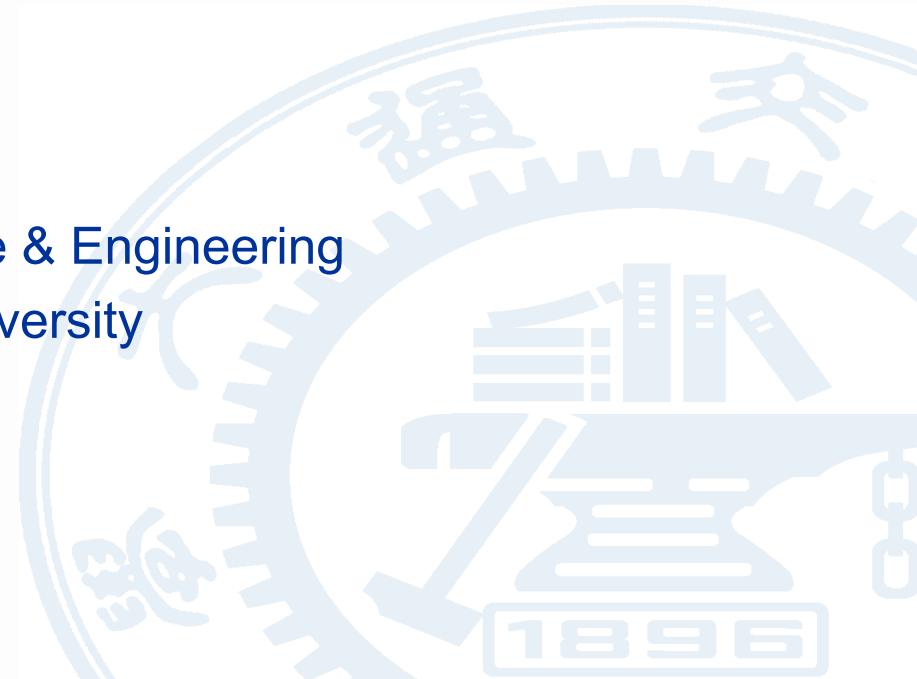




# Machine Learning

## Lecture &

Yang Yang  
Department of Computer Science & Engineering  
Shanghai Jiao Tong University



# Bayesian Learning



**Bayes, Thomas (1763)** An essay towards solving a problem in the doctrine of chances. *Philosophical Transactions of the Royal Society of London*, 53:370-418

# Your first consulting job

- A billionaire asks you a question:
  - He says: I have a coin, if I flip it, what's the probability it will fall with the head up?
  - You say: Please flip it a few times:



- You say: The probability is: **3/5**
- **He says: Why???**
- You say: Because...

# Bernoulli distribution

Data,  $D =$



$$D = \{X_i\}_{i=1}^n, X_i \in \{\text{H}, \text{T}\}$$

- $P(\text{Heads}) = \theta, P(\text{Tails}) = 1-\theta$
- Flips are **i.i.d.:**
  - **Independent** events
  - **Identically distributed** according to Bernoulli distribution

Choose  $\theta$  that maximizes the probability of observed data

# Maximum Likelihood Estimation

Choose  $\theta$  that maximizes the probability of observed data

$$\hat{\theta}_{MLE} = \arg \max_{\theta} P(D | \theta)$$

MLE of probability of head:

$$\hat{\theta}_{MLE} = \frac{\alpha_H}{\alpha_H + \alpha_T} = 3/5 \text{ "Frequency of heads"}$$

  
Number of heads    Number of tails

# Maximum Likelihood Estimation

Choose  $\theta$  that maximizes the probability of observed data

$$\begin{aligned}\hat{\theta}_{MLE} &= \arg \max_{\theta} P(D | \theta) \\ &= \arg \max_{\theta} \prod_{i=1}^n P(X_i | \theta) \quad \text{Independent draws} \\ &= \arg \max_{\theta} \prod_{i:X_i=H} \theta \prod_{i:X_i=T} (1 - \theta) \quad \text{Identically distributed} \\ &= \arg \max_{\theta} \underbrace{\theta^{\alpha_H} (1 - \theta)^{\alpha_T}}_{J(\theta)}\end{aligned}$$

# Maximum Likelihood Estimation

Choose  $\theta$  that maximizes the probability of observed data

$$\begin{aligned}\hat{\theta}_{MLE} &= \arg \max_{\theta} P(D | \theta) \\ &= \arg \max_{\theta} \underbrace{\theta^{\alpha_H} (1 - \theta)^{\alpha_T}}_{J(\theta)}\end{aligned}$$

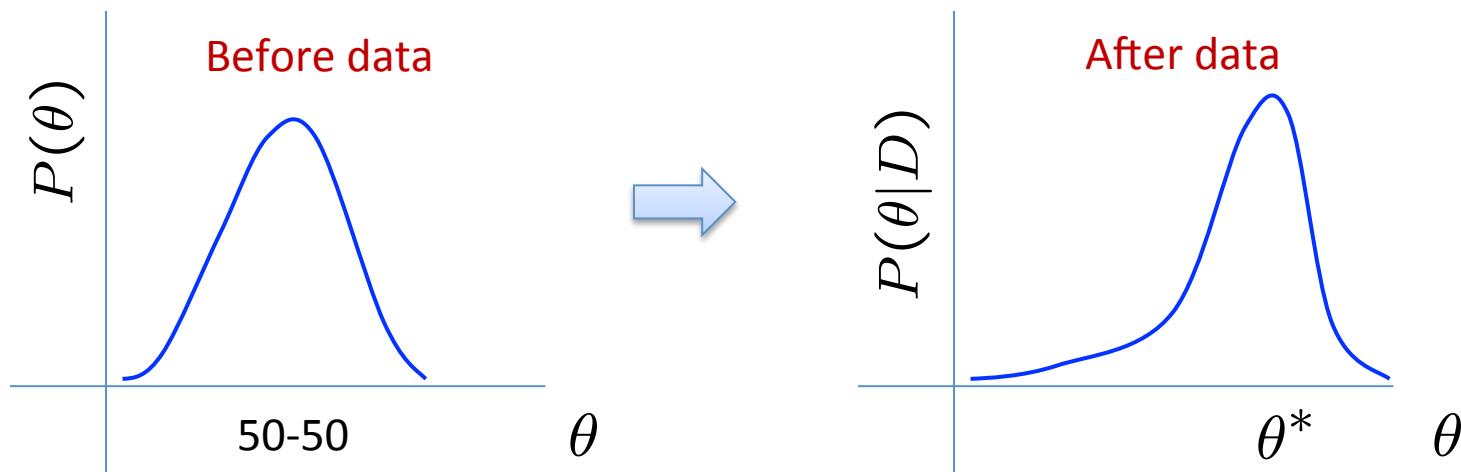
$$\frac{\partial J(\theta)}{\partial \theta} \Big|_{\theta=\hat{\theta}_{MLE}} = 0$$

$$\alpha_H(1 - \theta) - \alpha_T\theta \Big|_{\theta=\hat{\theta}_{MLE}} = 0$$

$$\boxed{\hat{\theta}_{MLE} = \frac{\alpha_H}{\alpha_H + \alpha_T}}$$

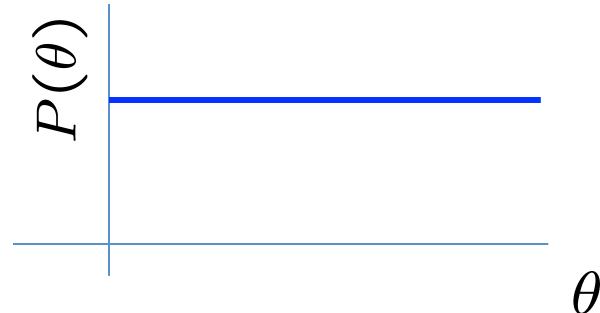
# What about prior knowledge?

- Billionaire says: Wait, I know that the coin is “close” to 50-50. What can you do for me now?
- You say: I can learn it the Bayesian way...
- Rather than estimating a single  $\theta$ , we obtain a distribution over possible values of  $\theta$



# Prior distribution

- What about prior?
  - Represents expert knowledge
  - Simple posterior form
- Uninformative priors:
  - Uniform distribution
- Conjugate priors:
  - Closed-form representation of posterior
  - $P(\theta)$  and  $P(\theta | D)$  have the same form



# Conjugate Prior

- $P(\theta)$  and  $P(\theta | D)$  have the same form

Eg. 1 Coin flip problem

Likelihood is  $\sim$  Binomial

$$P(\mathcal{D} | \theta) = \binom{n}{\alpha_H} \theta^{\alpha_H} (1 - \theta)^{\alpha_T}$$

If prior is Beta distribution,

$$P(\theta) = \frac{\theta^{\beta_H-1} (1 - \theta)^{\beta_T-1}}{B(\beta_H, \beta_T)} \sim Beta(\beta_H, \beta_T)$$

Then posterior is Beta distribution

$$P(\theta|D) \sim Beta(\beta_H + \alpha_H, \beta_T + \alpha_T)$$

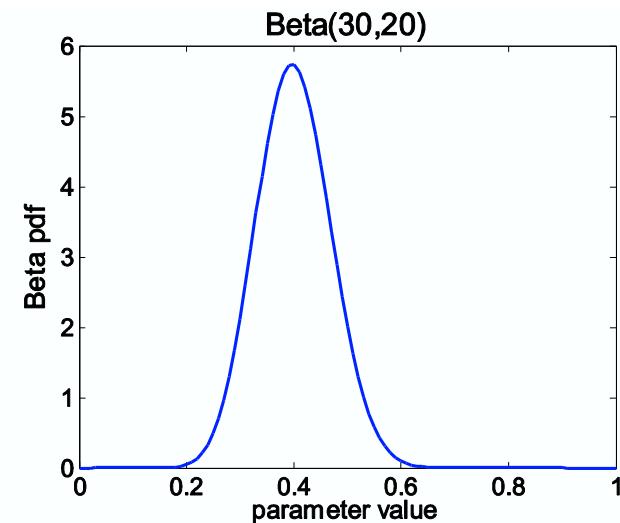
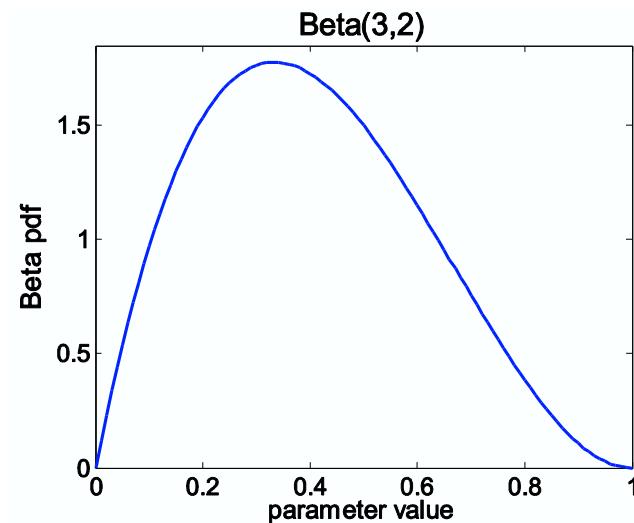
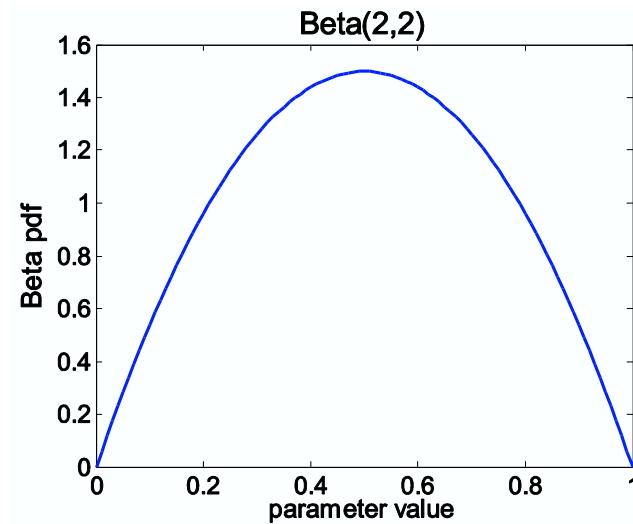
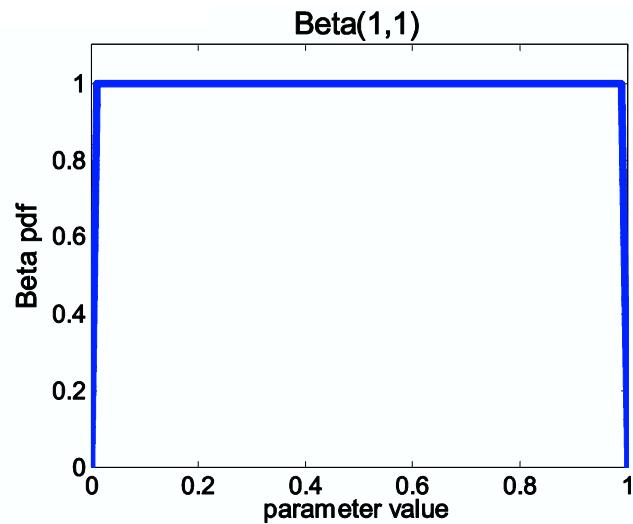


**For Binomial, conjugate prior is Beta distribution.**

# Beta distribution

$Beta(\beta_H, \beta_T)$

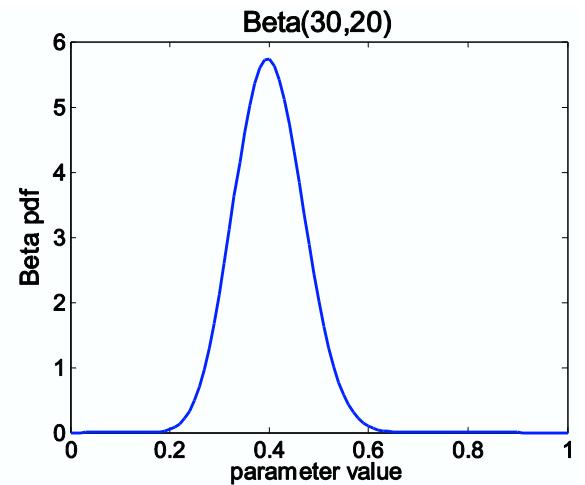
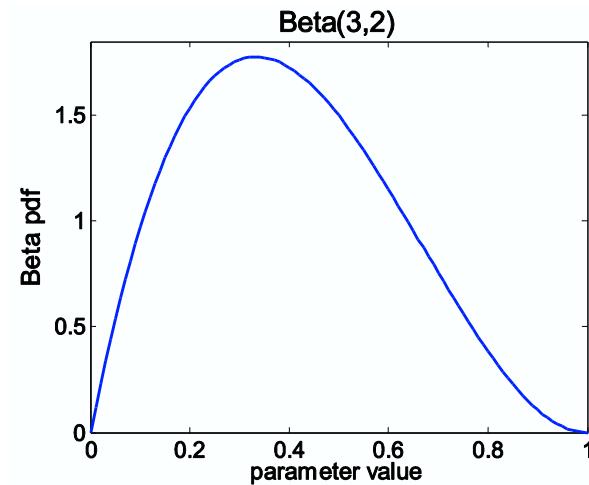
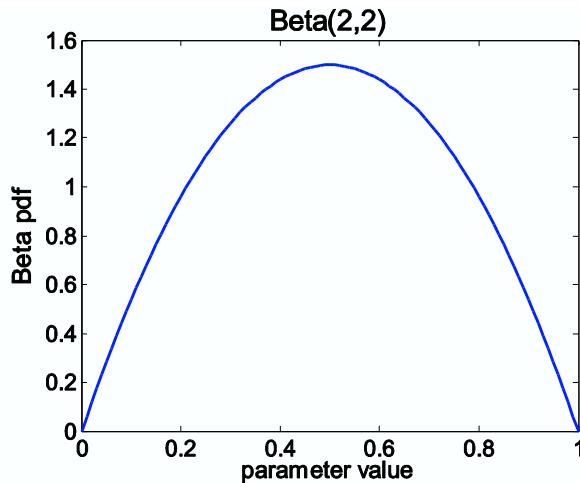
More concentrated as values of  $\beta_H, \beta_T$  increase



# Beta conjugate prior

$$P(\theta) \sim Beta(\beta_H, \beta_T)$$

$$P(\theta|D) \sim Beta(\beta_H + \alpha_H, \beta_T + \alpha_T)$$



As  $n = \alpha_H + \alpha_T$   
increases

As we get more samples, effect of prior is “washed out”

# Conjugate Prior

- $P(\theta)$  and  $P(\theta | D)$  have the same form

Eg. 2 Dice roll problem (6 outcomes instead of 2)



Likelihood is  $\sim \text{Multinomial}(\theta = \{\theta_1, \theta_2, \dots, \theta_k\})$

$$P(D | \theta) = \frac{n!}{\alpha_1! \alpha_2! \dots \alpha_k!} \theta_1^{\alpha_1} \theta_2^{\alpha_2} \dots \theta_k^{\alpha_k} \sum_{i=1}^k \alpha_i = n \quad \sum_{i=1}^k \theta_i = 1$$

If prior is Dirichlet distribution,

$$P(\theta) = \frac{\prod_{i=1}^k \theta_i^{\beta_i - 1}}{B(\beta_1, \dots, \beta_k)} \sim \text{Dirichlet}(\beta_1, \dots, \beta_k)$$

Then posterior is Dirichlet distribution

$$P(\theta | D) \sim \text{Dirichlet}(\beta_1 + \alpha_1, \dots, \beta_k + \alpha_k)$$

**For Multinomial, conjugate prior is Dirichlet distribution.**

# Maximum A Posteriori Estimation

Choose  $\theta$  that maximizes a posterior probability

$$\begin{aligned}\hat{\theta}_{MAP} &= \arg \max_{\theta} P(\theta | D) \\ &= \arg \max_{\theta} P(D | \theta)P(\theta)\end{aligned}$$

MAP estimate of probability of head:

$$P(\theta|D) \sim Beta(\beta_H + \alpha_H, \beta_T + \alpha_T)$$

$$\hat{\theta}_{MAP} = \frac{\alpha_H + \beta_H - 1}{\alpha_H + \beta_H + \alpha_T + \beta_T - 2} \quad \text{Mode of Beta distribution}$$

# MLE vs. MAP

- Maximum Likelihood estimation (MLE)  
Choose value that maximizes the probability of observed data

$$\hat{\theta}_{MLE} = \arg \max_{\theta} P(D|\theta)$$

- Maximum *a posteriori* (MAP) estimation  
Choose value that is most probable given observed data and prior belief

$$\begin{aligned}\hat{\theta}_{MAP} &= \arg \max_{\theta} P(\theta|D) \\ &= \arg \max_{\theta} P(D|\theta)P(\theta)\end{aligned}$$

When is MAP same as MLE?

# MLE vs. MAP

$$\hat{\theta}_{MLE} = \frac{\alpha_H}{\alpha_H + \alpha_T}$$

What if we toss the coin too few times?



- You say: Probability next toss is a head = 0
- Billionaire says: You're fired!

$$\hat{\theta}_{MAP} = \frac{\alpha_H + \beta_H - 1}{\alpha_H + \beta_H + \alpha_T + \beta_T - 2}$$

- Beta prior equivalent to extra coin flips
- As  $n \rightarrow \infty$ , prior is “forgotten”
- **But, for small sample size, prior is important!**

# Bayesians vs. Frequentists

You are no  
good when  
sample is  
small

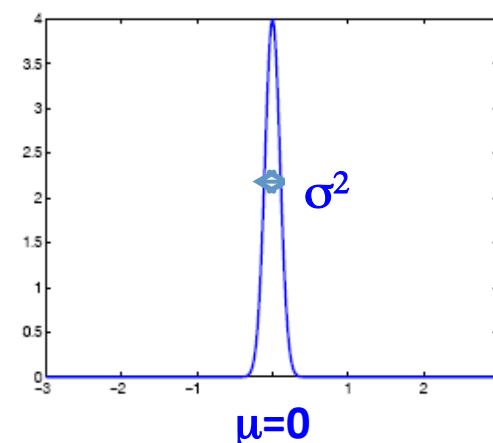
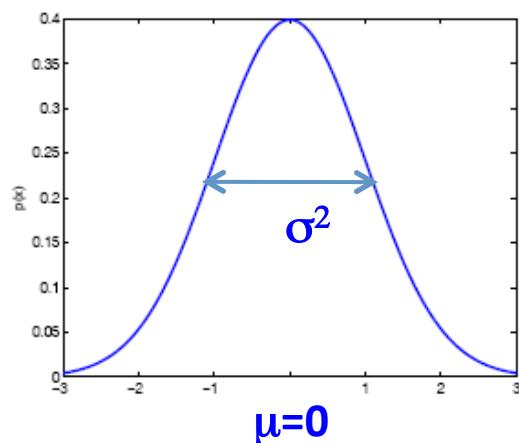


You give a  
different  
answer for  
different  
priors

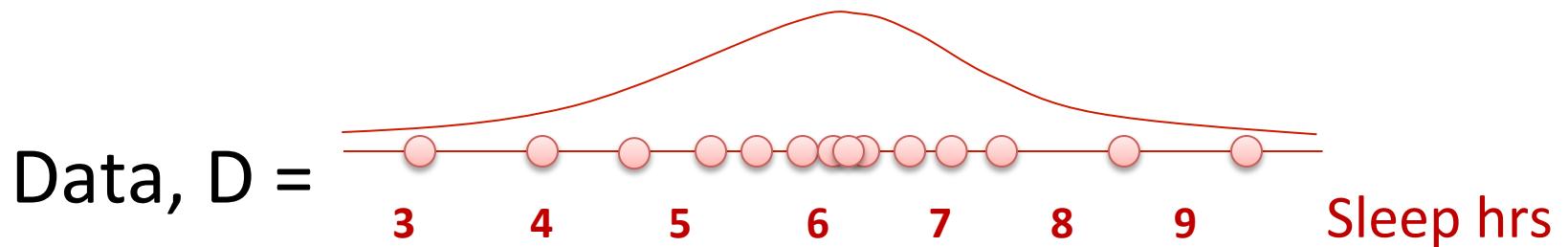
# What about continuous variables?

- Billionaire says: If I am measuring a continuous variable, what can you do for me?
- You say: Let me tell you about Gaussians...

$$P(x | \mu, \sigma) = \frac{1}{\sigma \sqrt{2\pi}} e^{\frac{-(x-\mu)^2}{2\sigma^2}} = N(\mu, \sigma^2)$$



# Gaussian distribution



- Parameters:  $\mu$  – mean,  $\sigma^2$  - variance
- Sleep hrs are **i.i.d.**:
  - **Independent** events
  - **Identically distributed** according to Gaussian distribution

# Properties of Gaussians

- affine transformation (multiplying by scalar and adding a constant)
  - $X \sim N(\mu, \sigma^2)$
  - $Y = aX + b \ ! Y \sim N(a\mu+b, a^2\sigma^2)$
- Sum of Gaussians
  - $X \sim N(\mu_X, \sigma^2_X)$
  - $Y \sim N(\mu_Y, \sigma^2_Y)$
  - $Z = X+Y \ ! Z \sim N(\mu_X+\mu_Y, \sigma^2_X+\sigma^2_Y)$

# MLE for Gaussian mean and variance

Choose  $\theta = (\mu, \sigma^2)$  that maximizes the probability of observed data

$$\begin{aligned}\hat{\theta}_{MLE} &= \arg \max_{\theta} P(D | \theta) \\ &= \arg \max_{\theta} \prod_{i=1}^n P(X_i | \theta)\end{aligned}$$

# MLE for Gaussian mean and variance

$$\hat{\mu}_{MLE} = \frac{1}{n} \sum_{i=1}^n x_i$$

$$\hat{\sigma}_{MLE}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \hat{\mu})^2$$

Note: MLE for the variance of a Gaussian is **biased**

- Expected result of estimation is **not** true parameter!
- Unbiased variance estimator:

$$\hat{\sigma}_{unbiased}^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \hat{\mu})^2$$

# MAP for Gaussian mean and variance

- Conjugate priors
  - Mean: Gaussian prior
  - Variance: Wishart Distribution
- Prior for mean:

$$P(\mu \mid \eta, \lambda) = \frac{1}{\lambda \sqrt{2\pi}} e^{\frac{-(\mu-\eta)^2}{2\lambda^2}}$$



# Bayes Optimal Classifier & Naive Bayes

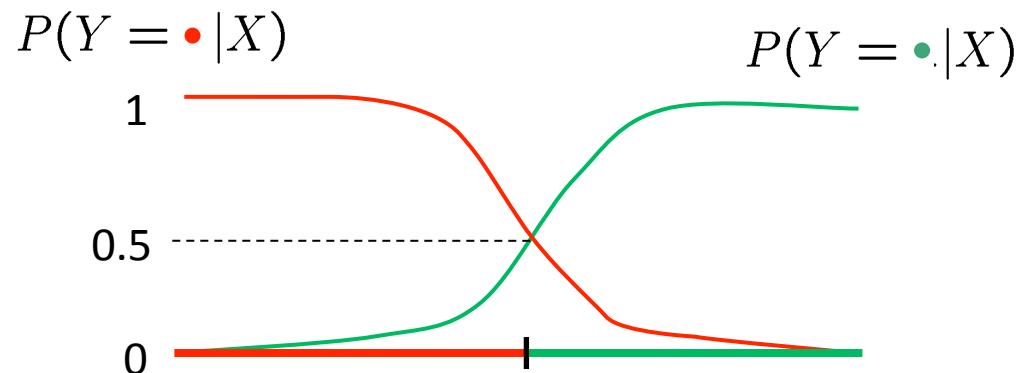


# Optimal Classification

Optimal predictor:  $f^* = \arg \min_f P(f(X) \neq Y)$   
(Bayes classifier)

Equivalently,

$$f^*(x) = \arg \max_{Y=y} P(Y = y | X = x)$$



# Optimal Classifier

**Bayes Rule:**  $P(Y|X) = \frac{P(X|Y)P(Y)}{P(X)}$

$$P(Y = y|X = x) = \frac{P(X = x|Y = y)P(Y = y)}{P(X = x)}$$

**Optimal classifier:**

$$\begin{aligned} f^*(x) &= \arg \max_{Y=y} P(Y = y|X = x) \\ &= \arg \max_{Y=y} P(X = x|Y = y)P(Y = y) \end{aligned}$$

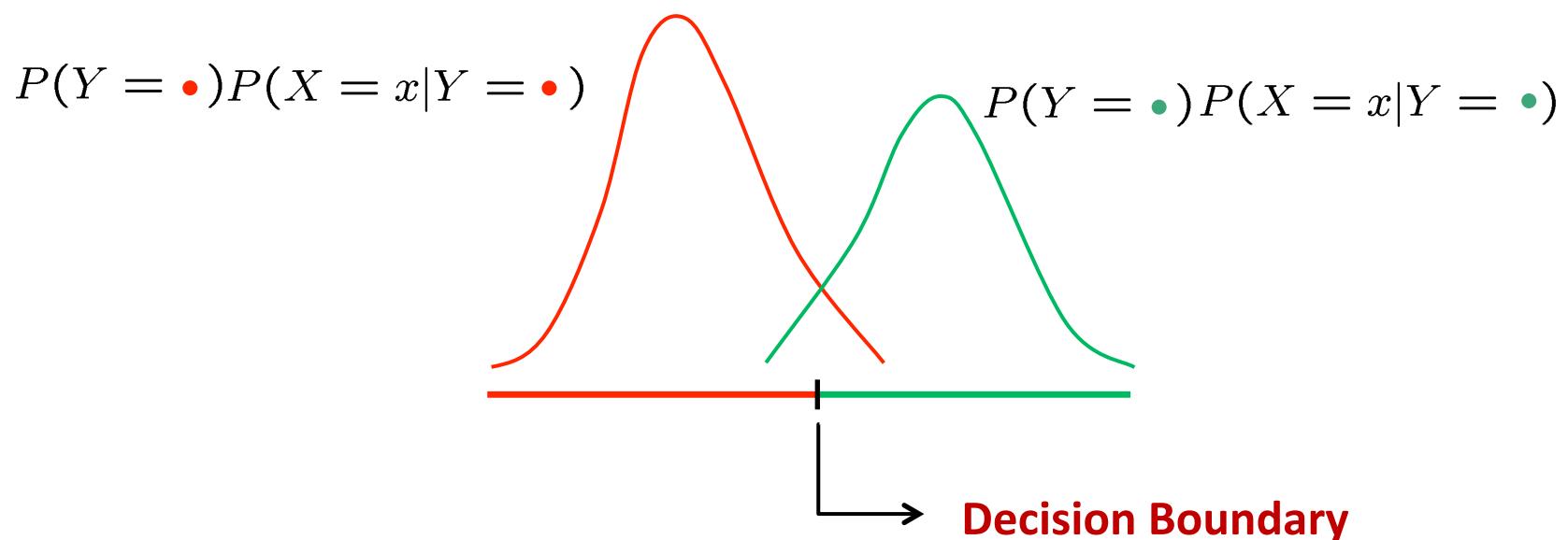
Class conditional    Class prior  
density

# Example Decision Boundaries

- Gaussian class conditional densities (1-dimension/feature)

$$P(X = x|Y = y) = \frac{1}{\sqrt{2\pi\sigma_y^2}} \exp\left(-\frac{(x - \mu_y)^2}{2\sigma_y^2}\right)$$

Binary Classification – two classes



# Learning the Optimal Classifier

**Optimal classifier:**

$$\begin{aligned} f^*(x) &= \arg \max_{Y=y} P(Y = y|X = x) \\ &= \arg \max_{Y=y} P(X = x|Y = y)P(Y = y) \end{aligned}$$


Class conditional      Class prior  
density

Need to know Prior  $P(Y = y)$  for all  $y$

Likelihood  $P(X=x|Y = y)$  for all  $x,y$

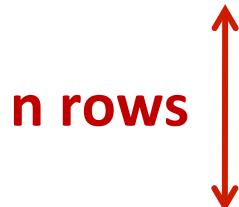
# Learning the Optimal Classifier

**Task:** Predict whether or not a picnic spot is enjoyable

**Training Data:**  $X = (X_1 \quad X_2 \quad X_3 \quad \dots \quad \dots \quad X_d)$   $Y$

Sky	Temp	Humid	Wind	Water	Forecast	EnjoySpt
Sunny	Warm	Normal	Strong	Warm	Same	Yes
Sunny	Warm	High	Strong	Warm	Same	Yes
Rainy	Cold	High	Strong	Warm	Change	No
Sunny	Warm	High	Strong	Cool	Change	Yes

n rows



Lets learn  $P(Y|X)$  – how many parameters?

Prior:  $P(Y = y)$  for all  $y$

K-1 if K labels

Likelihood:  $P(X=x|Y = y)$  for all  $x,y$

( $2^d - 1$ )K if d binary features

# Learning the Optimal Classifier

**Task:** Predict whether or not a picnic spot is enjoyable

**Training Data:**  $X = (X_1 \quad X_2 \quad X_3 \quad \dots \quad \dots \quad X_d)$   $Y$

Sky	Temp	Humid	Wind	Water	Forecast	EnjoySpt
Sunny	Warm	Normal	Strong	Warm	Same	Yes
Sunny	Warm	High	Strong	Warm	Same	Yes
Rainy	Cold	High	Strong	Warm	Change	No
Sunny	Warm	High	Strong	Cool	Change	Yes

**n rows**

Lets learn  $P(Y|X)$  – how many parameters?

$2^d K - 1$  ( $K$  classes,  $d$  binary features)

Need  $n >> 2^d K - 1$  number of training data to learn all parameters

# Conditional Independence

- X is **conditionally independent** of Y given Z:  
probability distribution governing X is independent of the value  
of Y, given the value of Z

$$(\forall x, y, z)P(X = x | Y = y, Z = z) = P(X = x | Z = z)$$

- Equivalent to:

$$P(X, Y | Z) = P(X | Z)P(Y | Z)$$

- e.g.,  $P(Thunder | Rain, Lightning) = P(Thunder | Lightning)$   
**Note:** does NOT mean Thunder is independent of Rain

# Prediction using Conditional Independence

- Predict Lightening
- From two **conditionally Independent** features
  - Thunder
  - Rain

# parameters needed to learn likelihood given L

$$P(T, R | L) \quad (2^2 - 1)2 = 6$$

With conditional independence assumption

$$P(T, R | L) = P(T | L) P(R | L) \quad (2-1)2 + (2-1)2 = 4$$

# Naïve Bayes Assumption

- Naïve Bayes assumption:
  - Features are conditionally independent given class:

$$\begin{aligned} P(X_1, X_2|Y) &= P(X_1|X_2, Y)P(X_2|Y) \\ &= P(X_1|Y)P(X_2|Y) \end{aligned}$$

- More generally:

$$P(X_1 \dots X_d|Y) = \prod_{i=1}^d P(X_i|Y)$$

- How many parameters now? **(2-1)dK vs. (2<sup>d</sup>-1)K**
  - Suppose  $\mathbf{X}$  is composed of  $d$  binary features

# Naïve Bayes Classifier

- Given:
  - Class Prior  $P(Y)$
  - $d$  conditionally independent features  $\mathbf{X}$  given the class  $Y$
  - For each  $X_i$ , we have likelihood  $P(X_i|Y)$

- Decision rule:

$$\begin{aligned}f_{NB}(\mathbf{x}) &= \arg \max_y P(x_1, \dots, x_d | y)P(y) \\&= \arg \max_y \prod_{i=1}^d P(x_i|y)P(y)\end{aligned}$$

- If conditional independence assumption holds, NB is optimal classifier! But worse otherwise.

# Naïve Bayes Algo – Discrete features

- Training Data  $\{(X^{(j)}, Y^{(j)})\}_{j=1}^n \quad X^{(j)} = (X_1^{(j)}, \dots, X_d^{(j)})$
- Maximum Likelihood Estimates

– For Class Prior

$$\hat{P}(y) = \frac{\{\#j : Y^{(j)} = y\}}{n}$$

– For Likelihood

$$\frac{\hat{P}(x_i, y)}{\hat{P}(y)} = \frac{\{\#j : X_i^{(j)} = x_i, Y^{(j)} = y\}/n}{\{\#j : Y^{(j)} = y\}/n}$$

- NB Prediction for test data  $X = (x_1, \dots, x_d)$

$$Y = \arg \max_y \hat{P}(y) \prod_{i=1}^d \frac{\hat{P}(x_i, y)}{\hat{P}(y)}$$

# Subtlety 1 – Violation of NB Assumption

- Usually, features are not conditionally independent:

$$P(X_1 \dots X_d | Y) \neq \prod_i P(X_i | Y)$$

- Nonetheless, NB is the single most used classifier out there
  - NB often performs well, even when assumption is violated
  - [Domingos & Pazzani '96] discuss some conditions for good performance

## Subtlety 2 – Insufficient training data

- What if you never see a training instance where  $X_1=a$  when  $Y=b$ ?
  - e.g.,  $Y=\{\text{SpamEmail}\}$ ,  $X_1=\{\text{'Earn'}\}$
  - $P(X_1=a \mid Y=b) = 0$
- Thus, no matter what the values  $X_2, \dots, X_d$  take:
  - $P(Y=b \mid X_1=a, X_2, \dots, X_d) = 0$

$$P(X_1 = a, X_2, \dots, X_n \mid Y) = P(X_1 = a \mid Y) \prod_{i=2}^d P(X_i \mid Y)$$

- What now???

# Naïve Bayes Algo – Discrete features

- Training Data  $\{(X^{(j)}, Y^{(j)})\}_{j=1}^n \quad X^{(j)} = (X_1^{(j)}, \dots, X_d^{(j)})$
- Maximum A Posteriori Estimates – add m “virtual” examples  
Assume priors

$$Q(Y = b) \quad Q(X_i = a, Y = b)$$

MAP Estimate

$$\hat{P}(X_i = a | Y = b) = \frac{\#\{j : X_i^{(j)} = a, Y^{(j)} = b\} + mQ(X_i = a, Y = b)}{\#\{j : Y^{(j)} = b\} + \underbrace{mQ(Y = b)}_{\text{\# virtual examples with } Y = b}}$$

Now, even if you never observe a class/feature posterior probability never zero.

# Case Study: Text Classification

- Classify e-mails
  - $Y = \{\text{Spam}, \text{NotSpam}\}$
- Classify news articles
  - $Y = \{\text{what is the topic of the article?}\}$
- Classify webpages
  - $Y = \{\text{Student, professor, project, ...}\}$
- What about the features  $X$ ?
  - The text!

# Features $X$ are entire document – $X_i$ for $i^{\text{th}}$ word in article

Article from rec.sport.hockey

---

Path: cantaloupe.srv.cs.cmu.edu!das-news.harvard.e  
From: xxx@yyy.zzz.edu (John Doe)  
Subject: Re: This year's biggest and worst (opinic  
Date: 5 Apr 93 09:53:39 GMT

I can only comment on the Kings, but the most obvious candidate for pleasant surprise is Alex Zhitnik. He came highly touted as a defensive defenseman, but he's clearly much more than that. Great skater and hard shot (though wish he were more accurate). In fact, he pretty much allowed the Kings to trade away that huge defensive liability Paul Coffey. Kelly Hrudey is only the biggest disappointment if you thought he was any good to begin with. But, at best, he's only a mediocre goaltender. A better choice would be Tomas Sandstrom, though not through any fault of his own, but because some thugs in Toronto decided

# NB for Text Classification

- $P(\mathbf{X}|\mathbf{Y})$  is huge!!!
  - Article at least 1000 words,  $\mathbf{X}=\{X_1, \dots, X_{1000}\}$
  - $X_i$  represents  $i^{\text{th}}$  word in document, i.e., the domain of  $X_i$  is entire vocabulary, e.g., Webster Dictionary (or more), 10,000 words, etc.
- NB assumption helps a lot!!!
  - $P(X_i=x_i|Y=y)$  is just the probability of observing word  $x_i$  at the  $i^{\text{th}}$  position in a document on topic  $y$

$$h_{NB}(\mathbf{x}) = \arg \max_y P(y) \prod_{i=1}^{LengthDoc} P(x_i|y)$$

# Bag of words model

- Typical additional assumption – **Position in document doesn't matter**:  $P(X_i=x_i | Y=y) = P(X_k=x_i | Y=y)$ 
  - “Bag of words” model – order of words on the page ignored
  - Sounds really silly, but often works very well!

$$\prod_{i=1}^{LengthDoc} P(x_i|y) = \prod_{w=1}^W P(w|y)^{count_w}$$

When the lecture is over, remember to wake up the person sitting next to you in the lecture room.

# Bag of words model

- Typical additional assumption – **Position in document doesn't matter**:  $P(X_i=x_i | Y=y) = P(X_k=x_i | Y=y)$ 
  - “Bag of words” model – order of words on the page ignored
  - Sounds really silly, but often works very well!

$$\prod_{i=1}^{LengthDoc} P(x_i|y) = \prod_{w=1}^W P(w|y)^{count_w}$$

in is lecture lecture next over person remember room  
sitting the the to to up wake when you

# Bag of words approach

*the world of*

**TOTAL**



***all about the company***

Our energy exploration, production, and distribution operations span the globe, with activities in more than 100 countries.

At TOTAL, we draw our greatest strength from our fast-growing oil and gas reserves. Our strategic emphasis on natural gas provides a strong position in a rapidly expanding market.

Our expanding refining and marketing operations in Asia and the Mediterranean Rim complement already solid positions in Europe, Africa, and the U.S.

Our growing specialty chemicals sector adds balance and profit to the core energy business.

**All About The Company**

- [Global Activities](#)
- [Corporate Structure](#)
- [TOTAL's Story](#)
- [Upstream Strategy](#)
- [Downstream Strategy](#)
- [Chemicals Strategy](#)
- [TOTAL Foundation](#)
- [Homepage](#)



aardvark	0
about	2
all	2
Africa	1
apple	0
anxious	0
...	
gas	1
...	
oil	1
...	
Zaire	0

# NB with Bag of Words for text classification

- Learning phase: using multiple training documents
  - Class Prior  $P(Y)$
  - $P(X_i|Y)$
- Test phase:
  - For each test document, use naïve Bayes decision rule:

$$\begin{aligned} h_{NB}(x) &= \arg \max_y P(y) \prod_{i=1}^{LengthDoc} P(x_i|y) \\ &= \arg \max_y P(y) \prod_{w=1}^W P(w|y)^{count_w} \end{aligned}$$

# Twenty news groups results

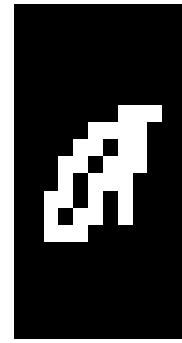
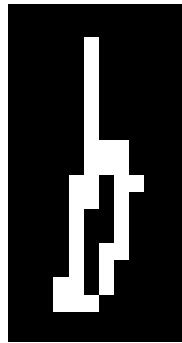
Given 1000 training documents from each group  
Learn to classify new documents according to  
which newsgroup it came from

comp.graphics	misc.forsale
comp.os.ms-windows.misc	rec.autos
comp.sys.ibm.pc.hardware	rec.motorcycles
comp.sys.mac.hardware	rec.sport.baseball
comp.windows.x	rec.sport.hockey
alt.atheism	sci.space
soc.religion.christian	sci.crypt
talk.religion.misc	sci.electronics
talk.politics.mideast	sci.med
talk.politics.misc	
talk.politics.guns	

Naive Bayes: 89% classification accuracy

# What if features are continuous?

Eg., character recognition:  $X_i$  is intensity at  $i^{\text{th}}$  pixel



Gaussian Naïve Bayes (GNB):

$$P(X_i = x \mid Y = y_k) = \frac{1}{\sigma_{ik}\sqrt{2\pi}} e^{\frac{-(x-\mu_{ik})^2}{2\sigma_{ik}^2}}$$

Different mean and variance for each class k and each pixel i.

Sometimes assume variance

- is independent of Y (i.e.,  $\sigma_i$ ),
- or independent of  $X_i$  (i.e.,  $\sigma_k$ )
- or both (i.e.,  $\sigma$ )

# Estimating parameters: Y discrete, X<sub>i</sub> continuous

Maximum likelihood estimates:

$$\hat{\mu}_{MLE} = \frac{1}{N} \sum_{j=1}^N x_j$$

$$\hat{\mu}_{ik} = \frac{1}{\sum_j \delta(Y^j = y_k)} \sum_j X_i^j \delta(Y^j = y_k)$$

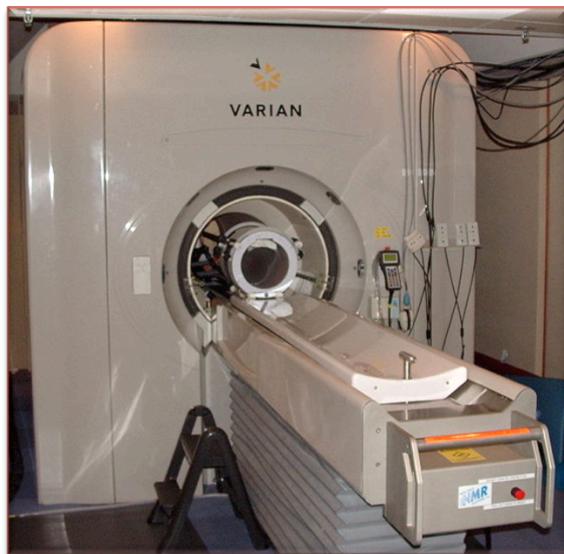
i<sup>th</sup> pixel in ←  
j<sup>th</sup> training image
→ k<sup>th</sup> class  
→ j<sup>th</sup> training image

$$\hat{\sigma}_{unbiased}^2 = \frac{1}{N-1} \sum_{j=1}^N (x_j - \hat{\mu})^2$$

$$\hat{\sigma}_{ik}^2 = \frac{1}{\sum_j \delta(Y^j = y_k) - 1} \sum_j (X_i^j - \hat{\mu}_{ik})^2 \delta(Y^j = y_k)$$

# Example: GNB for classifying mental states

[Mitchell et al.]



~1 mm resolution

~2 images per sec.

15,000 voxels/image

non-invasive, safe

measures Blood Oxygen Level Dependent (BOLD) response

