



Northeastern

**EECE2300**

# **Computational Methods for Data Analytics**

Lecture 3: Review of Probability Space and Set Theory

# Terms

1. **Set** is a collection of objects. These objects are called elements of the set
2. **Subset**  $b$  of a set  $a$  is a set whose elements are also elements of  $a$ ,  
*i.e.*,  $b \subset a$
3. **Space**  $S$  is the largest set; Thus, all other sets under consideration  $s_i \subset S$
4. **Null Set**  $O$  is an empty or null set.  $O$  contains no elements.

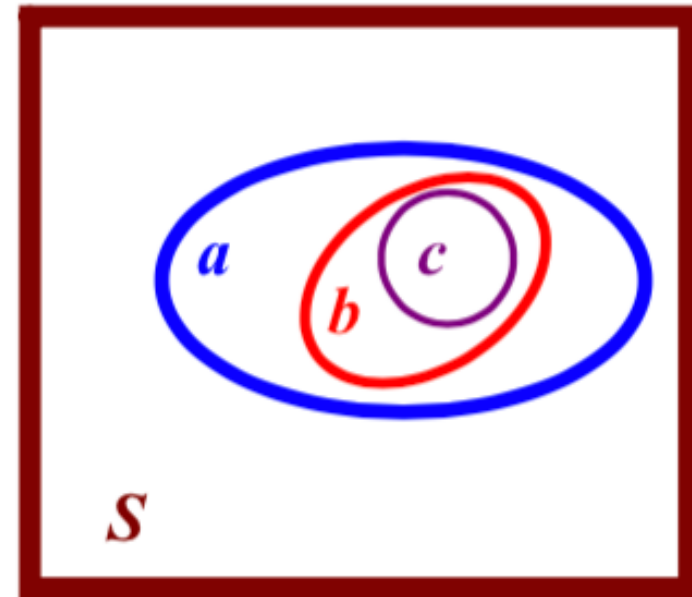


# Set Operations

Subset  $b \subset a$ , or the set  $a$  contains  $b$ ,  $a \supset b$ , if all elements of  $b$  are also elements of  $a$ . That is,

- If  $b \subset a$ , and  $c \subset b$ , then  $c \subset a$ .
- The following relationship holds:

$$a \subset a, 0 \subset a, a \subset S$$

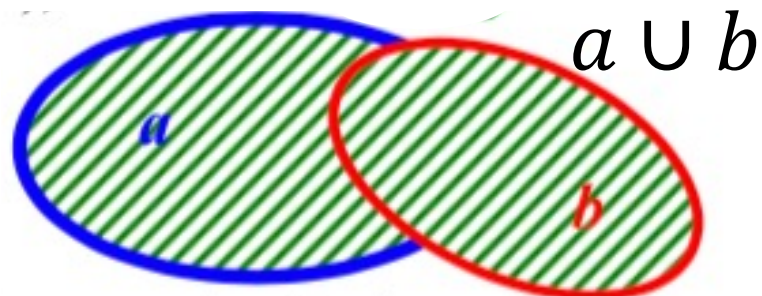


Example of subsets

# Set Operations

## 1. Equality

- $a = b$  iff  $a \subset b$  and  $b \subset a$



Example of union of sets  $a$  &  $b$

## 2. Union (Sum)

- The union of two sets  $a$  and  $b$  is a set consisting of all elements of  $a$  or of  $b$  or of both. The union operation satisfies the following properties:

$$a \cup b = b \cup a$$

(Commutative),

$$a \cup a = a$$

$$a \cup 0 = a$$

$$a \cup S = S$$

$$(a \cup b) \cup c = a \cup (b \cup c) = a \cup b \cup c$$

(Associative)

# Set Operations

3. **Intersection (Product)** of two sets  $a$  and  $b$  is a set consisting of all elements that are common to the sets  $a$  and  $b$ . The intersection operation satisfies the following properties:

$$a \cap b = b \cap a \quad (\text{Commutative}),$$

$$a \cap a = a$$

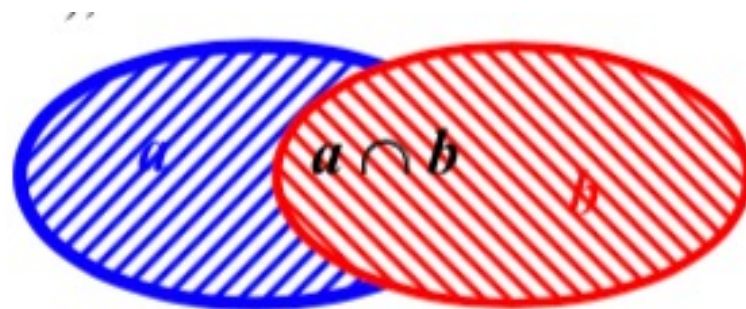
$$a \cap 0 = 0$$

$$a \cap S = a$$

$$(a \cap b) \cap c = a \cap (b \cap c) = a \cap b \cap c \quad (\text{Associative})$$

$$\text{If } b \subset a, b \cap a = b$$

$$\text{Also, } a \cap (b \cap c) = (a \cap b) \cap (a \cap c) \quad (\text{Distributive})$$



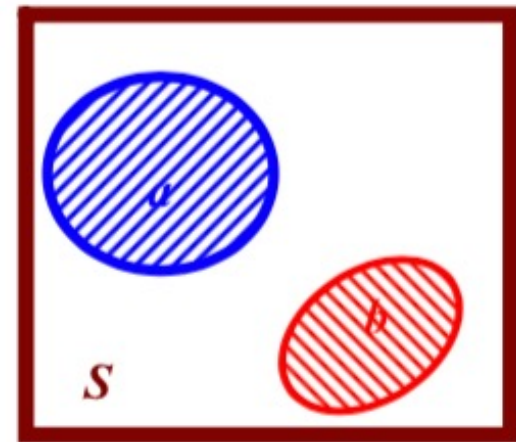
Example of intersection of sets  $a$  &  $b$

# Mutually Exclusive Sets

- Two sets  $a$  and  $b$  are called mutually exclusive or disjoint if they have no common elements, i.e.

$$a \cap b = 0$$

- The sets  $a_1, a_2, \dots$  are called mutually An of mutually sets.
- Exclusive if  $a_i \cap a_j = 0$  for every  $i \neq j$ .



Example of mutually exclusive sets

# Compliments

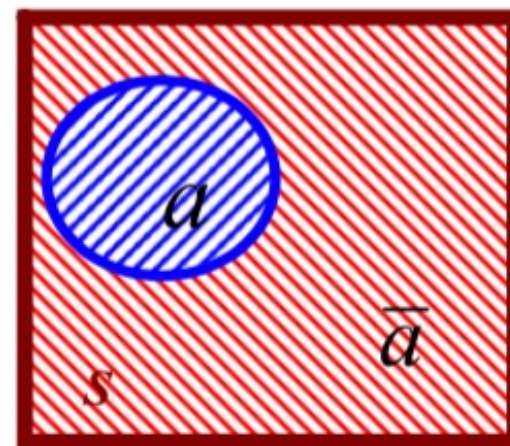
- The complement  $\bar{a}$  of a set  $a$  is defined as a set consisting of all elements of  $S$  that are not in  $a$ . Complement sets satisfy the following properties:

$$a \cup \bar{a} = S$$

$$a \cap \bar{a} = 0$$

$$\bar{0} = S, \bar{S} = 0$$

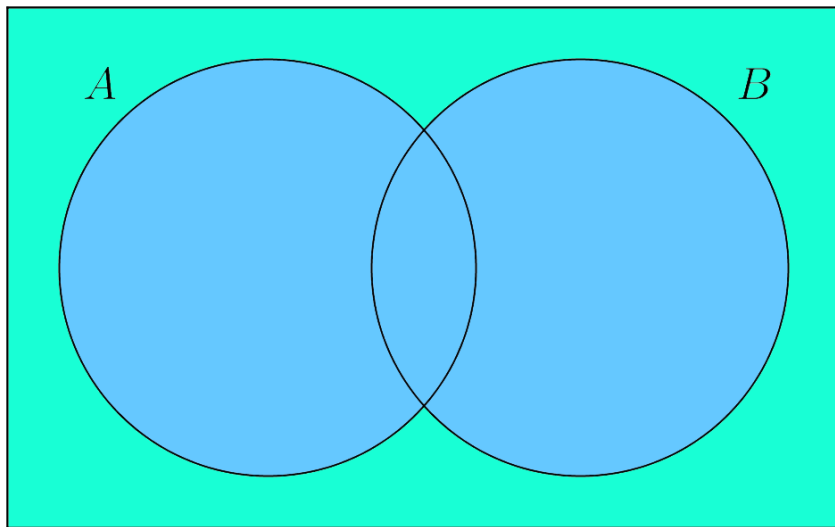
$$\text{if } b \subset a, \bar{b} \supset \bar{a}$$



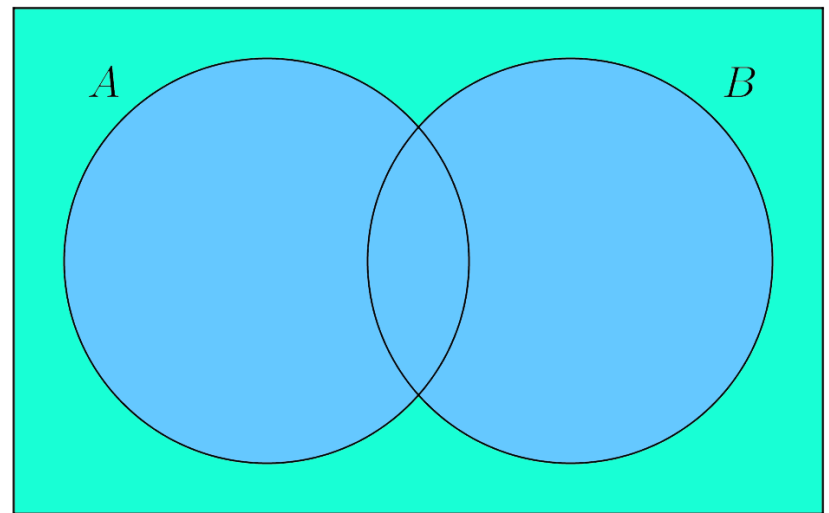
Example of complements

# De Morgan Law

$$\overline{a \cup b} = \bar{a} \cap \bar{b}.$$



$$\overline{a \cap b} = \bar{a} \cup \bar{b}$$





# Difference of Two Sets

The difference set of  $a - b$  is a set consisting of elements of  $a$  that are not in  $b$ . The difference satisfy the following properties:

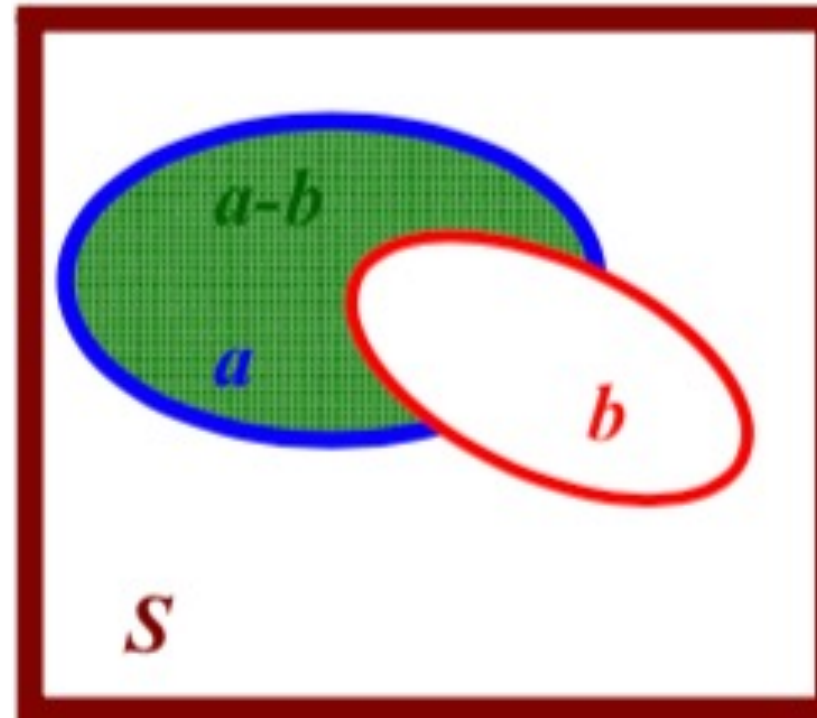
$$a - b = a \cap \bar{b} = a - a \cap b,$$

$$a \cup a = a,$$

$$(a - a) \cup a = a,$$

$$\bar{a} = S - a,$$

$$a = (a - b) \cup (b \cap a).$$



Example of difference of 2 sets

# Probability Space

## 1. *Random Experiment* $\xi$

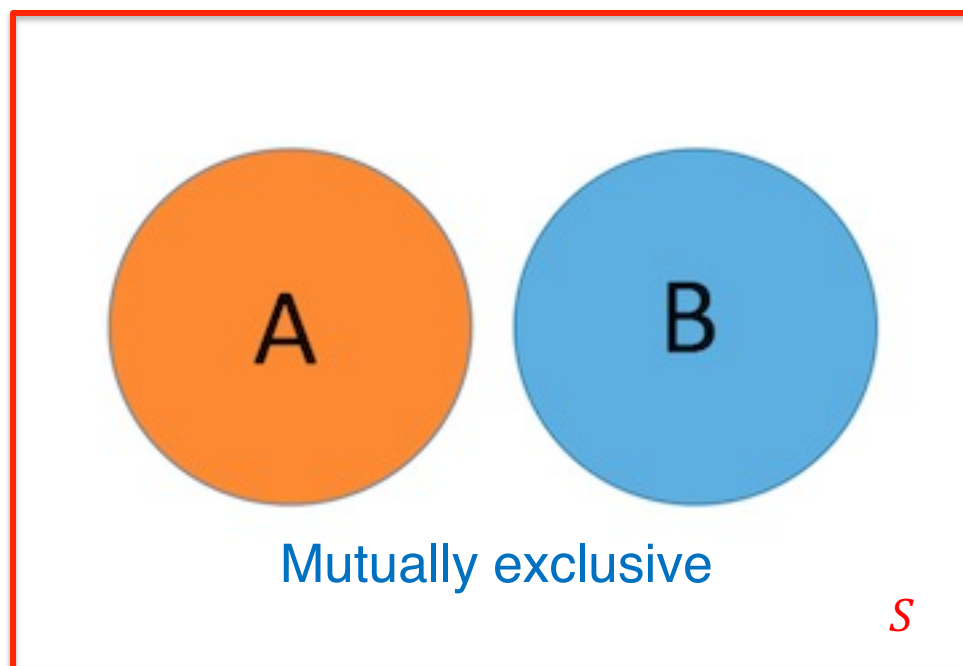
- By an experiment  $I$ , we mean a (set) space  $S$  of outcomes  $\xi$ .
- Elements of  $S$  are *outcomes* or *elementary events*.
- $S$  is a probability (sample) space. Subsets of  $S$  are called *events*.
  - Space  $S$  is the *sure (certain) event*.
  - Empty set  $\emptyset$  is the *impossible event*.



# Probability Space

## 2. *Mutually Exclusive Events*

- Two events  $a$  and  $b$  are mutually exclusive if  $a \cap b = \emptyset$ .



# Probability Space

## 3. Axioms of Probability

- Each event  $a_i$  is a measure (i.e., number)  $P(a_i)$  called the *probability of event*. Thus,  $a_i$  is assigned a  $P(a_i)$  subjected to the following 3 axioms:
  1.  $P(a_i) \geq 0$
  2.  $P(S) = 1$
  3. If  $a_1 \cap a_2 = \emptyset$ , then  $P(a_1 \cup a_2) = P(a_1) + P(a_2)$

### Corollaries:

$$P(\emptyset) = 0$$

$$P(a_i) = 1 - P(\bar{a}_i) \leq 1$$

$$\text{If } a_1 \cap a_2 \neq \emptyset, \text{ then } P(a_1 \cup a_2) = P(a_1) + P(a_2) - P(a_1 \cap a_2)$$

$$\text{If } b \subset a, P(a) = P(b) + P(a \cap \bar{b}) \geq P(b)$$



# Probability Space

**Field**  $F$  is a nonempty class of sets such that

1. If  $a \in F$ , then  $\bar{a} \in F$ ;
2. If  $a \in F$  and  $b \in F$ , then  $a \cup b \in F$ .

**Corollaries:**

If  $a \in F$  and  $b \in F$ , then  $a \cup b \in F$  and  $a - b \in F$ ;

Also  $0 \in F$  and  $S \in F$ .

**Borel Field** has the property that if sets  $a_1, a_2, \dots, a_n, \dots$  belong to it, and so does the set  $a_1 \cup a_2 \cup \dots \cup a_n \cup \dots$ , then the field is called a Borel field.

Note that the class of all subsets of  $S$  is called the Borel field.



# Probability Space

**Probability Experiment  $\xi$  is:**

1. A set  $S$  of outcomes  $\xi$ :
  - a. this set is called probability space
2. A Borel field  $F$  consisting of certain subsets of  $S$  called events
3. A measure (number)  $P(a_i)$  assigned to every event  $a_i$ :
  - a. Measure called *probability of event  $a_i$* , if satisfies axioms 1-3

- 
- Common to use the following notation for probability experiments:
    - $\xi : (S, F, P)$  identifies a probability experiment with space of outcomes  $S$ , & associated field  $F$  with  $P(a)$  for all outcomes assigned.

Example: Probability experiment of tossing a coin,  $\xi : (S, F, P)$ . The space is  $S = \{h, t\}$

The events are?

$$F: 0, \{h\}, \{t\}, \{h, t\}$$

Probability of the events?

$$P(h) = p, P\{t\} = q, p + q = 1$$



# Probability Theory

**Random variable (RV):** a variable whose possible values are numerical outcomes of a random phenomenon.

**Examples:** A person's height, the outcome of a coin toss

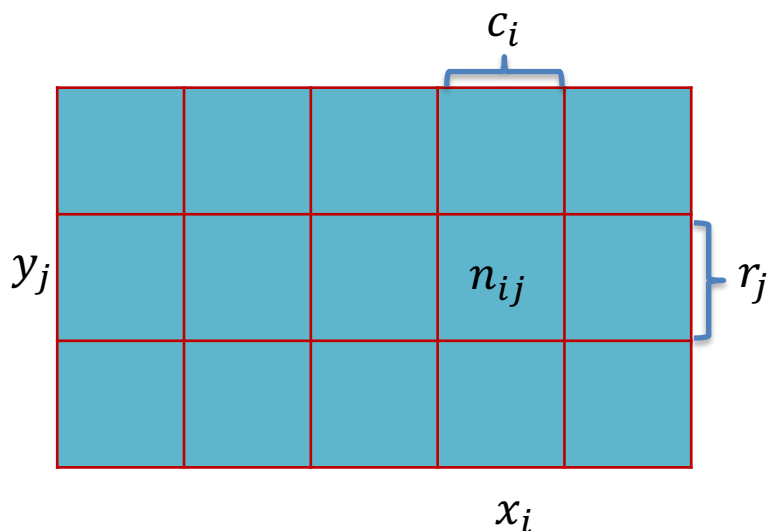
Distinguish between discrete and continuous variables.

The **distribution** of a discrete random variable:

- The probabilities of each value it can take.
- Notation:  $P(X = x_i)$ .
- These numbers satisfy:

$$\sum_i P(X = x_i) = 1$$

# Probability Theory



## Marginal Probability

$$p(X = x_i) = \frac{c_i}{N}$$

## Joint Probability

$$p(X = x_i, Y = y_j) = \frac{n_{ij}}{N}$$

## Conditional Probability

$$p(Y = y_j | X = x_i) = \frac{n_{ij}}{c_i}$$



# Probability Theory

A joint probability distribution for two RVs is a table

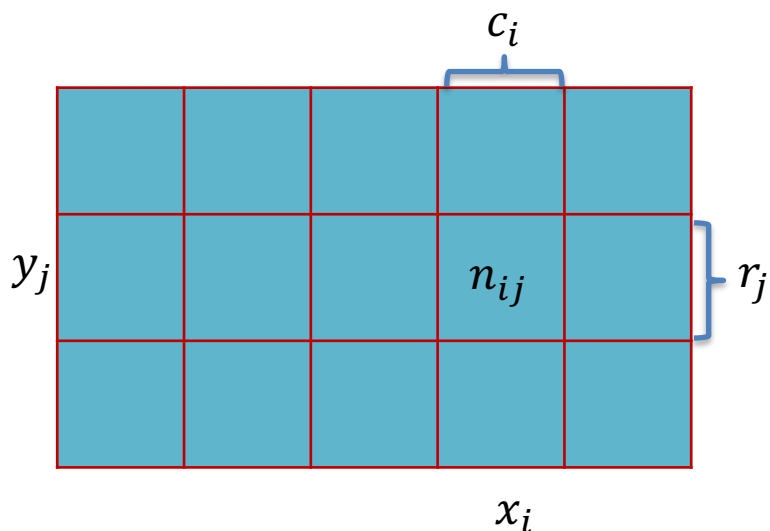
			$c_i$	
$y_j$			$n_{ij}$	$r_j$
			$x_i$	

If the two variables are binary, how many parameters does it have?

What about joint probability of  $d$  variables  $P(X_1, \dots, X_d)$ ?

How many parameters does it have if each variable is a binary?

# Probability Theory



## Marginalization

$$\begin{aligned} p(X = x_i) &= \frac{c_i}{N} = \frac{1}{N} \sum_{j=1}^L n_{ij} \\ &= \sum_{j=1}^L p(X = x_i, Y = y_j) \end{aligned}$$

## Product Rule

$$\begin{aligned} p(X = x_i, Y = y_j) &= \frac{n_{ij}}{N} = \frac{n_{ij}}{c_i} \cdot \frac{c_i}{N} \\ &= p(Y = y_j | X = x_i) p(X = x_i) \end{aligned}$$

# The Rules of Probability

**Marginalization**

$$p(X) = \sum_Y p(X, Y)$$

**Product Rule**

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

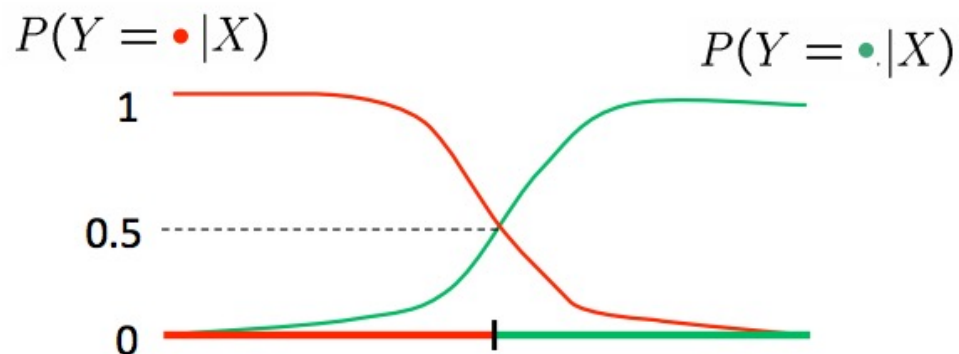
**Independence:**  $X$  and  $Y$  are independent if  $P(Y|X) = P(Y)$

This implies  $P(X, Y) = P(X)P(Y)$

# Using Probability in learning

**Interested in:**  $P(Y|X)$

$X$  – features    $Y$  – labels



For example, when classifying spam, we could estimate  $P(Y | \text{Viagara, lottery})$

We would then classify an example if  $P(Y | X) > 0.5$ .

However, it's usually easier to model  $P(X | Y)$

# Maximum Likelihood

Fit a probabilistic model  $P(X|\theta)$  to data

- Estimate  $\theta$

Given independent identically distributed (i.i.d.) data  $X = (x_1, x_2, \dots, x_n)$

- Likelihood

$$P(X|\theta) = P(x_1|\theta) P(x_2|\theta), \dots, P(x_n|\theta)$$

- Log Likelihood

$$\ln P(X|\theta) = \sum_{i=1}^n \ln P(x_i|\theta)$$

**Maximum likelihood** solution: parameters  $\theta$  that maximize  $\ln P(X|\theta)$

# Example

Example: coin toss

Estimate the probability  $p$  that a coin lands “Heads” using the result of  $n$  coin tosses,  $h$  of which resulted in heads.

The likelihood of the data:  $P(X|\theta) = p^h(1 - p)^{n-h}$

Log likelihood:  $\ln P(X|\theta) = h \ln(p) + (n - h) \ln(1 - p)$

Taking a derivative and setting to 0:

$$\frac{\partial \ln P(X|\theta)}{\partial p} = \frac{h}{p} - \frac{(n - h)}{(1 - p)} = 0$$

$$\Rightarrow p = \frac{h}{n}$$

# Bayes' rule

From the product rule

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

and

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

Therefore:

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

This is known as Baye's rule

# Bayes' rule

$$\underset{\text{Posterior}}{P(Y|X)} = \frac{\overset{\text{Likelihood}}{P(X|Y)} \overset{\text{Prior}}{P(Y)}}{\underset{\text{Evidence}}{P(X)}}$$

**Posterior  $\propto$  likelihood x prior**

$P(X)$  can be computed as

$$P(X) = \sum_Y P(X|Y)P(Y)$$

But is not important for inferring a label



# Maximum a-posteriori and maximum likelihood

The maximum a posteriori (MAP) rule:

$$y_{MAP} = \arg \max_Y P(Y|X) = \arg \max_Y \frac{P(X|Y)P(Y)}{P(X)} = \arg \max_Y P(X|Y)P(Y)$$

If we ignore the prior distribution or assume it is uniform we obtain the maximum likelihood rule:

$$y_{ML} = \arg \max_Y P(X|Y)$$

A classifier that has access to  $P(Y|X)$  is a Bayes optimal classifier.

# Naïve Bayes classifier

We would like to model  $P(X | Y)$ , where  $X$  is a feature vector, and  $Y$  is its associated label.

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

**Training Data:**

$X = (X_1$	$X_2$	$X_3$	...	...	$X_d)$	$Y$
Sky	Temp	Humid	Wind	Water	Forecst	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**

How many parameters?

Prior:  $P(Y)$   $k-1$  if  $k$  classes

Likelihood:  $P(X | Y)$   $(2^d - 1)k$  for binary features

# Naïve Bayes classifier

We would like to model  $P(\mathbf{X} | Y)$ , where  $\mathbf{X}$  is a feature vector, and  $Y$  is its associated label.

Simplifying assumption: conditional independence: given the class label the features are independent, i.e.

$$P(\mathbf{X}|Y) = P(x_1|Y)P(x_2|Y), \dots, P(x_d|Y)$$

How many parameters now?

# Naïve Bayes classifier

We would like to model  $P(\mathbf{X} | Y)$ , where  $\mathbf{X}$  is a feature vector, and  $Y$  is its associated label.

Simplifying assumption: conditional independence: given the class label the features are independent, i.e.

$$P(\mathbf{X}|Y) = P(x_1|Y)P(x_2|Y), \dots, P(x_d|Y)$$

How many parameters now?  $dk + k - 1$

# Training a Naïve Bayes classifier

Training data: Feature matrix  $X$  ( $n \times d$ ) and labels  $y_1, \dots, y_n$

Maximum likelihood estimates:

Class prior: 
$$\hat{P}(y) = \frac{|\{i : y_i = y\}|}{n}$$

Likelihood: 
$$\hat{P}(x_i|y) = \frac{\hat{P}(x_i, y)}{\hat{P}(y)} = \frac{|\{i : X_{ij} = x_i, y_i = y\}|/n}{|\{i : y_i = y\}|/n}$$

# Example

## Email classification

Suppose our vocabulary contains three words  $a$ ,  $b$  and  $c$ , and we use a multivariate Bernoulli model for our e-mails, with parameters

$$\theta^{\oplus} = (0.5, 0.67, 0.33)$$

$$\theta^{\ominus} = (0.67, 0.33, 0.33)$$

This means, for example, that the presence of  $b$  is twice as likely in spam (+), compared with ham.

The e-mail to be classified contains words  $a$  and  $b$  but not  $c$ , and hence is described by the bit vector  $\mathbf{x} = (1, 1, 0)$ . We obtain likelihoods

$$P(\mathbf{x}|\oplus) = 0.5 \cdot 0.67 \cdot (1 - 0.33) = 0.222$$

$$P(\mathbf{x}|\ominus) = 0.67 \cdot 0.33 \cdot (1 - 0.33) = 0.148$$

The ML classification of  $\mathbf{x}$  is thus spam.

# Example

Email classification: training data

E-mail	$a?$	$b?$	$c?$	Class
$e_1$	0	1	0	+
$e_2$	0	1	1	+
$e_3$	1	0	0	+
$e_4$	1	1	0	+
$e_5$	1	1	0	-
$e_6$	1	0	1	-
$e_7$	1	0	0	-
$e_8$	0	0	0	-

What are the parameters of the model?

# Example

## Email classification: training data

E-mail	$a?$	$b?$	$c?$	Class
$e_1$	0	1	0	+
$e_2$	0	1	1	+
$e_3$	1	0	0	+
$e_4$	1	1	0	+
$e_5$	1	1	0	-
$e_6$	1	0	1	-
$e_7$	1	0	0	-
$e_8$	0	0	0	-

What are the parameters of the model?

$$\hat{P}(y) = \frac{|\{i : y_i = y\}|}{n}$$

$$\hat{P}(x_i|y) = \frac{\hat{P}(x_i, y)}{\hat{P}(y)} = \frac{|\{i : X_{ij} = x_i, y_i = y\}|/n}{|\{i : y_i = y\}|/n}$$



# Example

## Email classification: training data

E-mail	$a?$	$b?$	$c?$	Class
$e_1$	0	1	0	+
$e_2$	0	1	1	+
$e_3$	1	0	0	+
$e_4$	1	1	0	+
$e_5$	1	1	0	-
$e_6$	1	0	1	-
$e_7$	1	0	0	-
$e_8$	0	0	0	-

What are the parameters of the model?

$$P(+) = 0.5, P(-) = 0.5$$

$$P(a|+) = 0.5, P(a|-) = 0.75$$

$$P(b|+) = 0.75, P(b|-) = 0.25$$

$$P(c|+) = 0.25, P(c|-) = 0.25$$

$$\hat{P}(y) = \frac{|\{i : y_i = y\}|}{n}$$

$$\hat{P}(x_i|y) = \frac{\hat{P}(x_i, y)}{\hat{P}(y)} = \frac{|\{i : X_{ij} = x_i, y_i = y\}|/n}{|\{i : y_i = y\}|/n}$$

# Comments on Naïve Bayes

Usually features are not conditionally independent, i.e.

$$P(\mathbf{X}|Y) \neq P(x_1|Y)P(x_2|Y), \dots, P(x_d|Y)$$

And yet, one of the most widely used classifiers. Easy to train!

It often performs well even when the assumption is violated.

Domingos, P., & Pazzani, M. (1997). Beyond Independence: Conditions for the Optimality of the Simple Bayesian Classifier. *Machine Learning*. 29, 103-130.

# When there are few training examples

What if you never see a training example where  $x_1=a$  when  $y=\text{spam}$ ?

$$P(x \mid \text{spam}) = P(a \mid \text{spam}) P(b \mid \text{spam}) P(c \mid \text{spam}) = 0$$

What to do?

# When there are few training examples

What if you never see a training example where  $x_1=a$  when  $y=\text{spam}$ ?

$$P(x \mid \text{spam}) = P(a \mid \text{spam}) P(b \mid \text{spam}) P(c \mid \text{spam}) = 0$$

What to do?

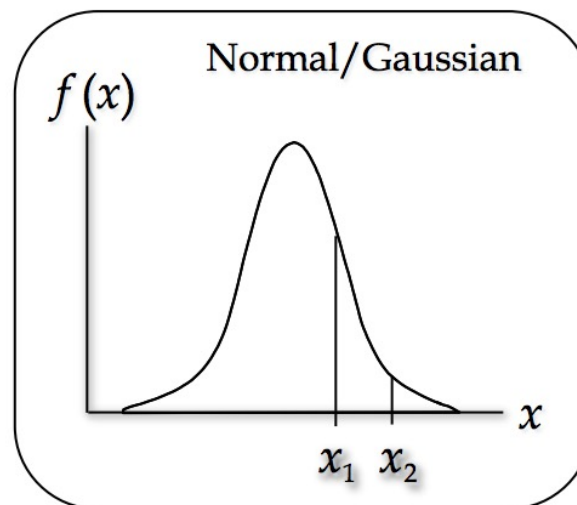
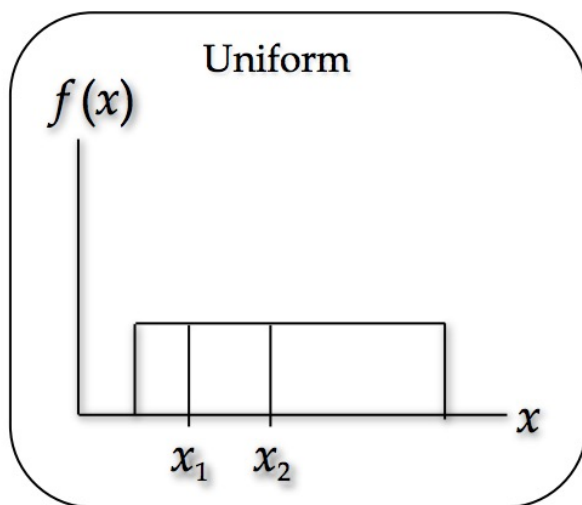
Add “virtual” examples for which  $x_1=a$  when  $y=\text{spam}$ .

# Naïve Bayes for continuous variables

**Need to talk about continuous distributions!**

# Continuous Probability Distributions

The probability of the random variable assuming a value within some given interval from  $x_1$  to  $x_2$  is defined to be the area under the graph of the probability density function between  $x_1$  and  $x_2$ .



# Expectations

Discrete variables

$$\mathbb{E}[f] = \sum_x p(x) f(x)$$

$$\mathbb{E}_x[f|y] = \sum_x p(x|y) f(x)$$

↑    - - - - - > x

$$\mathbb{E}[f] \simeq \frac{1}{N} \sum_{n=1}^N f(x_n)$$

Continuous variables

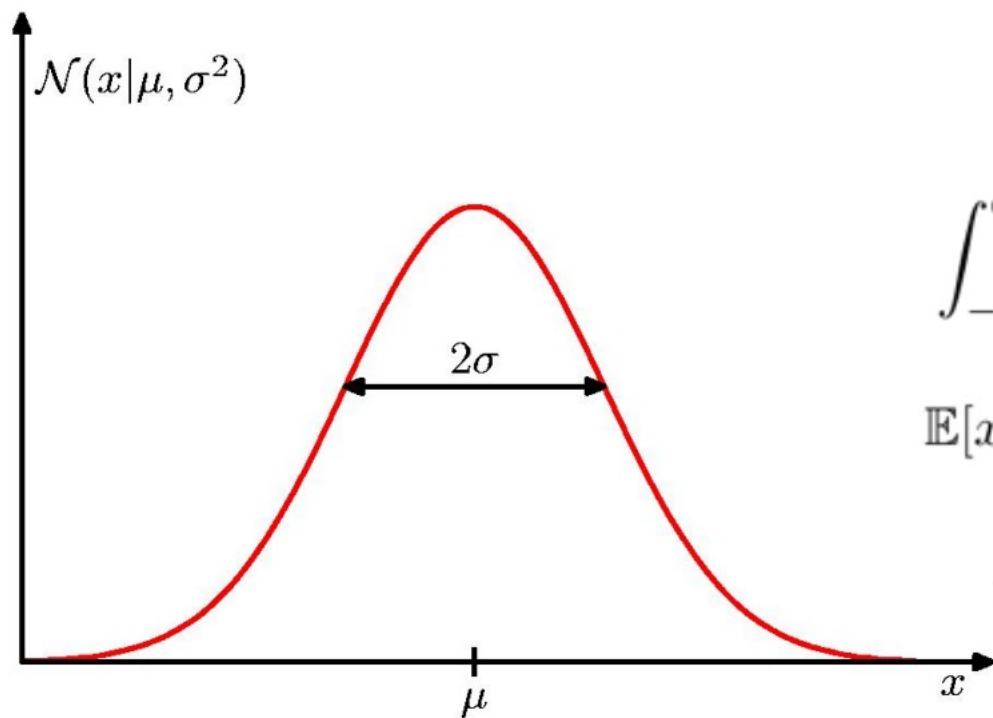
$$\mathbb{E}[f] = \int p(x) f(x) dx$$

Conditional expectation  
(discrete)

Approximate expectation  
(discrete and continuous)

# The Gaussian (normal) distribution

$$\mathcal{N}(x|\mu, \sigma^2) = \frac{1}{(2\pi\sigma^2)^{1/2}} \exp \left\{ -\frac{1}{2\sigma^2} (x - \mu)^2 \right\}$$



$$\mathcal{N}(x|\mu, \sigma^2) > 0$$

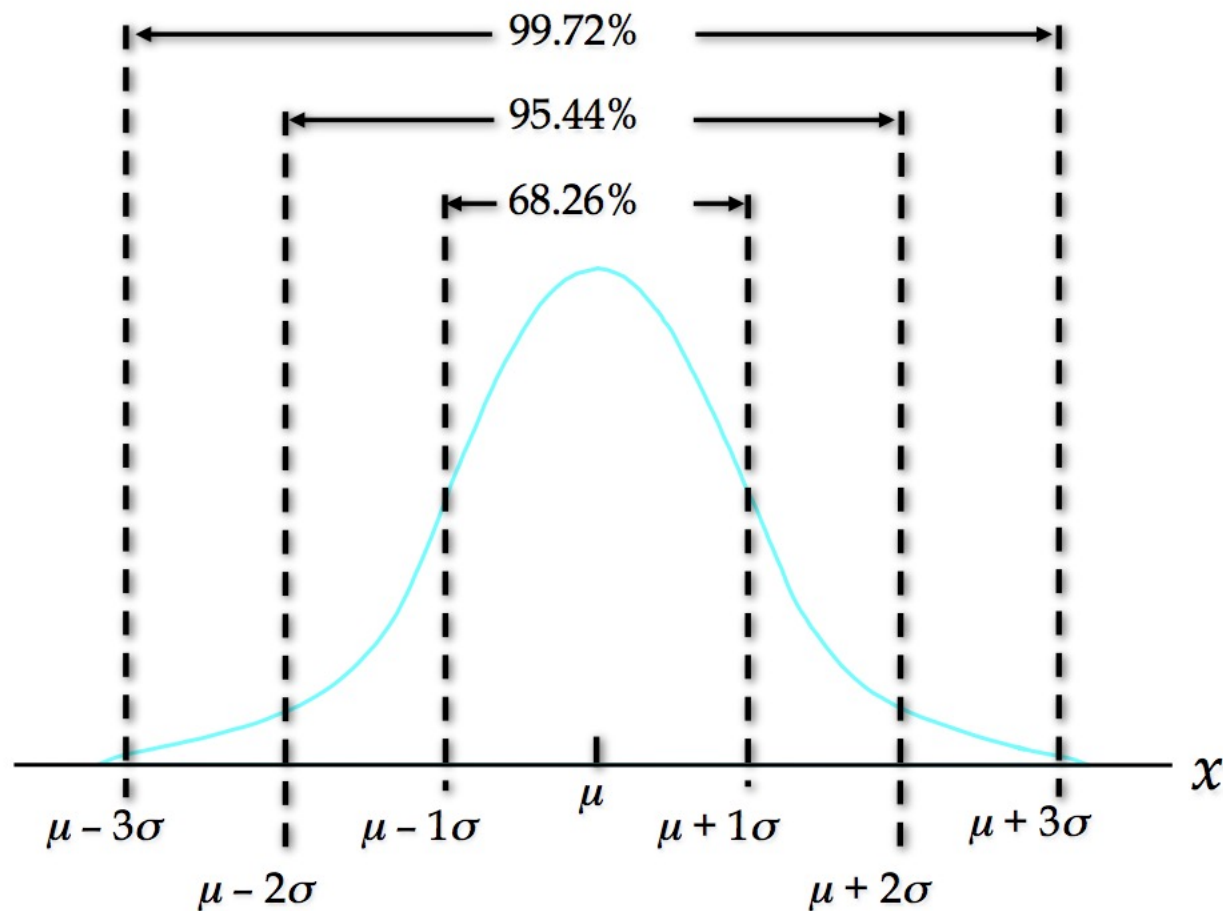
$$\int_{-\infty}^{\infty} \mathcal{N}(x|\mu, \sigma^2) dx = 1$$

$$\mathbb{E}[x] = \int_{-\infty}^{\infty} \mathcal{N}(x|\mu, \sigma^2) x dx = \mu$$

$$\text{var}[x] = \mathbb{E}[x^2] - \mathbb{E}[x]^2 = \sigma^2$$



# Properties of the Gaussian distribution




# Standard Normal Distribution

▶ A random variable having a normal distribution with a mean of 0 and a standard deviation of 1 is said to have a standard normal probability distribution.

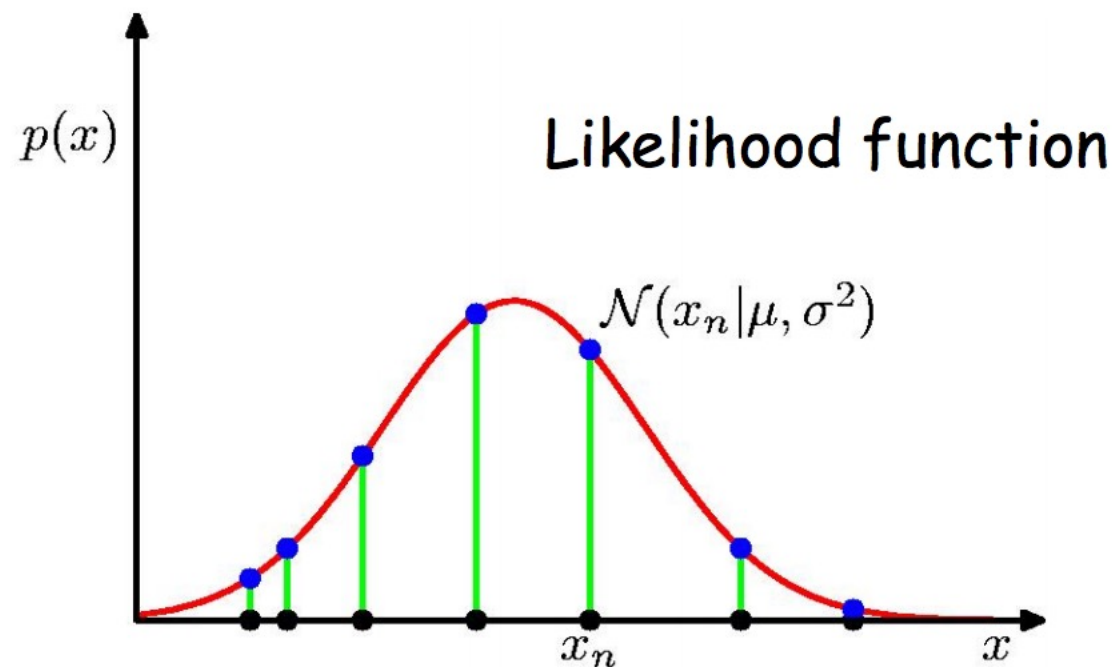
# Standard Normal Probability Distribution

## ■ Converting to the Standard Normal Distribution


$$z = \frac{x - \mu}{\sigma}$$

We can think of  $z$  as a measure of the number of standard deviations  $x$  is from  $\mu$ .

# Gaussian Parameter Estimation



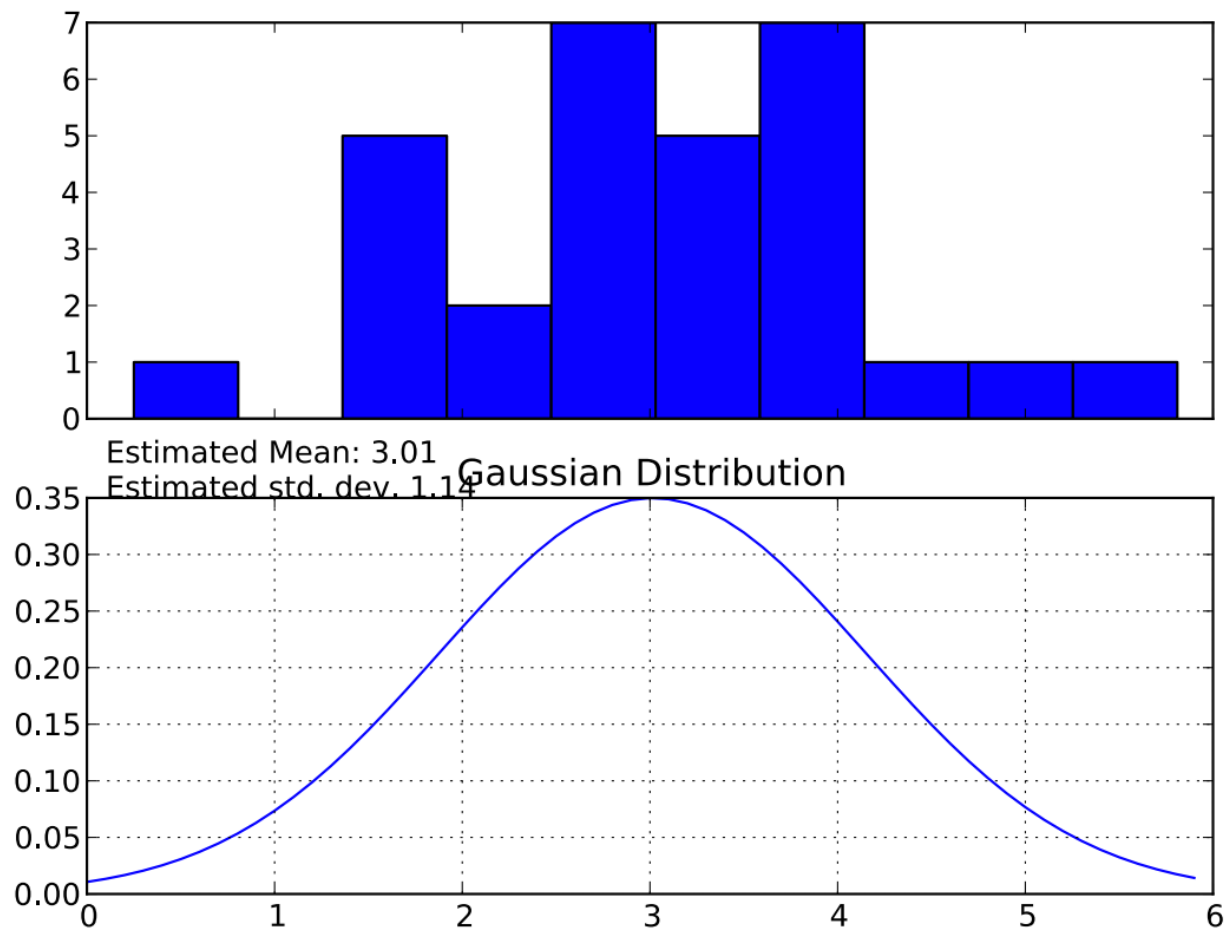
$$p(\mathbf{x} | \mu, \sigma^2) = \prod_{n=1}^N \mathcal{N}(x_n | \mu, \sigma^2)$$

# Maximum (Log) Likelihood

$$\ln p(\mathbf{x}|\mu, \sigma^2) = -\frac{1}{2\sigma^2} \sum_{n=1}^N (x_n - \mu)^2 - \frac{N}{2} \ln \sigma^2 - \frac{N}{2} \ln(2\pi)$$

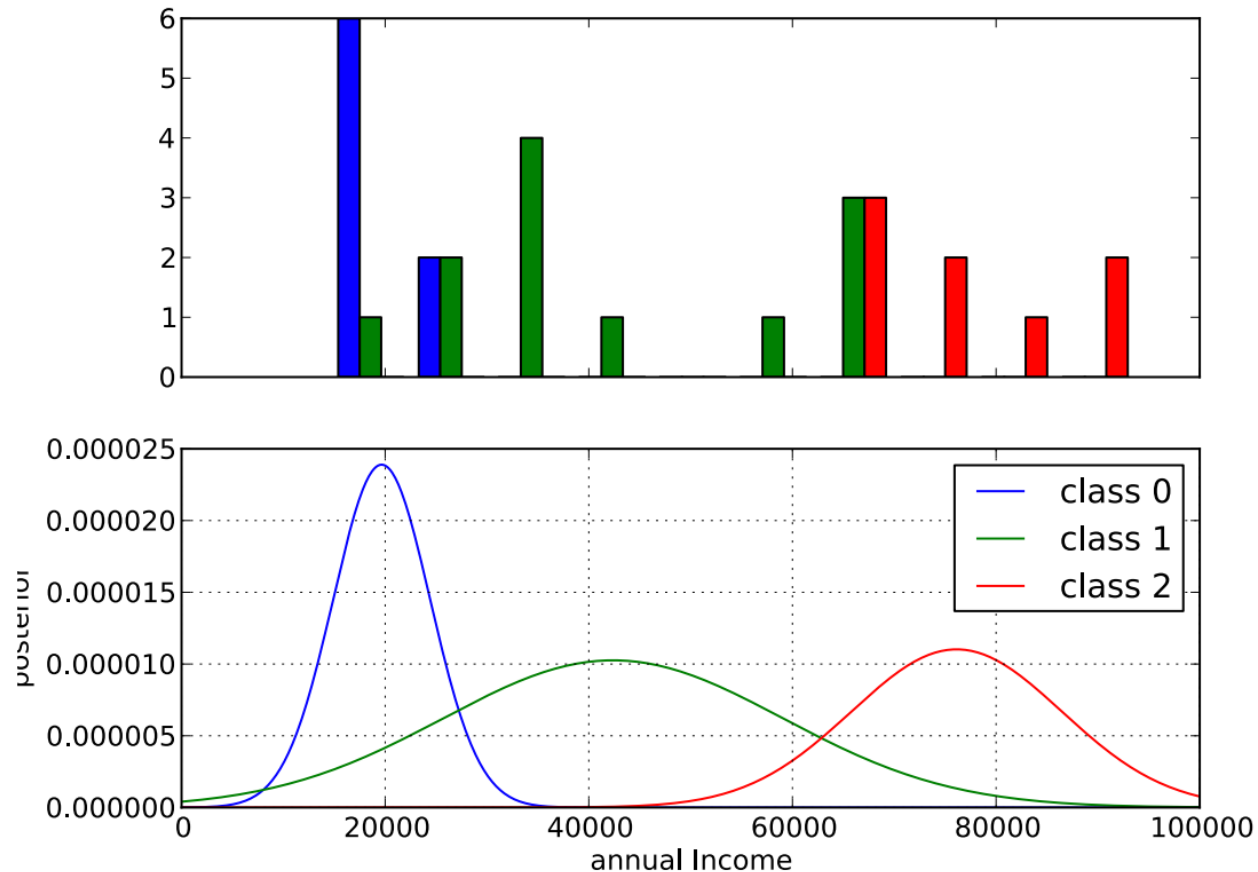
$$\mu_{\text{ML}} = \frac{1}{N} \sum_{n=1}^N x_n \qquad \sigma_{\text{ML}}^2 = \frac{1}{N} \sum_{n=1}^N (x_n - \mu_{\text{ML}})^2$$

# Example



# Gaussian models

Assume we have data that belongs to three classes, and assume a likelihood that follows a Gaussian distribution



# Gaussian Naïve Bayes

Likelihood function:

$$P(X_i = x|Y = y_k) = \frac{1}{\sqrt{2\pi}\sigma_{ik}} \exp\left(-\frac{(x - \mu_{ik})^2}{2\sigma_{ik}^2}\right)$$

Need to estimate mean and variance for each feature in each class.



# Summary

Naïve Bayes classifier:

- ✧ What's the assumption
- ✧ Why we make it
- ✧ How we learn it

Naïve Bayes for discrete data

Gaussian naïve Bayes