

---

# Probability Theory Primer plus Naïve Bayes all ways



James G. Shanahan <sup>1,2</sup>

<sup>1</sup>Church and Duncan Group,

<sup>2</sup>*School of Informatics, Computing and Engineering, Indiana University*

*EMAIL: James\_DOT\_Shanahan\_AT\_gmail\_DOT\_com*

# References and Resources

---

- **Manning, Raghavan, Schutz, IRBook**
- **Sebastian Raschka - Naive Bayes & Text Classification**
  - [http://sebastianraschka.com/Articles/2014\\_naive\\_bayes\\_1.html](http://sebastianraschka.com/Articles/2014_naive_bayes_1.html)
  - <https://onlinecourses.science.psu.edu/stat200/node/42>

# Probability Basics → Naïve Bayes Models

- **Introduction**
- **Probability Basics**
  - Probability Axioms
  - Conditional probabilities
  - Product Rule, Chain Rule, Bayes Rule
- **Bayes Nets And Naïve Bayes**
  - Learning
  - Independence
  - Conditional independence
  - Naïve Bayes derivation (discrete case plus smoothing)
- **Naïve Bayes Flavors**
  - Continuous input variables
  - Discrete input variables (2 flavors: Bernoulli, multinomial)
- **Case Study: Spam detector in Naïve Bayes**
- **Summary**

# 0-1 Loss function for Naïve Bayes

---

- The 0-1 loss function penalizes misclassification, i.e. it assigns the smallest loss to the solution that has greatest number of correct classifications. So in both cases we are talking about estimating mode. Recall that mode is the *most common value* in the dataset, or the *most probable value*, so both maximizing the posterior probability and minimizing the 0-1 loss leads to estimating the mode.
- If you need a formal proof, the one is given in the *Introduction to Bayesian Decision Theory* paper by Angela J. Yu:

<https://stats.stackexchange.com/questions/296014/why-is-the-naive-bayes-classifier-optimal-for-0-1-loss>

# Supervised ML: Build Classifiers

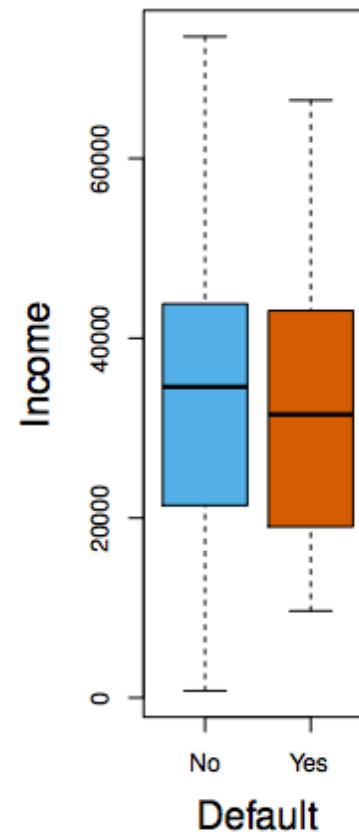
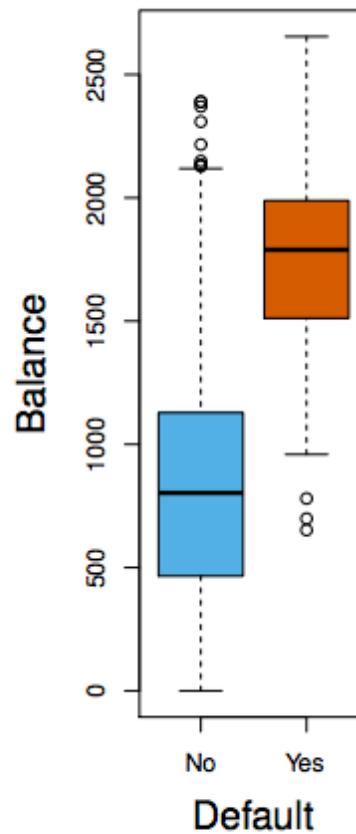
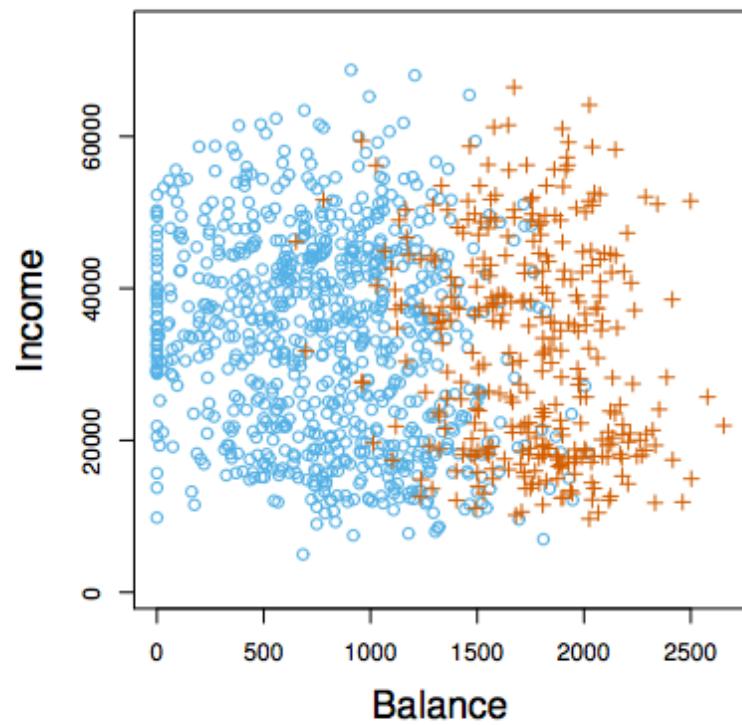
---

- Qualitative variables take values in an unordered set  $\mathcal{C}$ , such as:  
 $\text{eye color} \in \{\text{brown}, \text{blue}, \text{green}\}$   
 $\text{email} \in \{\text{spam}, \text{ham}\}.$
- Given a feature vector  $X$  and a qualitative response  $Y$  taking values in the set  $\mathcal{C}$ , the classification task is to build a function  $C(X)$  that takes as input the feature vector  $X$  and predicts its value for  $Y$ ; i.e.  $C(X) \in \mathcal{C}$ .
- Often we are more interested in estimating the *probabilities* that  $X$  belongs to each category in  $\mathcal{C}$ .

For example the SPAMiness of an email

# Exploratory Data Analysis

## Example: Credit Card Defualt



# **From complex class conditional joint probability distributions of the order $2^n$ to n**

- **From complex class conditional joint probability distributions of the order  $2^n$  to n**
- **$2^n$  Possibilities -> n combinations that we have to estimate class conditional probabilities for**

# Summary

---

## ■ Bayesian prediction:

- requires solving density estimation problems.
- often difficult to estimate  $\Pr[\mathbf{x} \mid y]$  for  $\mathbf{x} \in \mathbb{R}^N$ .
- but, simple and easy to apply; widely used.

## ■ Naive Bayes:

- strong assumption.
- straightforward estimation problem.
- specific linear classifier.
- sometimes surprisingly good performance.

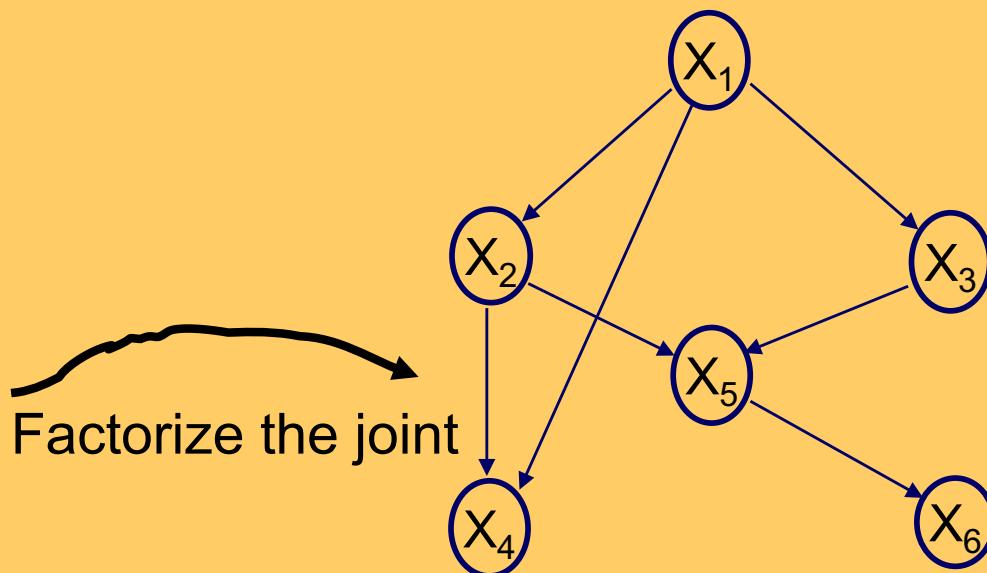
# Decompose Bayesian Networks

GOAL:  $2^N$  Possibilities  $\rightarrow$   $1^N$

P: Joint Probability Distribution

#	X1	X2	X3	X4	X5	X6	Pr(X1,X2, X3, X4, X5, X6)
1	1						
2	1						
..	1						
4..	1						
5	1						
6	1						
7	1						
64							

G: Directed Acyclic Graph



$$p(x_1, x_2, x_3, x_4, x_5, x_6)$$

**1. Partial Order**

$$= p(x_1) p(x_2 | x_1) p(x_3 | x_1, x_2) p(x_4 | x_1, x_2, x_3) p(x_5 | x_4, x_3, x_2, x_1) p(x_6 | x_5, x_4, x_3, x_2, x_1)$$

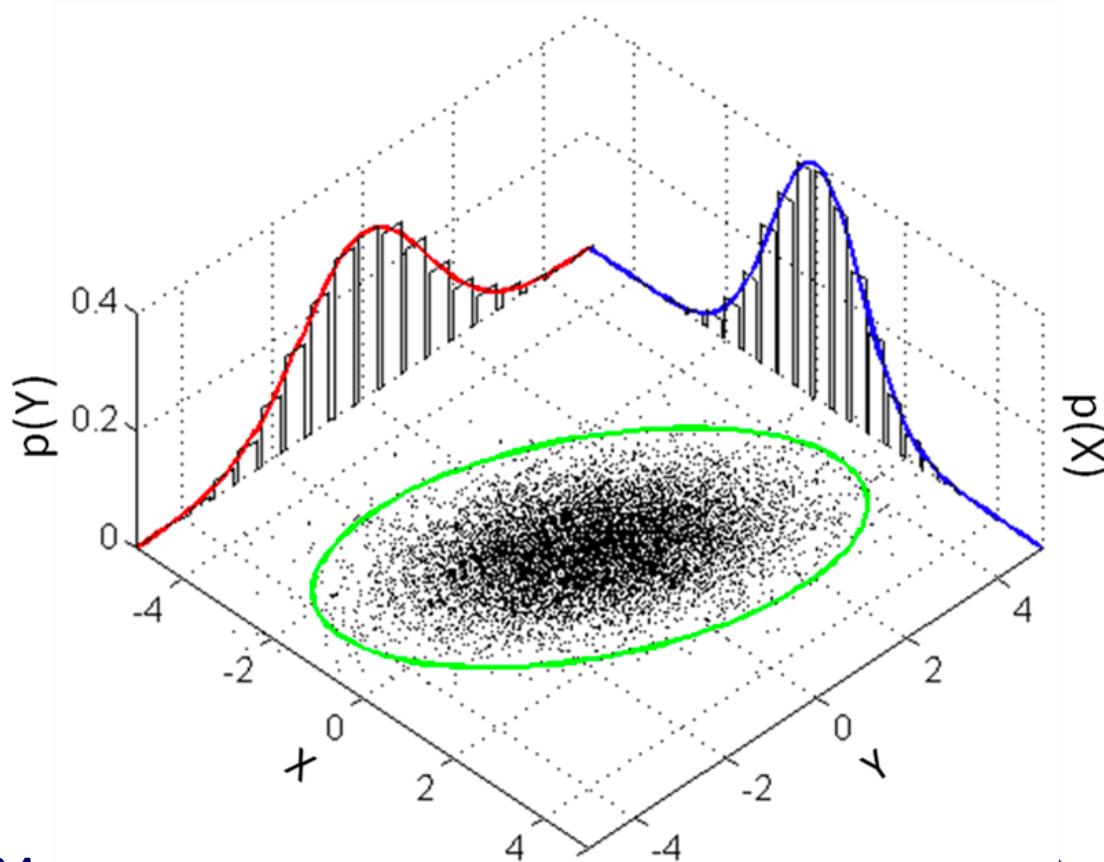
$$= p(x_1) p(x_2 | x_1) p(x_3 | x_1) p(x_4 | x_2, x_1) p(x_5 | x_3, x_2) p(x_6 | x_5) \quad \text{3. Markov Property}$$

$$= p(x_1) p(x_2) p(x_3) p(x_4) p(x_5) p(x_6) \quad \text{4: Independence (see next section)}$$

$$P(y|x_6, x_5, x_4, x_3, x_2, x_1) = p(x_1|y) p(x_2|y) p(x_3|y) p(x_4|y) p(x_5|y) p(x_6|y) \quad \text{5: Naïve Bayes via Cond. Independence}$$

# Decompose large joint distributions

- Decompose large joint distributions of many variables into production many univariate distributions



# Probability Basics → Naïve Bayes Models

- **Introduction**
- **Probability Basics**
  - Probability Axioms
  - Conditional probabilities
  - Product Rule, Chain Rule, Bayes Rule
- **Bayes Nets And Naïve Bayes**
  - Learning
  - Independence
  - Conditional independence
  - Naïve Bayes derivation (discrete case plus smoothing)
- **Naïve Bayes Flavors**
  - Continuous input variables
  - Discrete input variables (2 flavors: Bernoulli, multinomial)
- **Case Study: Spam detector in Naïve Bayes**
- **Summary**

# Probability theory 1/3

---

- **Probability theory is the branch of mathematics concerned with probability, the analysis of random phenomena.**
- **The central objects of probability theory are random variables, stochastic processes, and events:**
  - mathematical abstractions of non-deterministic events or measured quantities that may either be single occurrences or evolve over time in an apparently random fashion.

# Probability theory 2/3

---

- **It is not possible to predict precisely results of random events.**
- **However, if a sequence of individual events, such as coin flipping or the roll of dice, is influenced by other factors, such as friction, it will exhibit certain patterns, which can be studied and predicted.**
- **Two representative mathematical results describing such patterns are the**
  - (1) law of large numbers (LoLN) and the (see next slide)
  - (2) central limit theorem. (central tendencies) (see next slide)

# CLT and LoLN

---

- **Both the law of large numbers and central limit theorem are about many independent samples from same distribution**
- **CLT is about shape: shape is bell-like the more samples we draw**
  - The Central Limit Theorem tell us that as the sample size tends to infinity, the of the distribution of sample means approaches the normal distribution.
  - This is a statement about the SHAPE of the distribution. A normal distribution is bell shaped so the shape of the distribution of sample means begins to look bell shaped as the sample size increases.
- **The Law of Large Numbers tells us where the center (maximum point) of the bell is located.**
  - Again, as the sample size approaches infinity the center of the distribution of the sample means becomes very close to the population mean.

Population of IQ  
scores, 10-year olds

$$\mu = 100$$

$$\sigma = 16$$

$$n = 64$$

Sample  
1

$$\bar{X}_1 = 103.70$$

Sample  
2

$$\bar{X}_2 = 98.58$$

Sample  
3

$$\bar{X}_3 = 100.11$$

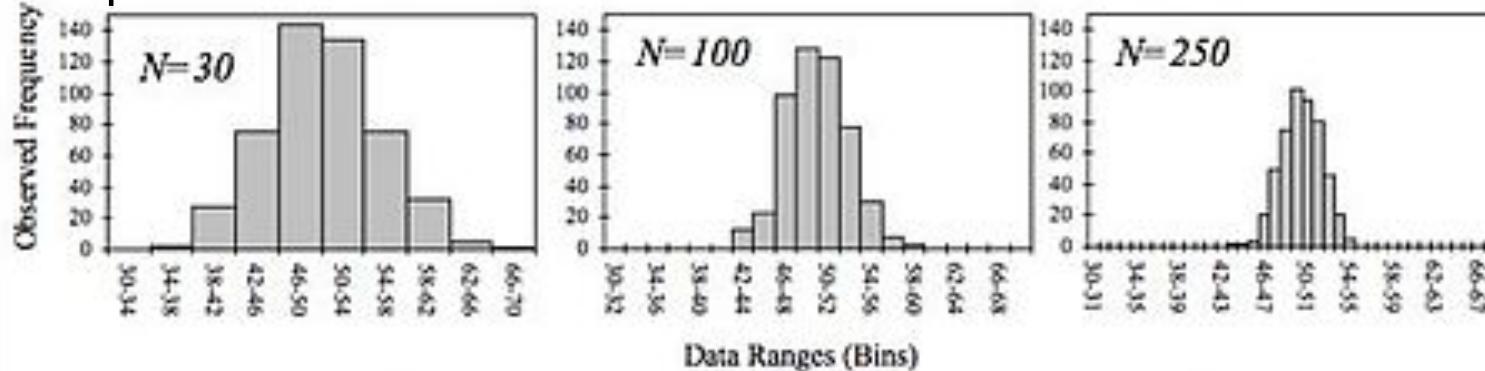
Etc 

Is sample 2 a likely  
representation  
of our population?

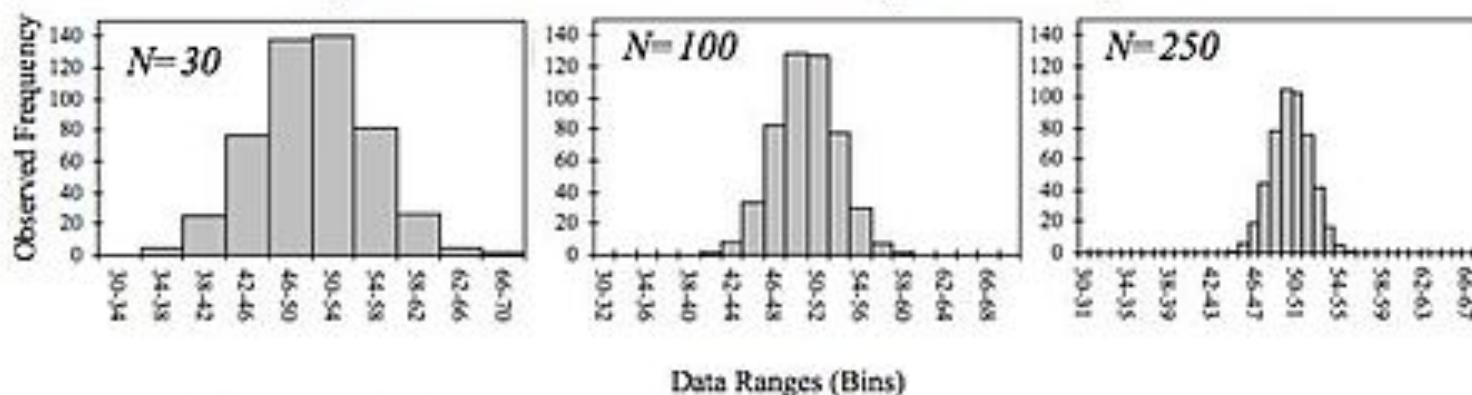
# Uniform distribution [0,100], Pop Mean is 50. vs. Gaussian based population

Histograms of 500 Observed Sample Means Randomly Drawn from a Population (0 to 100) with a Uniform Distribution for Various Sample Sizes ( $N$ )

Sample Size =  $N$



Histograms of ~500 Expected Values for the Normalized Gaussian Distribution  
Using the Best Estimates from the Sample Data as Input Parameters



$$\tilde{\chi}^2_{n=30} \approx 0.33$$

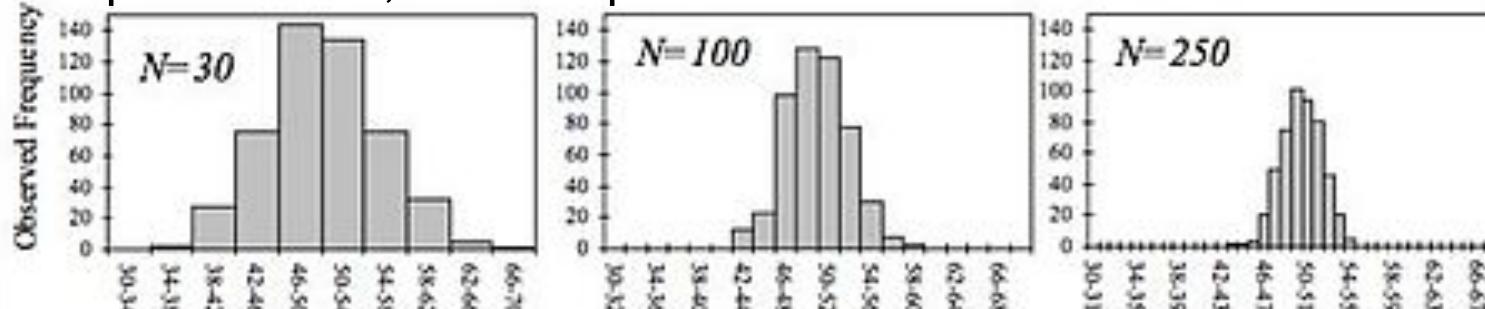
$$\tilde{\chi}^2_{n=100} \approx 0.95$$

$$\tilde{\chi}^2_{n=250} \approx 0.41$$

This figure demonstrates the central limit theorem. The sample means are generated using a random number generator, which draws numbers between 0 and 100 from a uniform probability distribution. It illustrates that increasing sample sizes result in the 500 measured sample means being more closely distributed about the population mean (50 in this case). It also compares the observed distributions with the distributions that would be expected for a normalized Gaussian distribution, and shows the chi-squared values that quantify the goodness of the fit (the fit is good if the reduced chi-squared value is less than or approximately equal to one). The input into the normalized Gaussian function is the mean of sample means (~50) and the mean sample standard deviation divided by the square root of the sample size (~ $28.87/\sqrt{n}$ ), which is called the standard deviation of the mean (since it refers to the spread of sample means).

# Uniform distribution [0,100], Pop Mean is 50. vs. Gaussian based population

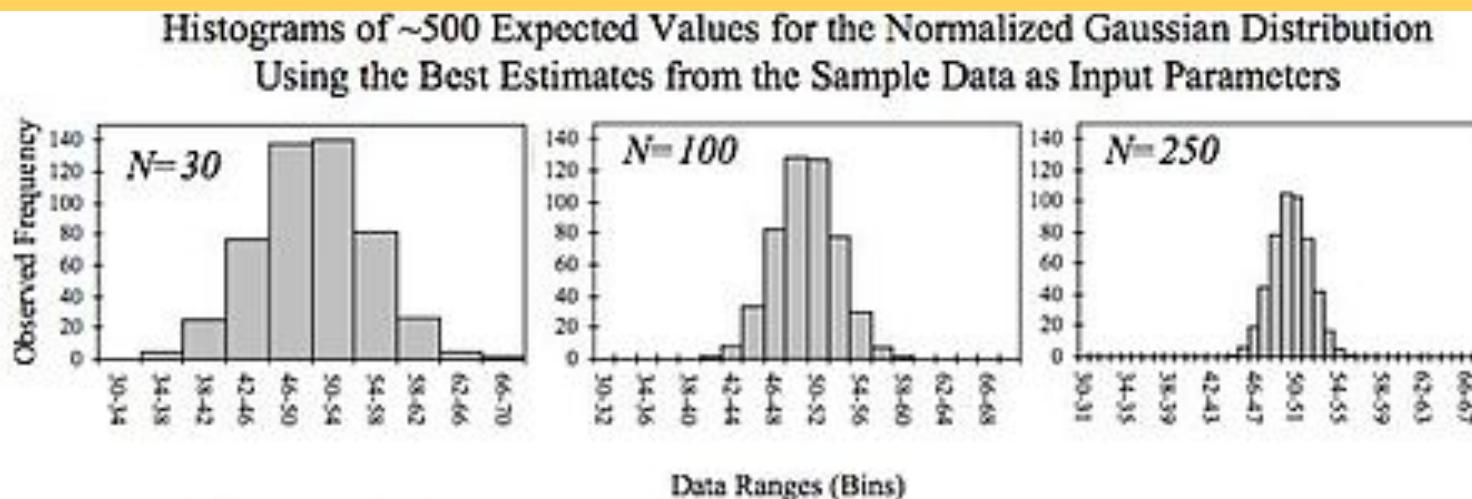
Histograms of 500 Observed Sample Means Randomly Drawn from a Population (0 to 100) with a Uniform Distribution for Various Sample Sizes ( $N$ )  
Sample Size =  $N$ ; 500 samples



This figure demonstrates the central limit theorem. The sample means are generated using a random number generator, which draws numbers between 0 and 100 from a uniform probability distribution. It illustrates that increasing sample sizes result in the 500 measured sample means being more closely distributed about the population mean (50 in this case). It also compares the

Mean of means is getting closer to the population mean: LoLN

The shape of the distribution is gaussian or Bellshape as the sample size  $N$  increases



Gaussian distribution, and shows the chi-squared values that quantify the goodness of the fit (the fit is good if the reduced chi-squared value is less than or approximately equal to one). The input into the normalized Gaussian function is the mean of sample means (~50) and the mean sample standard deviation divided by the square root of the sample size (~ $28.87/\sqrt{n}$ ), which is called the standard deviation of the mean (since it refers to the spread of sample means).

$$\tilde{\chi}^2_{n=30} \approx 0.33$$

$$\tilde{\chi}^2_{n=100} \approx 0.95$$

$$\tilde{\chi}^2_{n=250} \approx 0.41$$

# Probability theory 3/3: many apps

---

- **As a mathematical foundation for statistics, probability theory is essential to many human activities that involve quantitative analysis of large sets of data.**
- **Many applications**
  - Machine learning
  - statistical mechanics: Methods of probability theory also apply to descriptions of complex systems given only partial knowledge of their state, as in statistical mechanics.
  - quantum mechanics: A great discovery of twentieth century physics was the probabilistic nature of physical phenomena at atomic scales, described in quantum mechanics.

# law of large numbers

---

- In probability theory, the law of large numbers (LLN) is a theorem that describes the result of performing the same experiment a large number of times.
- According to the law, the average of the results obtained from a large number of trials should be close to the expected value, and will tend to become closer as more trials are performed.

# CLT: mean will be approximately normally distributed

---

- In probability theory, the central limit theorem (CLT) states that, given certain conditions, the arithmetic mean of a sufficiently large number of iterates of independent random variables, each with a well-defined expected value and well-defined variance, will be approximately normally distributed, regardless of the underlying distribution.[1][2]
- To illustrate what this means, suppose that a sample is obtained containing a large number of observations, each observation being randomly generated in a way that does not depend on the values of the other observations, and that the arithmetic average of the observed values is computed.
- If this procedure is performed many times, the central limit theorem says that the computed values of the average will be distributed according to the normal distribution (commonly known as a "bell curve").
- A simple example of this is that if one flips a coin many times, the probability of getting a given number of heads should follow a normal curve, with mean equal to half the total number of flips.

# Sampling Distributions of Sample Statistics

---

- **Sample Statistics**
  - Two common statistics are the sample proportion,  $\hat{p}$  ("p-hat"), and the sample mean,  $\bar{x}$  ("x-bar").
  - Sample statistics are random variables because they vary from sample to sample.
- **Sampling distribution**
  - As a result, sample statistics also have a distribution called the sampling distribution.
  - These sampling distributions, similar to distributions discussed previously, have a mean and standard deviation.
- **Standard error**
  - We refer to the standard deviation of a sampling distribution as the standard error. Thus, the standard error is simply the standard deviation of a sampling distribution. Often times statisticians will interchange these two terms.

<https://onlinecourses.science.psu.edu/stat200/node/42>

# Standard Error

---

- The standard error is the standard deviation of the sampling distribution of a statistic.<sup>[1]</sup>
- The term may also be used to refer to an estimate of that standard deviation, derived from a particular sample used to compute the estimate.
- For example, the sample mean is the usual estimator of a population mean. However, different samples drawn from that same population would in general have different values of the sample mean. The standard error of the mean (i.e., of using the sample mean as a method of estimating the population mean) is the standard deviation of those sample means over all possible samples (of a given size) drawn from the population. Secondly, the standard error of the mean can refer to an estimate of that standard deviation, computed from the sample of data being analyzed at the time.

Population of IQ  
scores, 10-year olds

$$\mu = 100$$

$$\sigma = 16$$

$$n = 64$$

Sample  
1

$$\bar{X}_1 = 103.70$$

Sample  
2

$$\bar{X}_2 = 98.58$$

Sample  
3

$$\bar{X}_3 = 100.11$$

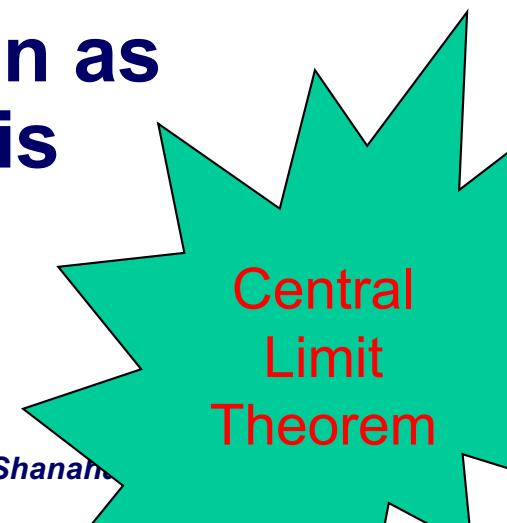
Etc 

Is sample 2 a likely  
representation  
of our population?

# Distribution of Sample Means

---

1. The mean of a sampling distribution is identical to mean of raw scores in the population ( $\mu$ )
2. If the population is Normal, the distribution of sample means is also Normal
3. If the population is not Normal, the distribution of sample means approaches Normal distribution as the size of sample on which it is based gets larger



Central  
Limit  
Theorem

# Standard Error of the Mean

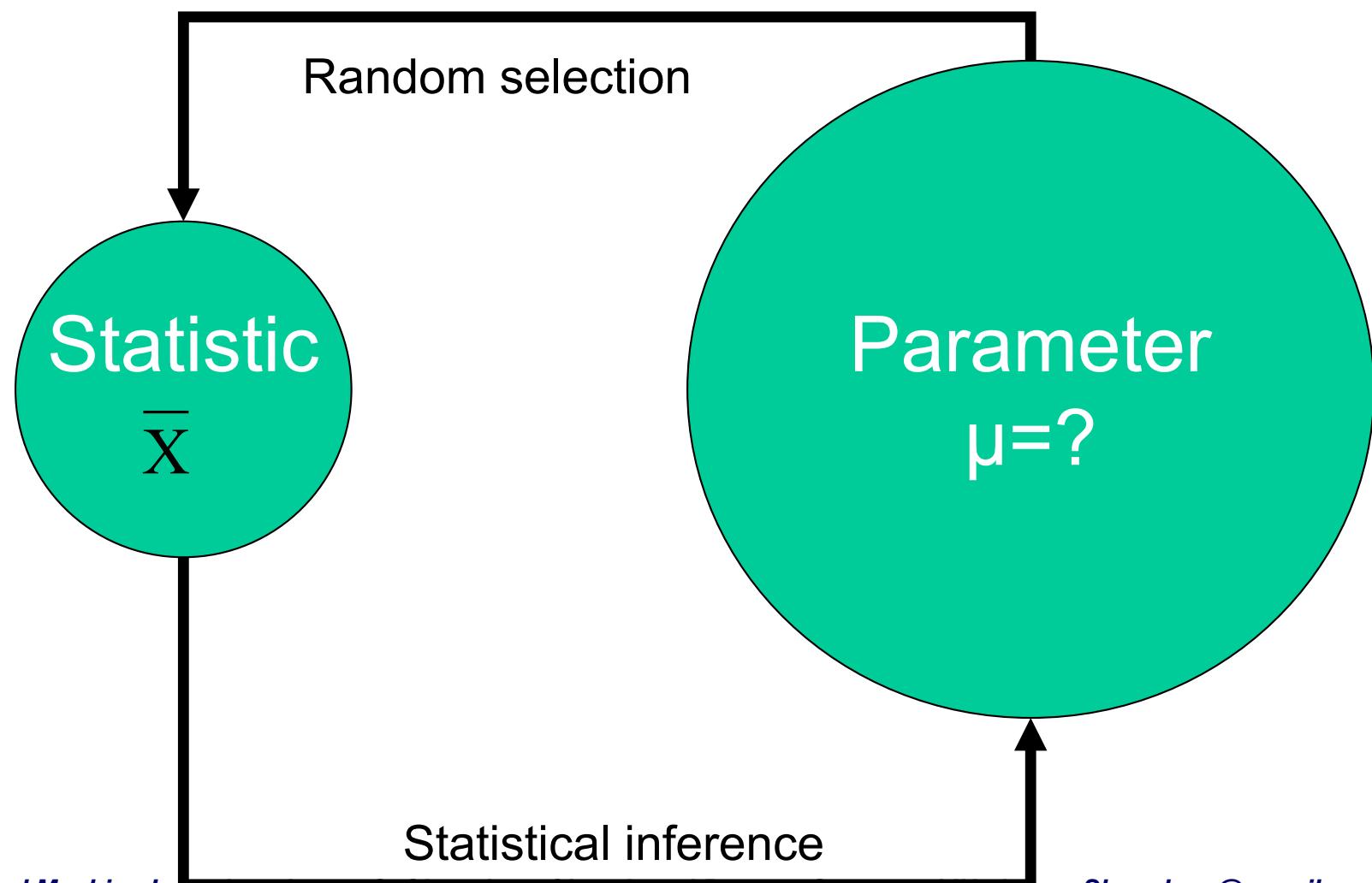
---

- The standard deviation of means in a sampling distribution is known as the **standard error of the mean**.
- It can be calculated from the standard deviation of observations.

$$S_{\bar{X}} = \frac{s}{\sqrt{n}}$$

$$\text{Standard Error of } \bar{X} : \sigma_{\bar{X}} = \frac{\sigma}{\sqrt{n}}$$

3. The larger our sample size, the smaller our standard error



## STANDARD ERROR CALCULATION

### Procedure:

Step 1: Calculate the mean (Total of all samples divided by the number of samples).

Step 2: Calculate each measurement's deviation from the mean (Mean minus the individual measurement).

Step 3: Square each deviation from mean. Squared negatives become positive.

Step 4: Sum the squared deviations (Add up the numbers from step 3).

Step 5: Divide that sum from step 4 by one less than the sample size ( $n-1$ , that is, the number of measurements minus one)

Step 6: Take the square root of the number in step 5. That gives you the "standard deviation (S.D.)."

Step 7: Divide the standard deviation by the square root of the sample size ( $n$ ). That gives you the "standard error".

Step 8: Subtract the standard error from the mean and record that number. Then add the standard error to the mean and record that number. You have plotted  $\text{mean} \pm 1$  standard error (S. E.), the distance from 1 standard error below the mean to 1 standard error above the mean

### Example:

Name	Height to nearest 0.5 cm	2 Deviations ( $m-i$ )	3 Squared deviations ( $(m-i)^2$ )
1. Waldo	150.5	11.9	141.61
2. Finn	170.0	-7.6	57.76
3. Henry	160.0	2.4	5.76
4. Alfie	161.0	1.4	1.96
5. Shane	170.5	-8.1	65.61
$n=5$	<b>1 Mean <math>m = 162.4</math> cm</b>		<b>4 Sum of squared deviations <math>\sum(m-i)^2 = 272.70</math></b>

**5** Divide by number of measurements-1.  $\sum (m-i)^2 / (n-1) = 272.70 / 4 = 68.175$

**6 Standard deviation** = square root of  $\sum (m-i)^2 / n-1 = \sqrt{68.175} = 8.257$

**7 Standard error** = Standard deviation/ $\sqrt{n}$  =  $8.257 / \sqrt{5} = 3.69$

**8**  $m \pm 1\text{SE} = 162 \pm 3.7$  or 159cm to 166cm for the men ( $162.4 - 3.7$  to  $162.4 + 3.7$ ).

Draw a sample of size N from the population.

In the SE calculation N is the sample size (not the number of samples drawn)

# Estimation Procedures

---

- **Point estimates**
  - For example mean of a sample of 25 patients
    - No information regarding probability of accuracy
  - Interval estimates
  - Estimate a range of values that is likely
    - Confidence interval between two limit values
      - The degree of confidence depends on the probability of including the population mean

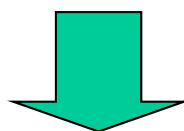
$$95\% \text{ CI} = \bar{X} \pm 1.96 \frac{\sigma}{\sqrt{n}}$$

$$99\% \text{ CI} = \bar{X} \pm 2.58 \frac{\sigma}{\sqrt{n}}$$

# When Sample size is small ...

---

$$95\% \text{ CI} = \bar{X} \pm 1.96 \frac{\sigma}{\sqrt{n}}$$



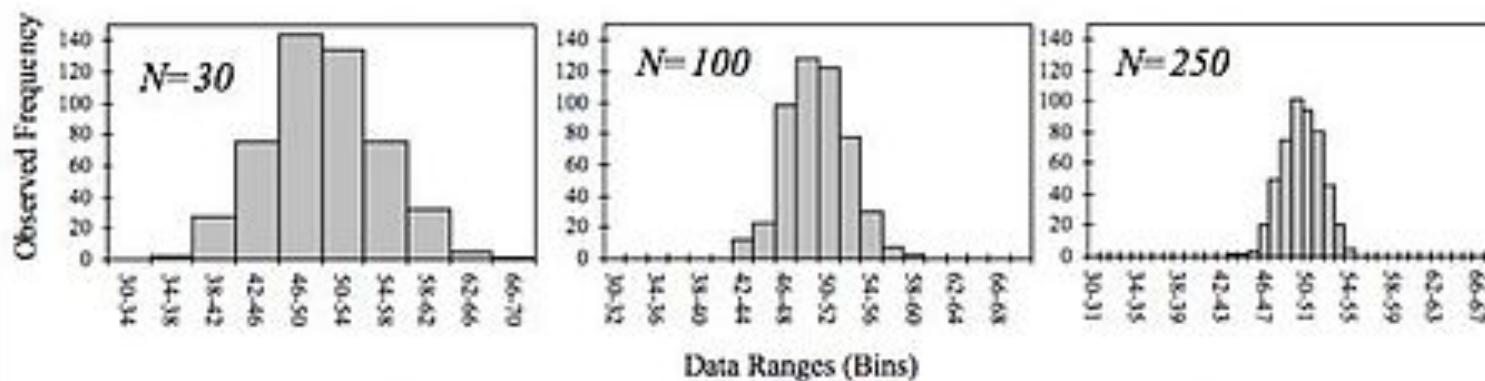
$$95\% \text{ CI} = \bar{X} \pm t \frac{S}{\sqrt{n}}$$



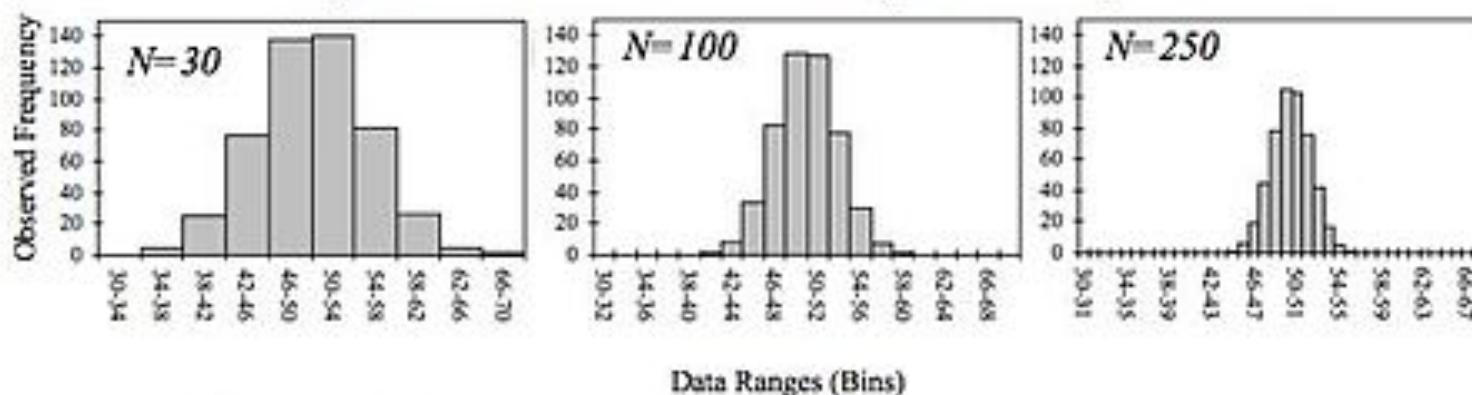
A constant from  
Student t Distribution  
that depends on confidence  
interval and sample size

# Uniform distribution [0,100], Pop Mean is 50. vs. Gaussian based population

Histograms of 500 Observed Sample Means Randomly Drawn from a Population (0 to 100) with a Uniform Distribution for Various Sample Sizes ( $N$ )



Histograms of ~500 Expected Values for the Normalized Gaussian Distribution  
Using the Best Estimates from the Sample Data as Input Parameters



$$\tilde{\chi}^2_{n=30} \approx 0.33$$

$$\tilde{\chi}^2_{n=100} \approx 0.95$$

$$\tilde{\chi}^2_{n=250} \approx 0.41$$

This figure demonstrates the central limit theorem. The sample means are generated using a random number generator, which draws numbers between 0 and 100 from a uniform probability distribution. It illustrates that increasing sample sizes result in the 500 measured sample means being more closely distributed about the population mean (50 in this case). It also compares the observed distributions with the distributions that would be expected for a normalized Gaussian distribution, and shows the chi-squared values that quantify the goodness of the fit (the fit is good if the reduced chi-squared value is less than or approximately equal to one). The input into the normalized Gaussian function is the mean of sample means (~50) and the mean sample standard deviation divided by the square root of the sample size (~ $28.87/\sqrt{n}$ ), which is called the standard deviation of the mean (since it refers to the spread of sample means).

# Notation

---

- Proposition - statement or assertion about a state of the world
- Variable  $X$  is a set of mutually exclusive propositions  $x_i$
- Variables – upper-case
- Propositions – lowercase
  - Example ( $X=x$ ,  $Y=y$ ,  $Z=z$ )
  - Shortened:  $(x,y,z)$
- Sets of variables – bold
  - Example:  $(X, Y, Z)$
- Latent/Hidden variable – states are inferred but never observed directly

# From Axioms: deduce theorems and propositions

---

- **Math**
  - One strategy in mathematics is to start with a few statements, then build up more mathematics from these statements.
  - The beginning statements are known as axioms.
    - An axiom is typically something that is mathematically self evident.
    - From a relatively short list of axioms, deductive logic is used to prove other statements, called theorems or propositions.
- **Probability Theory has 3 axioms → theorems or propositions**
  - The area of mathematics known as probability is no different.
  - Underlying probability is a handful of axioms from which we can derive all sorts of results. But what are these probability axioms?
  - Probability can be reduced to three axioms.
  - It presupposes that we have a set of outcomes called the sample space  $S$  comprised of subsets called events  $E_1, E_2, \dots, E_n$  and a way of assigning a probability to any event  $E$ . The probability of the event  $E$  is denoted by  $P(E)$ .

# Axioms of Probabilities

---

<http://statistics.about.com/od/Mathstat/a/what-is-the-power-set.htm>

Axiom 1 :  $0 \leq P(A) \leq 1$

Axiom 2:  $P(\text{Sure Proposition}) = 1$

Axiom 3:  $P(A \text{ or } B) = P(A) + P(B) - P(A \text{ and } B)$

Marginal probability:  $P(A) = P(A, B) + P(A, \neg B)$

$$P(A) = \sum_i P(A, B_i)$$

The first axiom of probability is that the probability of any event is a nonnegative real number. This means that the smallest that a probability can ever be is zero, and that it cannot be infinite.

The third axiom of probability deals with mutually exclusive events. If E1 and E2 are mutually exclusive, meaning that they have an empty intersection and we use U to denote the union, then  $P(E1 \cup E2) = P(E1) + P(E2)$ .

# Axiom Applications: Pr(an impossible event)

---

The three axioms set an upper bound for the probability of any event. We denote the complement of the event  $E$  by  $E^C$ . From set theory  $E$  and  $E^C$  have empty intersection and are mutually exclusive. Furthermore  $E \cup E^C = S$ , the entire sample space.

These facts, combined with the axioms give us:

$$1 = P(S) = P(E \cup E^C) = P(E) + P(E^C).$$

We rearrange the above equation and see that  $P(E) = 1 - P(E^C)$ . Since we know that probabilities must be nonnegative, we now have that an upper bound for the probability of any event is 1.

By rearranging the formula again we have  $P(E^C) = 1 - P(E)$ . We also can deduce from this formula that the probability of an event not occurring is one minus the probability that it does occur.

The above equation also provides us a way to calculate the probability of the impossible event, denoted by the empty set. To see this, recall that the empty set is the complement of the universal set, in this case  $S^C$ . Since  $1 = P(S) + P(S^C) = 1 + P(S^C)$ , by algebra we have  $P(S^C) = 0$ .

# Complement space

---

2<sup>nd</sup> Axiom: states that the probability of all the events, i.e., the probability of the entire sample space is 1.

Mathematically, if S represents the Sample space, then  $P(S)=1$ .

This means that there are no events outside the sample space and it includes all possible events in it.

$$P(A) + P(\neg A) = 1$$

$$P(A | B)$$

$$\Pr(A) \rightarrow P(A|B)$$

- Belief in A under the assumption that B is known with absolute certainty
- A is conditioned on B

# Probability Basics → Naïve Bayes Models

- **Introduction**
- **Probability Basics**
  - Probability Axioms
  - Conditional probabilities
  - Product Rule, Chain Rule, Bayes Rule
- **Bayes Nets And Naïve Bayes**
  - Learning
  - Independence
  - Conditional independence
  - Naïve Bayes derivation (discrete case plus smoothing)
- **Naïve Bayes Flavors**
  - Continuous input variables
  - Discrete input variables (2 flavors: Bernoulli, multinomial)
- **Case Study: Spam detector in Naïve Bayes**
- **Summary**

# Sample Space: possible outcomes

- 
- In probability theory, the sample space of an experiment or random trial is the set of all possible outcomes or results of that experiment. A sample space is usually denoted using set notation, and the possible outcomes are listed as elements in the set. It is common to refer to a sample space by the labels  $S$ ,  $\Omega$ , or  $U$  (for "universal set").
  - E.g.,
    - For example, if the experiment is tossing a coin, the sample space is typically the set {head, tail}.
    - For tossing two coins, the corresponding sample space would be {(head,head), (head,tail), (tail,head), (tail,tail)}.
    - For tossing a single six-sided die, the typical sample space is {1, 2, 3, 4, 5, 6} (in which the result of interest is the number of pips facing up).

{head, tail}

{(head,head),  
(head,tail),  
(tail,head),  
(tail,tail)}

# Multiple sample spaces (views)

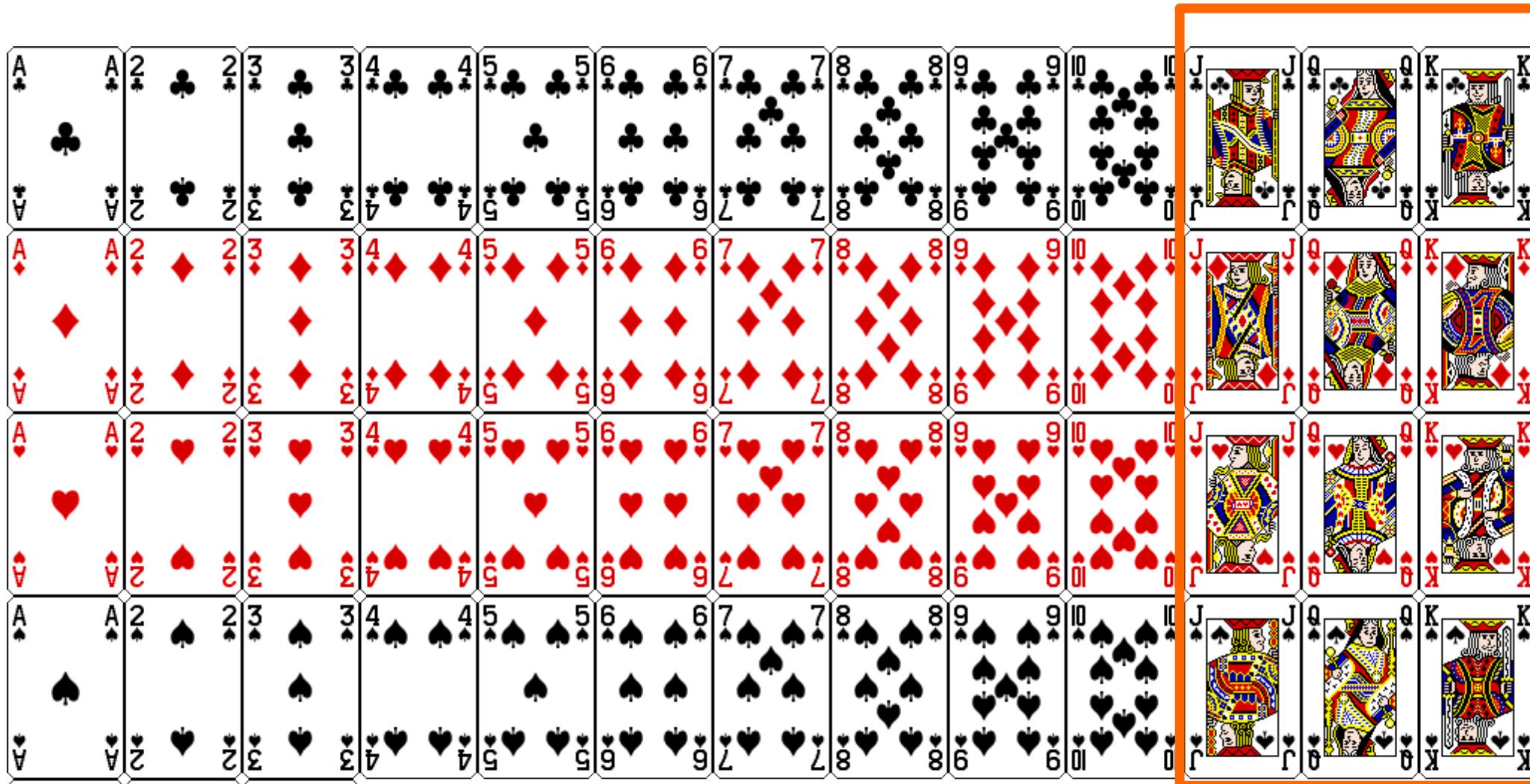
---

- For many experiments, there may be more than one plausible sample space available, depending on what result is of interest to the experimenter.
- For example, when drawing a card from a standard deck of fifty-two playing cards,
  - Rank sample space (picture cards)
    - one possibility for the sample space could be the various ranks (Ace through King),
  - Suits sample space
    - while another could be the suits (clubs, diamonds, hearts, or spades).
  - Cartesian Sample Space
    - A more complete description of outcomes, however, could specify both the denomination and the suit, and a sample space describing each individual card can be constructed as the Cartesian product of the two sample spaces noted above (this space would contain fifty-two equally likely outcomes).
  - Still other sample spaces are possible, such as {right-side up, up-side down} if some cards have been flipped when shuffling.

PictureCards X Suit = 4 X 4 = 16

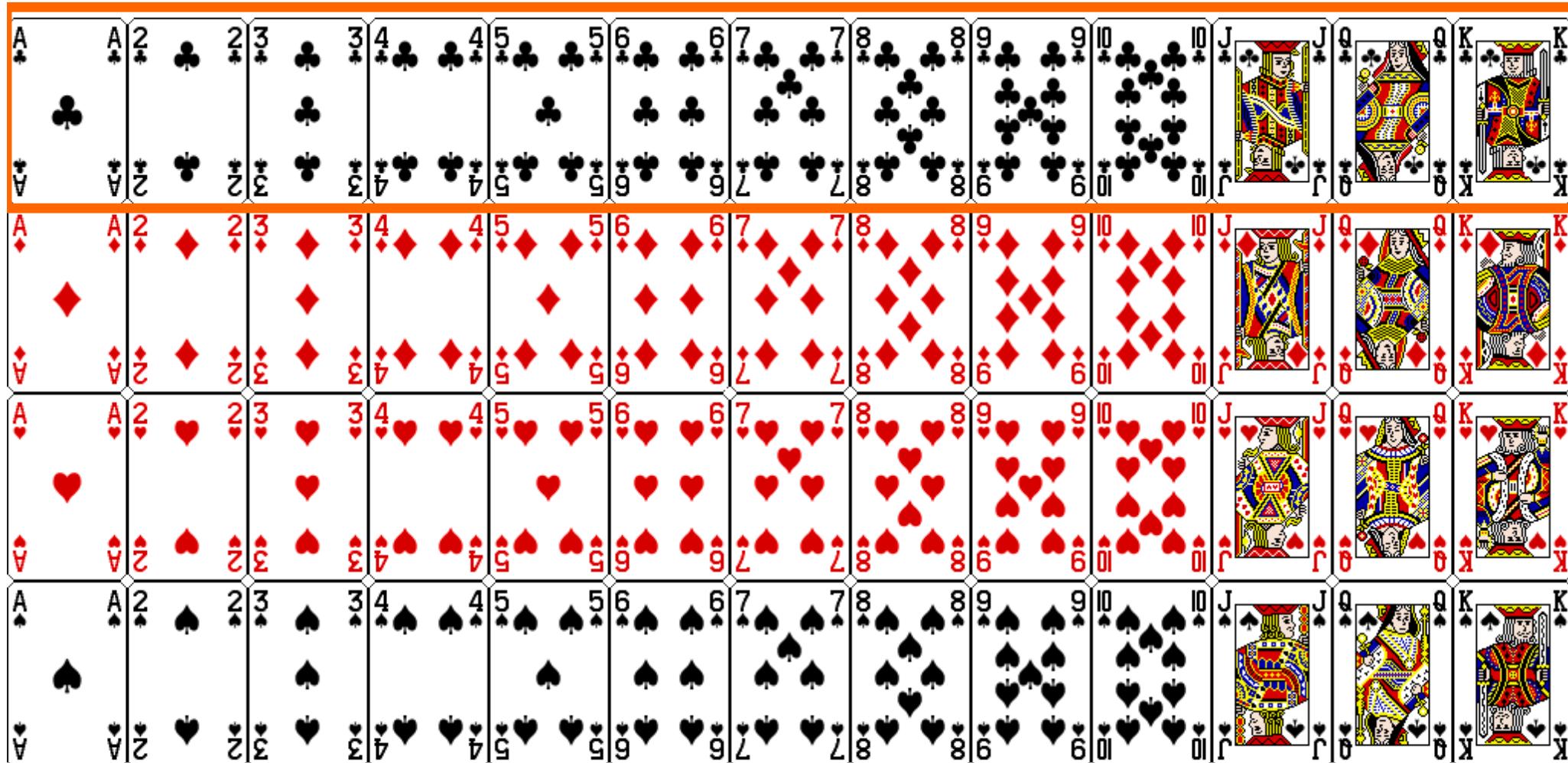
# Sample space

Royal Cards = {Jack, Queen, King}



# Sample space

Clubs

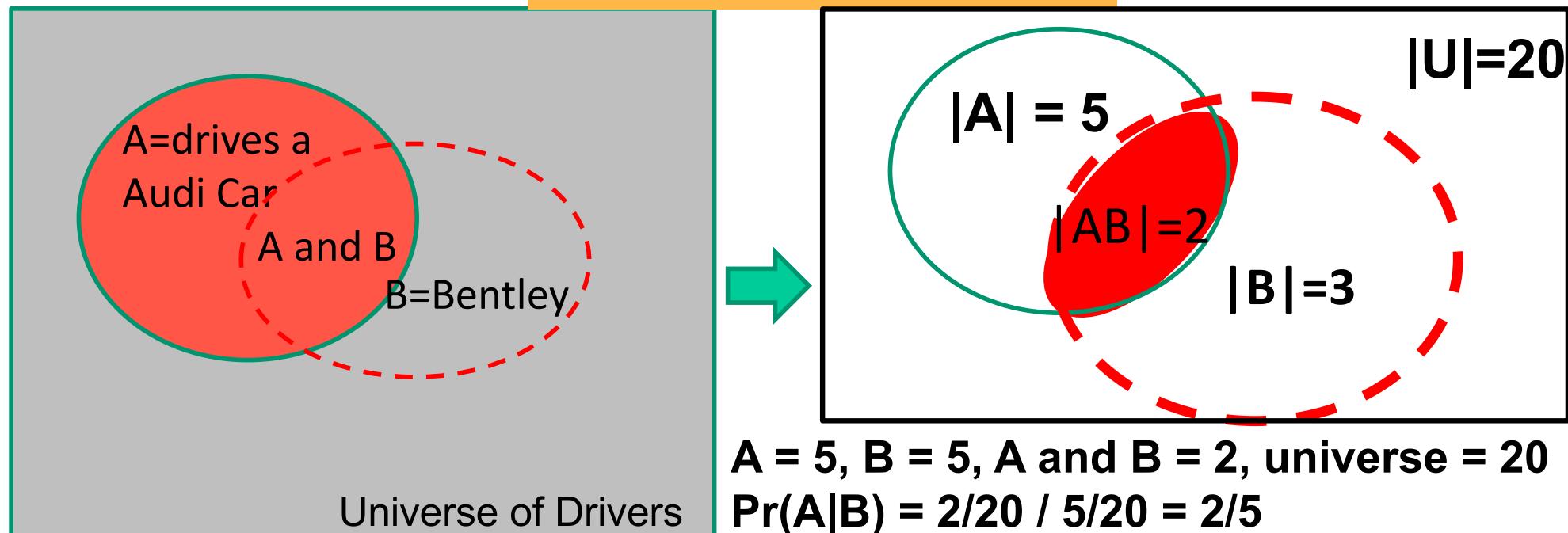


# Conditional Probability as probability wrt a reduced sample space

$$P(A|B) = \frac{P(A \cap B)}{P(B)}$$

Conditional probability can be seen to be the probability with respect to a reduced sample space. We can illustrate the conditional probability with the Venn diagram.

$$\Pr(A) \rightarrow P(A|B)$$



# Calculating Conditional Probability

$$P(A|B) = \frac{P(A \cap B)}{P(B)}$$

Conditional probability can be seen to be the probability with respect to a reduced sample space. We can illustrate the conditional probability with the Venn diagram.

$$\Pr(A) \rightarrow P(A|B) = \frac{P(A \cap B)}{P(B)} = \frac{2}{5}$$

A=drives a Audi Car  
A and B  
B=Bentley

Universe of Drivers

Given: A = 5 B = 5, A and B = 2, universe = 20

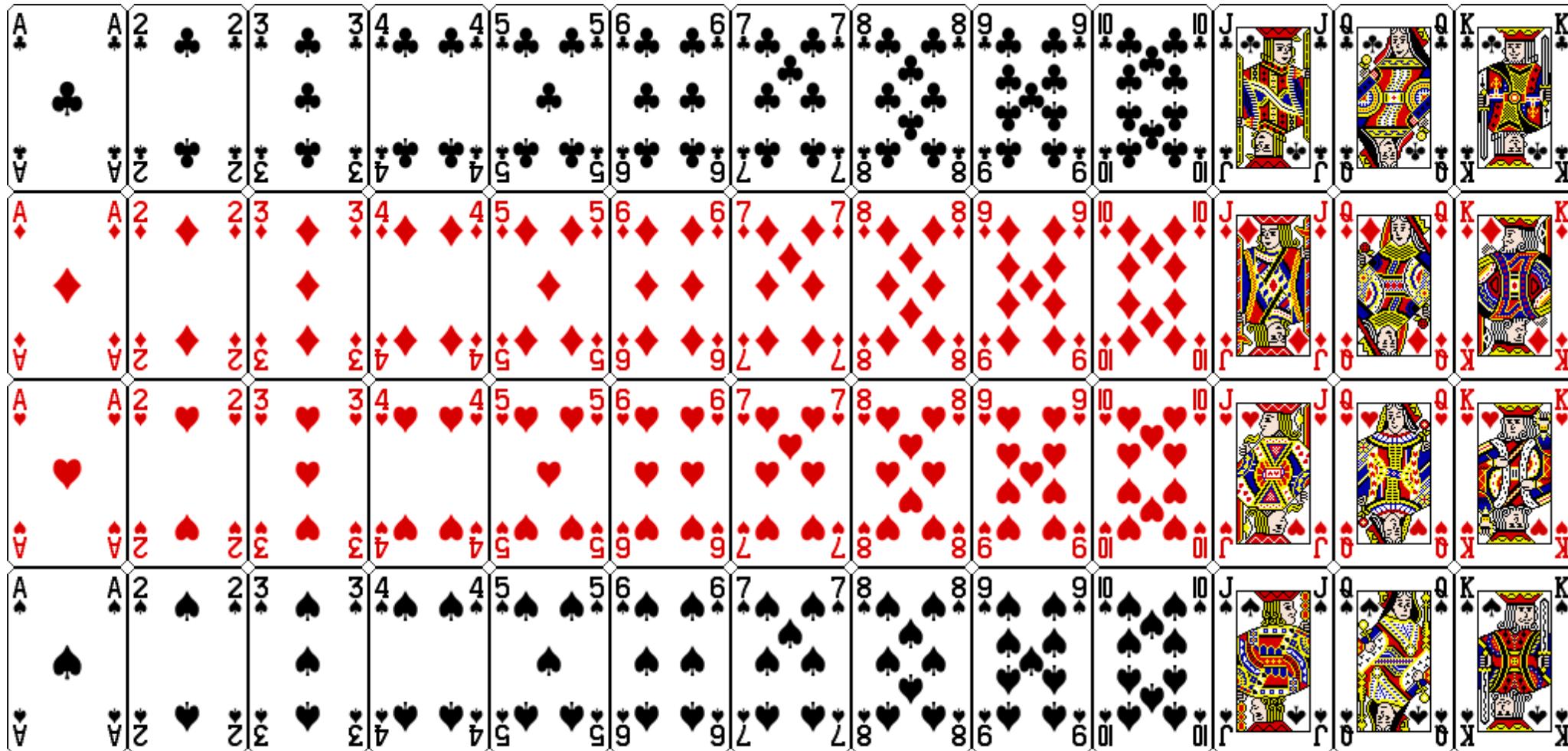
$$\Pr(A|B) = 2/20 / 5/20 = 2/5$$

$$\Pr(A|\text{not } B) = 3/20 / 15/20 = 3/5$$

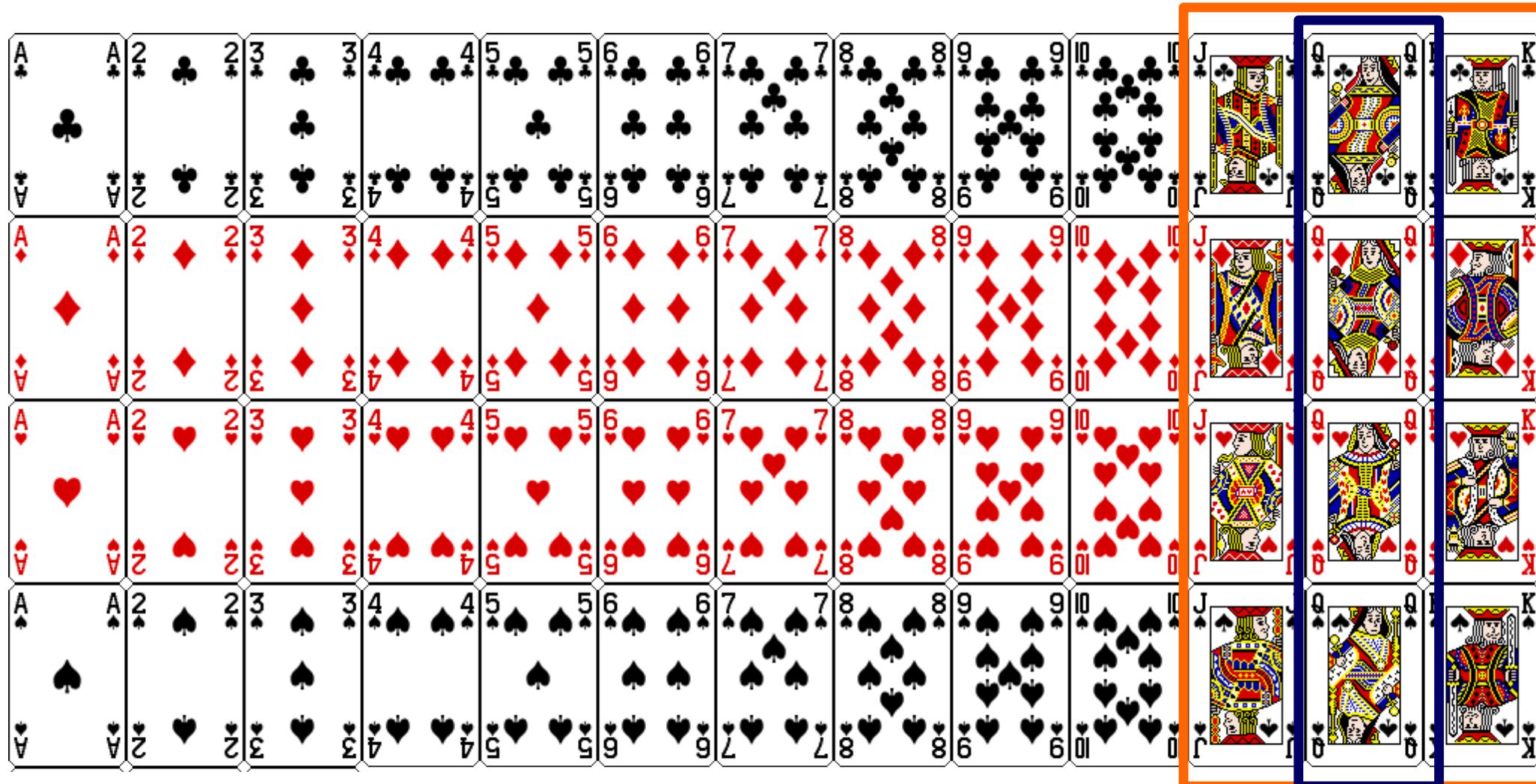
$$\Pr(A) = \Pr(A|B) + \Pr(A|\text{not } B) = 2/5 + 3/5 = 1$$

# Conditional probability example

Royal Cards = {Jack, Queen, King}

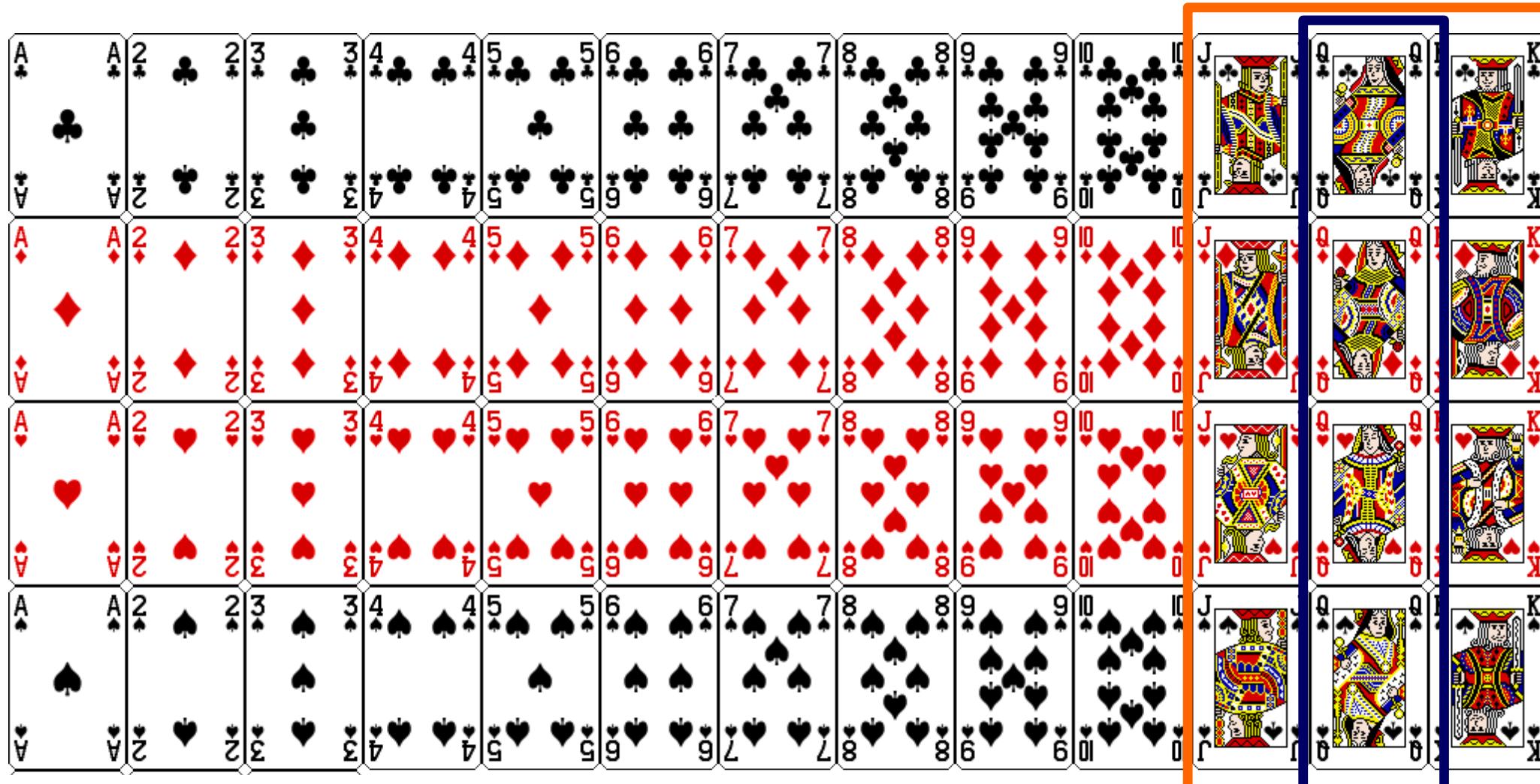


# Conditional probability example



$$P(\text{Queen} | \text{RoyalCard}) = \frac{P(\text{Queen} \cap \text{RoyalCard})}{P(\text{RoyalCard})} = \frac{1/13}{3/13} = \frac{1}{3}$$

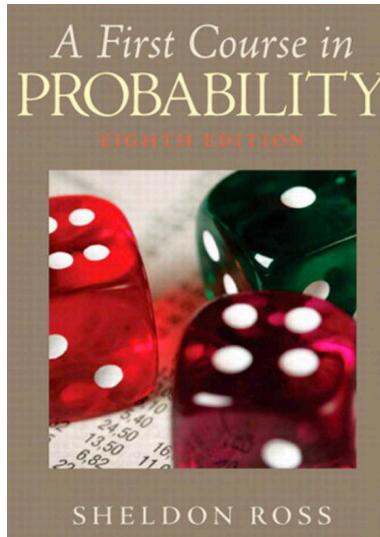
$\Pr(\text{Queen}|\text{Club})$  is independent?



$$P(\text{Queen}|\text{RoyalCard}) = \frac{P(\text{Queen} \cap \text{RoyalCard})}{P(\text{RoyalCard})} = \frac{1/13}{3/13} = \frac{1}{3}$$

# Conditional Probabilities Example

## 3.2 CONDITIONAL PROBABILITIES



Suppose that we toss 2 dice, and suppose that each of the 36 possible outcomes is equally likely to occur and hence has probability  $\frac{1}{36}$ . Suppose further that we observe that the first die is a 3. Then, given this information, what is the probability that the sum of the 2 dice equals 8? To calculate this probability, we reason as follows: Given that the initial die is a 3, there can be at most 6 possible outcomes of our experiment, namely, (3, 1), (3, 2), (3, 3), (3, 4), (3, 5), and (3, 6). Since each of these outcomes originally had the same probability of occurring, the outcomes should still have equal probabilities. That is, given that the first die is a 3, the (conditional) probability of each of the outcomes (3, 1), (3, 2), (3, 3), (3, 4), (3, 5), and (3, 6) is  $\frac{1}{6}$ , whereas the (conditional) probability of the other 30 points in the sample space is 0. Hence, the desired probability will be  $\frac{1}{6}$ .

If we let  $E$  and  $F$  denote, respectively, the event that the sum of the dice is 8 and the event that the first die is a 3, then the probability just obtained is called the *conditional probability that  $E$  occurs given that  $F$  has occurred* and is denoted by

$$P(E|F)$$

A general formula for  $P(E|F)$  that is valid for all events  $E$  and  $F$  is derived in the same manner: If the event  $F$  occurs, then, in order for  $E$  to occur, it is necessary that the actual occurrence be a point both in  $E$  and in  $F$ ; that is, it must be in  $EF$ . Now, since we know that  $F$  has occurred, it follows that  $F$  becomes our new, or reduced, sample space; hence, the probability that the event  $EF$  occurs will equal the probability of  $EF$  relative to the probability of  $F$ . That is, we have the following definition.

58

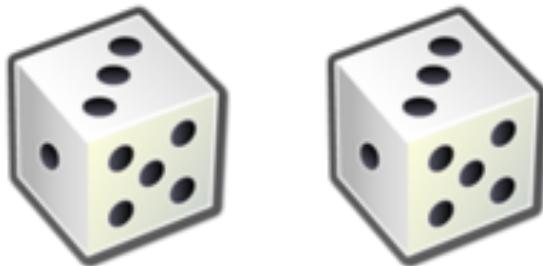
### Definition

If  $P(F) > 0$ , then

$$P(E|F) = \frac{P(EF)}{P(F)} \quad (2.1)$$

# Quiz: Conditional Probabilities Puzzle

Roll a die twice



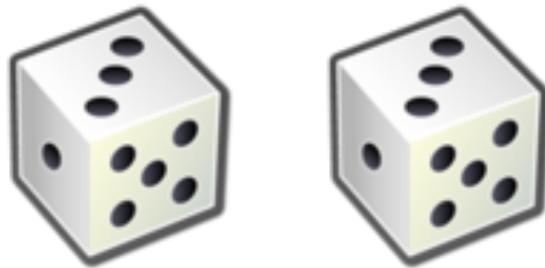
**Q1**

**2 Questions**

Roll a die twice. What is the probability that the total we get is 3? **Q2** Given the information that the first number is 1, what is probability that the total we get is 3?

# Conditional Probabilities Puzzle

Roll a die twice



## Question

Roll a die twice. What is the probability that the total we get is 3? Given the information that the first number is 1, what is probability that the total we get is 3?

Answer:  $2/36$  and  $1/6$ .

2 of  $6 \times 6 = 36$   
equillikely events:  
 $1+2$ ,  $2+1$

Since 1<sup>st</sup> number =1, 2<sup>nd</sup> number must =2 to get a sum of 3. There are 6 possible 2<sup>nd</sup> outcomes, 1 of which is to get a 2, hence  $1/6$

# Conditional Probabilities Example 2

---

Example 1: If the probability that a research project will be well planned is 0.8 and the probability that it will be well planned and well executed is 0.72, what is the probability that a well planned research project will be well executed?

# Conditional Probabilities Examples

---

Example 1: If the probability that a research project will be well planned is 0.8 and the probability that it will be well planned and well executed is 0.72, what is the probability that a well planned research project will be well executed?

Answer:  $0.72/0.8 = 0.9$ .  $P[\text{Exec}|\text{Plan}] = P[\text{Exec} \& \text{Plan}] / P[\text{Plan}]$

# Conditional probability as an easier path to an answer

Sometimes the conditional probability can be determined easily, so we can actually use the conditional probability to calculate probability.

The multiplication rule:  $P(A \cap B) = P(A|B)P(B) = P(B|A)P(A)$ .

Example: There are 3 red balls and 2 blue balls in a box. Randomly take 2 balls from the box. What is the probability that both are red?

Answer:  $P(R_1 \cap R_2) = P(R_1)P(R_2|R_1) = (3/5)(2/4) = 0.3$ .

3 of the 5 are red

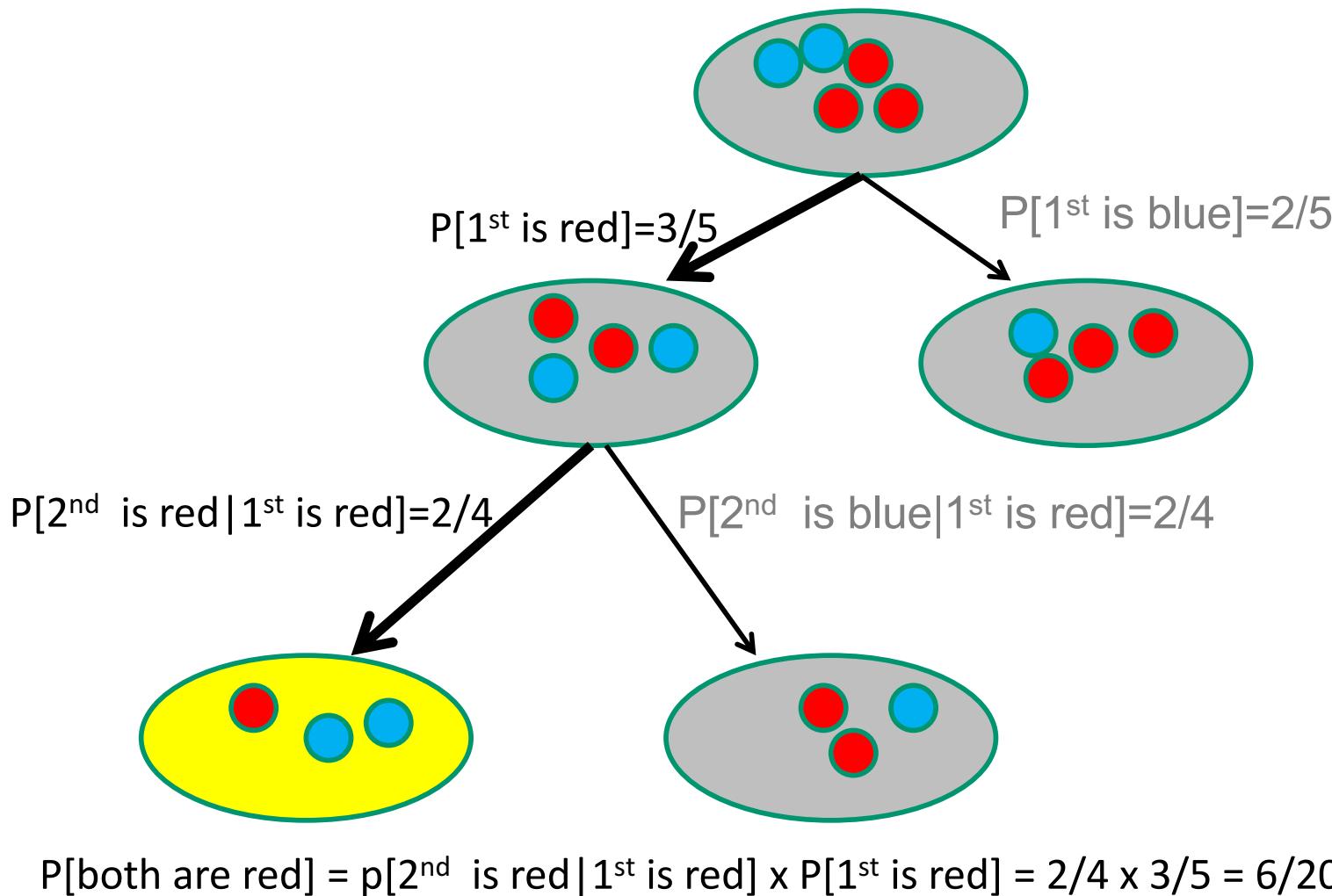
With 1 red ball gone, 2 of the remaining 4 balls are red

Alternative solution:  $P[2 \text{ red in 2 tries}] =$

$$\text{Comb}(3,2)*\text{Comb}(2,0)/\text{Comb}(5,2)$$

$$= 0.3$$

# Solution via Tree Diagram



# Conditional Probability Example

---

The initial intuition for conditional probability comes from considering probabilities that are ratios. In the case of ratios,  $P(E|F)$ , as defined above, is the fraction of items in  $F$  that are also in  $E$ . We show this as follows. Let  $n$  be the number of items in the sample space,  $n_F$  be the number of items in  $F$ , and  $n_{EF}$  be the number of items in  $E \cap F$ . Then

$$\frac{P(E \cap F)}{P(F)} = \frac{n_{EF}/n}{n_F/n} = \frac{n_{EF}}{n_F},$$

which is the fraction of items in  $F$  that are also in  $E$ . As far as meaning,  $P(E|F)$  means the probability of  $E$  occurring given that we know  $F$  has occurred.

**Example 1.6** Again consider drawing the top card from a deck of cards, let Queen be the set of the 4 queens, RoyalCard be the set of the 12 royal cards, and Spade be the set of the 13 spades. Then

$$P(\text{Queen}) = \frac{1}{13}$$

$$P(\text{Queen}|\text{RoyalCard}) = \frac{P(\text{Queen} \cap \text{RoyalCard})}{P(\text{RoyalCard})} = \frac{1/13}{3/13} = \frac{1}{3}$$

$$P(\text{Queen}|\text{Spade}) = \frac{P(\text{Queen} \cap \text{Spade})}{P(\text{Spade})} = \frac{1/52}{1/4} = \frac{1}{13}.$$

# Probability Basics → Naïve Bayes Models

- **Introduction**
- **Probability Basics**
  - Probability Axioms
  - Conditional probabilities
  - Product Rule, Chain Rule, Bayes Rule
- **Bayes Nets And Naïve Bayes**
  - Learning
  - Independence
  - Conditional independence
  - Naïve Bayes derivation (discrete case plus smoothing)
- **Naïve Bayes Flavors**
  - Continuous input variables
  - Discrete input variables (2 flavors: Bernoulli, multinomial)
- **Case Study: Spam detector in Naïve Bayes**
- **Summary**

# Probability Theory Rules

- **cond<sup>al</sup> probability**  $P(A | B) = \frac{P(A, B)}{P(B)}$
- **PRODUCT RULE and Chain Rule:**

$$P(A, B) = P(A | B)P(B)$$

$$P(A, B, C, D) = P(A | B, C, D)P(B | C, D)P(C | D)P(D)$$

- **Total Probabilities :**
  - $\Pr(A) = \Pr(A) = \Pr(A|B) + \Pr(A|\text{not } B)$
- **Bayes Rules**

Posterior      Likelihood      Prior

$$P(A | B) = \frac{P(B | A)P(A)}{P(B)}$$

$$P(A | B) = \frac{P(B | A)P(A)}{\sum_{i=1}^n P(B | A = a_i)P(A = a_i)}$$

# Product Rule is Fundamental and follows from cond<sup>al</sup> probability

$$P(Queen \mid Club) = \frac{P(Queen, Club)}{P(Club)} = \frac{\frac{1}{52}}{\frac{13}{52}} = \frac{1}{13}$$

Cond<sup>al</sup> Prob:

$$P(A \mid B) = \frac{P(A, B)}{P(B)}$$

Chain Rule, Bayes Rule  
follow from the Product  
Rule

WHERE  $P(A, B)$  is the belief in the joint event of A and B  
(joint probability )

$P(B)$  is the marginal probability of B

Rewrite the cond<sup>al</sup> probability to give us the product rule

$$P(A, B) = P(A \mid B)P(B)$$

PRODUCT RULE :

Divide and conquer

# Product Rule Part 2

$$P(Queen \mid Club) = \frac{P(Queen, Club)}{P(Club)} = \frac{\frac{1}{52}}{\frac{13}{52}} = \frac{1}{13}$$

$$P(A \mid B) = \frac{P(A, B)}{P(B)}$$

Chain Rule, Bayes Rule  
follow from the Product  
Rule

PRODUCT RULE :

$$P(A, B) = P(A \mid B)P(B)$$

$$P(A \mid B) = \frac{P(A, B)}{P(B)} \text{ and } P(B \mid A) = \frac{P(A, B)}{P(A)}$$

**Product Rule Part 2!**  $P(A, B) = P(A \mid B)P(B) = P(B \mid A)P(A)$   
(Often left to your imaginations)

# Law of Total Probabilities: Event and partition

The law of total probability is<sup>[1]</sup> the proposition that if  $\{B_n : n = 1, 2, 3, \dots\}$  is a finite or countably infinite partition of a sample space (in other words, a set of pairwise disjoint events whose union is the entire sample space) and each event  $B_n$  is measurable, then for any event  $A$  of the same probability space:

$$\Pr(A) = \sum_n \Pr(A \cap B_n)$$

or, alternatively,<sup>[1]</sup>

$$\Pr(A) = \sum_n \Pr(A | B_n) \Pr(B_n),$$

where, for any  $n$  for which  $\Pr(B_n) = 0$  these terms are simply omitted from the summation, because  $\Pr(A | B_n)$  is finite.

The summation can be interpreted as a weighted average, and consequently the marginal probability,  $\Pr(A)$ , is sometimes called "average probability".<sup>[2]</sup> "overall probability" is sometimes used in less formal writings.<sup>[3]</sup>

The law of total probability can also be stated for conditional probabilities. Taking the  $B_n$  as above, and assuming  $C$  is an event independent with any of the  $B_n$ :

$$\Pr(A | C) = \sum_n \Pr(A | C \cap B_n) \Pr(B_n | C) = \sum_n \Pr(A | C \cap B_n) \Pr(B_n)$$

Boole's inequality

Venn diagram · Tree diagram

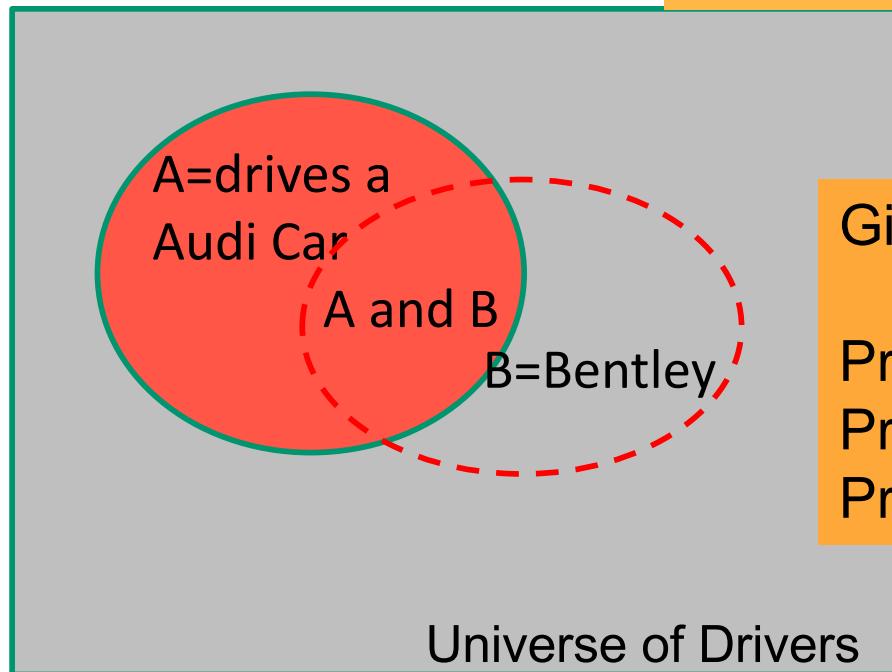
Given an sample space S, event A and partition of S, B  
For any event A we can compute its marginal probability  
with respect to a partition B as the sum of Pr(A and B<sub>n</sub>)

# Total Probabilities Example: $\Pr(A) = \Pr(A) = \Pr(A|B) + \Pr(A|\text{not } B)$

$$P(A|B) = \frac{P(A \cap B)}{P(B)}$$

Conditional probability can be seen to be the probability with respect to a reduced sample space. We can illustrate the conditional probability with the Venn diagram.

$\Pr(A) \rightarrow P(A|B)$



Given:  $A = 5$   $B = 5$ ,  $A \text{ and } B = 2$ , universe = 20

$$\Pr(A|B) = 2/20 / 5/20 = 2/5$$

$$\Pr(A|\text{not } B) = 3/20 / 15/20 = 3/5$$

$$\Pr(A) = \Pr(A|B) + \Pr(A|\text{not } B) = 2/5 + 3/5 = 1$$

# Marginalize via Conditional Probability

Simpler calculation: get marginal from joint; decompose joint using product rule

FROM  $P(A) = \sum_i P(A, B_i)$  Marginalize A from joint A,B  
AND PRODUCT RULE  $P(A, B) = P(A | B)P(B)$



$$P(A) = \sum_i P(A | B_i)P(B_i)$$

the belief in any event  $A$  is a weighted sum over the beliefs in all the distinct ways that  $A$  might be realized.

# Chain Rule is a generalization of the Product Rule

Decompose joint probabilities into a product of simpler conditionals

PRODUCT RULE:  $P(A, B) = P(A | B)P(B)$

1. Chain Rule = Generalization of PRODUCT RULE.
2. It permits the calculation of any member of the joint distribution of a set of random variables using only conditional probabilities

$$P(E_1, E_2 \dots E_n) = P(E_n | E_{n-1}, \dots, E_2, E_1)P(E_{n-1}, \dots, E_2, E_1)$$

let's further decompose!

$$\dots = P(E_n | E_{n-1}, \dots, E_2, E_1) \dots P(E_2 | E_1)P(E_1)$$

- **For example, derive using the chain rule**

$$P(A, B, C, D) = P(A | B, C, D)P(B | C, D)P(C | D)P(D)$$

# Chain Rule Example: Use to calculate joint prob.

---

- The rule is useful in the study of Bayesian networks, which describe a probability distribution in terms of conditional probabilities.
- Assume Urn 1 has 1 black balls and 2 white balls and Urn 2 has 1 black ball and 3 white balls. Suppose we pick an urn at random and then select a ball from that urn.
- Let event A be choosing the first urn:  $P(A) = P(\sim A) = 1/2$ . Let event B be the chance we choose a white ball. Chance of choosing a white ball, given that we've chose the first urn, is  $P(B|A) = 2/3$ . Chance of choosing a white ball, given that we've chosen the second urn is  $P(B|\sim A) = 3/4$ .
- Event A, B would be their intersection; choosing the first urn and a white ball from it. The probability can be found by the chain rule for probability:

$$P(A, B) = P(B | A)P(A) = 2/3 \times 1/2 = 1/3$$

[[Russell, Stuart J.](#); [Norvig, Peter](#) (2003), [\*Artificial Intelligence: A Modern Approach\*](#) (2nd ed.), Upper Saddle River, NJ: Prentice Hall, [ISBN 0-13-790395-2](#), <http://aima.cs.berkeley.edu/>, p. 496.]

# Bayes Rule: 1, 2, 3

$$1. P(A | B) = \frac{P(A, B)}{P(B)}$$

$$P(B | A) = \frac{P(A, B)}{P(A)} \quad \text{Cond<sup>al</sup> Prob}$$

$$2. P(A, B) = P(A | B)P(B) \quad P(A, B) = P(B | A)P(A) \quad \text{Product Rule}$$

$$3. P(A | B)P(B) = P(B | A)P(A)$$

Prior  
Probability

$$P(A | B) = \frac{P(B | A)P(A)}{P(B)}$$

Posterior probability

Posterior	Likelihood	Prior
-----------	------------	-------

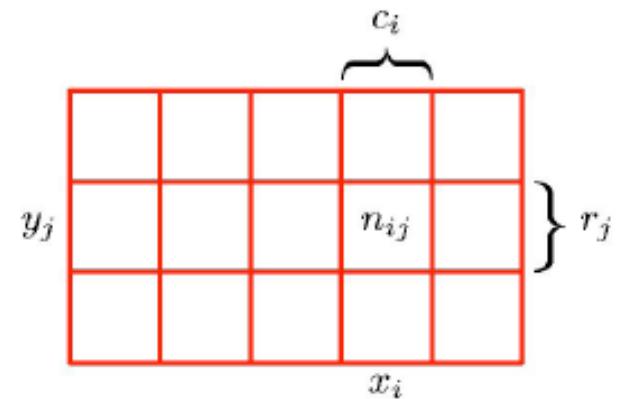
$$P(A | B) = \frac{P(B | A)P(A)}{\sum_{i=1}^n P(B | A = a_i)P(A = a_i)}$$

Marginalize B  
Total probabilities

# Rules of Probability

- Given random variables  $X$  and  $Y$
- Sum Rule gives Marginal Probability

$$p(X = x_i) = \sum_{j=1}^L p(X = x_i, Y = y_j) = \frac{c_i}{N}$$



- Product Rule: joint probability in terms of conditional and marginal

$$p(X, Y) = \frac{n_{ij}}{N} = p(Y | X)p(X) = \frac{n_{ij}}{c_i} \times \frac{c_i}{N}$$

- Combining we get Bayes Rule

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

where

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

Viewed as

Posterior  $\propto$  likelihood  $\times$  prior

[Srihari]



# Probability Theory Summary

---

1. **cond<sup>al</sup> probability**  $P(A | B) = \frac{P(A, B)}{P(B)}$

2. **PRODUCT RULE and Chain Rule:**

$$P(A, B) = P(A | B)P(B)$$

$$P(A, B, C, D) = P(A | B, C, D)P(B | C, D)P(C | D)P(D)$$

3. **Total Probabilities :**

$$1. \Pr(A) = \Pr(A) = \Pr(A|B) + \Pr(A|\text{not } B)$$

4. **Bayes Rules**

Posterior      Likelihood      Prior

$$P(A | B) = \frac{P(B | A)P(A)}{P(B)}$$

$$P(A | B) = \frac{P(B | A)P(A)}{\sum_{i=1}^n P(B | A = a_i)P(A = a_i)}$$

# Maximum a Posteriori (MAP)

---

- Any such maximally probable hypothesis is called a *maximum a posteriori (MAP) hypothesis*

# Bayes Rule Puzzle

---

- Great example to test your understanding of Bayes Rule in practice
- Please test yourself

# Bayes Rule

---

## Puzzle

test with two possible outcomes:  $\oplus$  (positive) and  $\ominus$  (negative). We have prior knowledge that over the entire population of people only .008 have this disease. Furthermore, the lab test is only an imperfect indicator of the disease. The test returns a correct positive result in only 98% of the cases in which the disease is actually present and a correct negative result in only 97% of the cases in which the disease is not present. In other cases, the test returns the opposite result. The above situation can be summarized by the following probabilities:

### PUZZLE

Fill in the probabilities based on the  
above text

$$\Pr(\text{Cancer}) = ??$$

$$\Pr(+|\text{Cancer}) = ??$$

$$\Pr(?????|????)=???$$

# Bayes Rule

---

## Example

test with two possible outcomes:  $\oplus$  (positive) and  $\ominus$  (negative). We have prior knowledge that over the entire population of people only .008 have this disease. Furthermore, the lab test is only an imperfect indicator of the disease. The test returns a correct positive result in only 98% of the cases in which the disease is actually present and a correct negative result in only 97% of the cases in which the disease is not present. In other cases, the test returns the opposite result. The above situation can be summarized by the following probabilities:

$$\begin{aligned} P(\text{cancer}) &= .008, & P(\neg\text{cancer}) &= .992 \\ P(\oplus|\text{cancer}) &= .98, & P(\ominus|\text{cancer}) &= .02 \\ P(\oplus|\neg\text{cancer}) &= .03, & P(\ominus|\neg\text{cancer}) &= .97 \end{aligned}$$

Patient has cancer

### PUZZLE:

**Observe a new patient for whom the lab tests return a positive result**

**Should we diagnose the patient as having cancer or not?**

# Bayes Rule

---

## Example

test with two possible outcomes:  $\oplus$  (positive) and  $\ominus$  (negative). We have prior knowledge that over the entire population of people only .008 have this disease. Furthermore, the lab test is only an imperfect indicator of the disease. The test returns a correct positive result in only 98% of the cases in which the disease is actually present and a correct negative result in only 97% of the cases in which the disease is not present. In other cases, the test returns the opposite result. The above situation can be summarized by the following probabilities:

$$\begin{aligned}P(\text{cancer}) &= .008, & P(\neg\text{cancer}) &= .992 \\P(\oplus|\text{cancer}) &= .98, & P(\ominus|\text{cancer}) &= .02 \\P(\oplus|\neg\text{cancer}) &= .03, & P(\ominus|\neg\text{cancer}) &= .97\end{aligned}$$

Patient has cancer

### PUZZLE:

Observe a new patient for whom the lab tests return a positive result

Should we diagnose the patient as having cancer or not?

Use Bayes theorem to calculate

Given  $\text{Pr}(\text{Cancer}|+) = \text{Pr}(+|\text{Cancer}) \times \text{Pr}(\text{Cancer}) / \text{Pr}(+)$

Then...

# Bayes Rule

---

## Example

test with two possible outcomes:  $\oplus$  (positive) and  $\ominus$  (negative). We have prior knowledge that over the entire population of people only .008 have this disease. Furthermore, the lab test is only an imperfect indicator of the disease. The test returns a correct positive result in only 98% of the cases in which the disease is actually present and a correct negative result in only 97% of the cases in which the disease is not present. In other cases, the test returns the opposite result. The above situation can be summarized by the following probabilities:

$$P(\text{cancer}) = .008, \quad P(\neg\text{cancer}) = .992$$

$$P(\oplus|\text{cancer}) = .98, \quad P(\ominus|\text{cancer}) = .02$$

$$P(\oplus|\neg\text{cancer}) = .03, \quad P(\ominus|\neg\text{cancer}) = .97$$

Patient has cancer

Suppose we now observe a new patient for whom the lab test returns a positive result. Should we diagnose the patient as having cancer or not? The maximum a posteriori hypothesis can be found using Equation (6.3):

$$\text{Numerator for } P(\text{cancer}|\oplus) \quad P(\oplus|\text{cancer})P(\text{cancer}) = (.98).008 = .0078$$

$$P(\text{Cancer}|\text{Test} = +) = \frac{P(\text{Test} = +|\text{Cancer})P(\text{Cancer})}{P(\text{Test} = +)}$$

$$P(\oplus|\neg\text{cancer})P(\neg\text{cancer}) = (.03).992 = .0298$$

Numerator for  $P(\text{NOTcancer}|\oplus)$

Thus,  $n_{MAP} = \neg\text{cancer}$ . The exact posterior probabilities can also be determined by normalizing the above quantities so that they sum to 1 (e.g.,  $P(\text{cancer}|\oplus) = \frac{.0078}{.0078+.0298} = .21$ ). This step is warranted because Bayes theorem states that the posterior probabilities are just the above quantities divided by the probability of the data,  $P(\oplus)$ . Although  $P(\oplus)$  was not provided directly as part of the problem statement, we can calculate it in this fashion because we know that  $P(\text{cancer}|\oplus)$  and  $P(\neg\text{cancer}|\oplus)$  must sum to 1 (i.e., either the patient has cancer or they do not). Notice that while the posterior probability of *cancer* is significantly higher than its prior probability, the most probable hypothesis is still that the patient does

# Probability Basics → Naïve Bayes Models

- **Introduction**
- **Probability Basics**
  - Probability Axioms
  - Conditional probabilities
  - Product Rule, Chain Rule, Bayes Rule
- **Bayes Nets And Naïve Bayes**
  - Learning
  - Independence
  - Conditional independence
  - Naïve Bayes derivation (discrete case plus smoothing)
- **Naïve Bayes Flavors**
  - Continuous input variables
  - Discrete input variables (2 flavors: Bernoulli, multinomial)
- **Case Study: Spam detector in Naïve Bayes**
- **Summary**

# Digression

---

- **Digress to cover Confidence intervals**

# Confidence interval for a binomial distribution

---

- In statistics, a binomial proportion confidence interval is a confidence interval for a proportion in a statistical population. It uses the proportion estimated in a statistical sample and allows for sampling error. There are several formulas for a binomial confidence interval, but all of them rely on the assumption of a binomial distribution.
- In general, a binomial distribution applies when an experiment is repeated a fixed number of times, each trial of the experiment has two possible outcomes (labeled arbitrarily success and failure), the probability of success is the same for each trial, and the trials are statistically independent.
- A simple example of a binomial distribution is the set of various possible outcomes, and their probabilities, for the number of heads observed when a (not necessarily fair) coin is flipped ten times.
- The observed binomial proportion is the fraction of the flips which turn out to be heads.
- Given this observed proportion, the confidence interval for the true proportion innate in that coin is a range of possible proportions which may contain the true proportion.
- A 95% confidence interval for the proportion, for instance, will contain the true proportion 95% of the times that the procedure for constructing the confidence interval is employed.
- Note that this does not mean that a calculated 95% confidence interval will contain the true proportion with 95% probability. Instead, one should interpret it as follows: the process of drawing a random sample and calculating an accompanying 95% confidence interval will generate a confidence interval that contains the true proportion in 95% of all cases.
- There are several ways to compute a confidence interval for a binomial proportion. The normal approximation interval is the simplest formula, and the one introduced in most basic Statistics classes and textbooks. This formula, however, is based on an approximation that does not always work well.

$$\hat{p} \pm z_{1-\frac{\alpha}{2}} \sqrt{\frac{1}{n} \hat{p} (1 - \hat{p})}$$

[https://en.wikipedia.org/wiki/Binomial\\_proportion\\_confidence\\_interval](https://en.wikipedia.org/wiki/Binomial_proportion_confidence_interval)

# Confidence interval for a binomial distribution

---

- In statistics, a binomial proportion confidence interval is a confidence interval for a proportion in a statistical population.
  - The number of people who voted for Donald Trump in a poll
- It uses the proportion estimated in a statistical sample and allows for sampling error.
- There are several formulas for a binomial confidence interval, but all of them rely on the assumption of a binomial distribution.

Confidence interval for a binomial distribution

$$\hat{p} \pm z_{1-\frac{\alpha}{2}} \sqrt{\frac{1}{n} \hat{p} (1 - \hat{p})}$$

# Confidence interval for a binomial distribution

[https://en.wikipedia.org/wiki/Binomial\\_proportion\\_confidence\\_interval](https://en.wikipedia.org/wiki/Binomial_proportion_confidence_interval)

Confidence interval for a binomial distribution

$$\hat{p} \pm z_{1-\frac{\alpha}{2}} \sqrt{\frac{1}{n} \hat{p}(1 - \hat{p})}$$

$$\hat{p} \pm \frac{z^*}{2\sqrt{n}}$$

Example: Poll of 1,000 people  
 $2 * 1 / (2 * \text{sqrt}(n))$ , e.g.,  $1 / \text{sqrt}(1000) = \pm 3\%$

Election Poll of 1000 people ( $N=1000$ )

Assume 45% ( $P=0.45$ ,  $Q=1-P=55\%$ ) favor candidate A (Donald) versus 49% for Candidate B (Hillary)

Standard error about the mean =  $\text{SQRT}(PQ/N) = \text{SQRT}(0.45 * 0.55 / 1000) = 0.015 = 1.5\%$

So the 95% confidence interval surrounding Donald's support of 45% is  $45\% \pm 2 * 1.5 = [42, \dots 48]$

In machine learning:

$$\Pr(Y=\text{Business}) = 100/1000 = 0.1$$

$$\text{SQRT}(0.1 * 0.9 / 1000) = 0.01 = 1\%$$

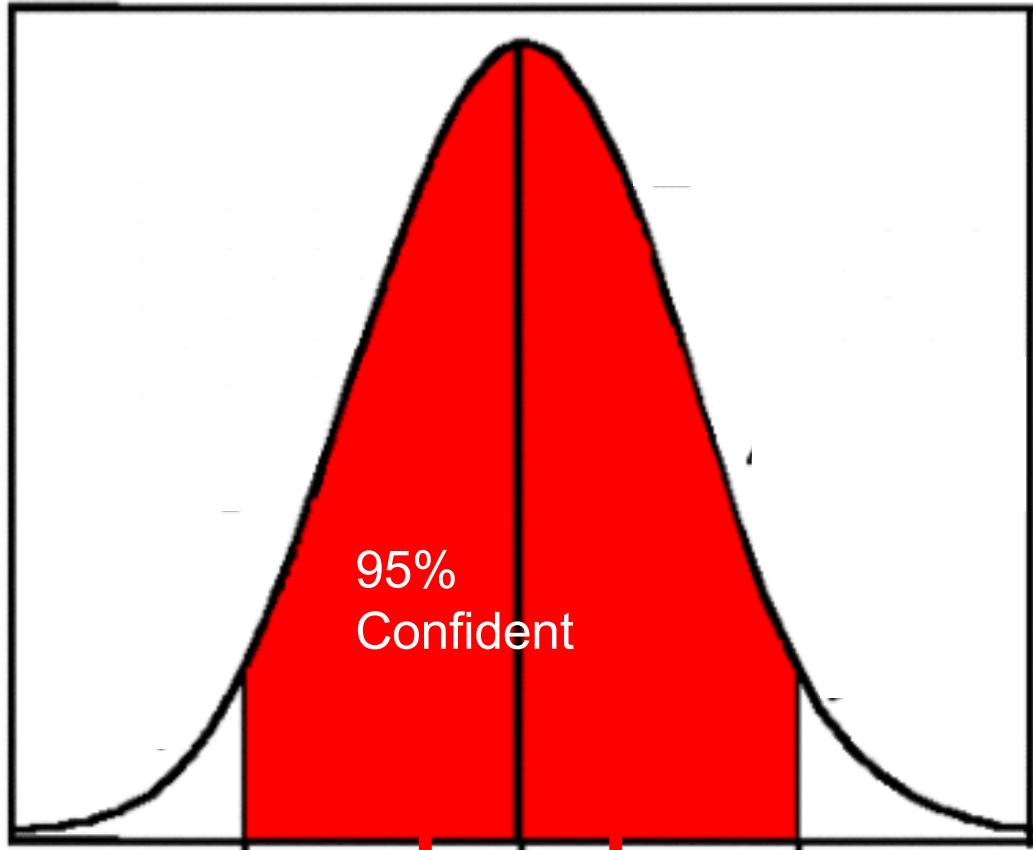
yielding  $\Pr(Y=\text{Business})$  Confidence interval =  $0.1 \pm 0.01 * 2$

[0.08, 0.12] Note we need 1,000 examples to get this tight CI;

How much data do we need to estimate  $\Pr(Y|X)$  with confidence?

# Estimating Reliable Probabilities

Estimate using Binomial  
MLE Estimates  
I.e., #Clicks/#Impression



1.6%    2.6%  
2.3%    2.9%

\$40/1,000 @CPC of \$1.60  
\$400/10,000

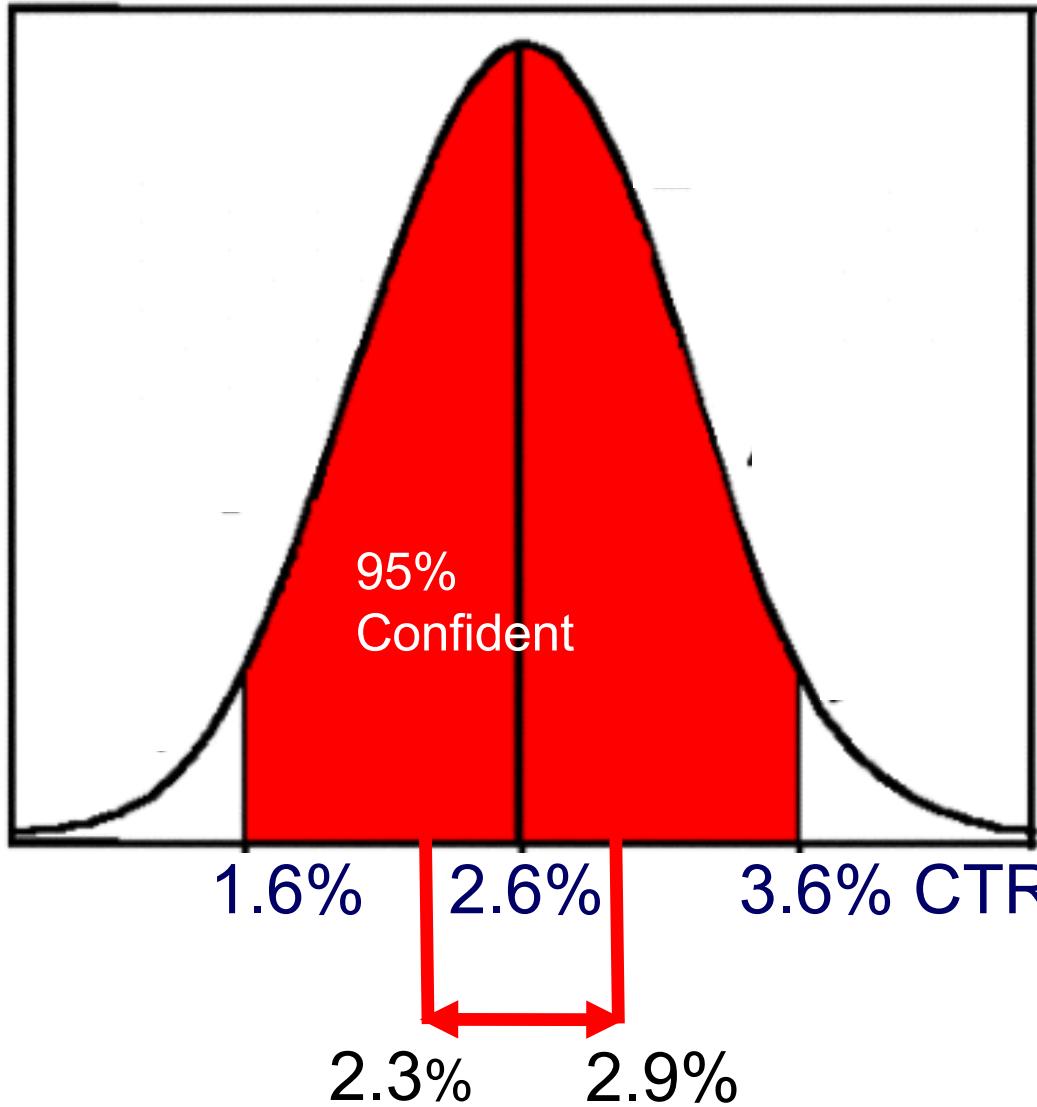
Standard error of the mean of one sample is the estimate of the standard deviation that would be obtained from the means of a large number of samples drawn from that population.

$$\text{StdErr} = \sqrt{.026 * (1 - 0.026) / 1000} * 1.96 = \pm 1\%$$

CTR (after 1,000 impressions)

(after 10,000 impressions)

# Estimating CTR (and later AR)



For a network of  
~ $10^9$  target pages,  
~ $10^6$  ads  
~ $10^7$  users

- .....
- Cannot afford this evaluation/auditioning
  - Borrow strength, marginalize
  - CoD (curse of dimensionality)

# Confidence Intervals in a Nutshell

---

- **The assumptions required for CI for a population proportion to be valid:**
  - the sample size  $n$  is large enough (check:  $\hat{p} \geq 10$  and  $n(1 - \hat{p}) \geq 10$ )
  - the data are a *random sample* from that population
- **General Confidence Interval for the Population Proportion  $p$ :**

$$\hat{p} \pm 2\sqrt{\frac{\hat{p}(1 - \hat{p})}{n}}$$

- **Approximate 95% Confidence Interval for the Population Proportion  $p$ :**

$$\hat{p} \pm z^* \sqrt{\frac{\hat{p}(1 - \hat{p})}{n}} > 2 * \text{sqrt}(.4 * .6 / 1000) \text{ #president's satisfaction rating}$$

[1] 0.03098387

- Note: standard error of  $\hat{p}$  which is largest when  $\hat{p} = \frac{1}{2}$
- **Conservative Confidence Interval for the Population Proportion  $p$ :**

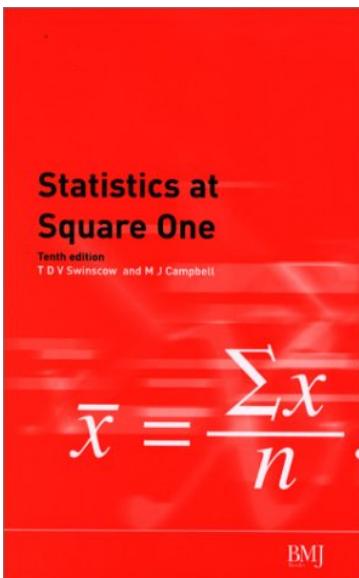
$$\hat{p} \pm \frac{z^*}{2\sqrt{n}} \quad 2 * 1 / (2 * \text{sqrt}(n)), \text{ e.g., } 1 / \text{sqrt}(1000) = \pm 3\%$$

# Sample Size Needed

- Sample Size Needed for Desired Confidence Level and Error Margin where  $m$  is the desired margin of error.

$$n = \left( \frac{z^*}{2m} \right)^2$$

SE= =1/sqrt(n)



Measure CTRs confidently  $p \pm 0.001$   
 $> (1.96/(2*.001))^2$   
[1] 960,400 sample size

# Bayes Theorem for Machine Learning

---

$$P(h | D) = \frac{P(D | h)P(h)}{P(D)}$$

- **P(h) = Prior probability of hypothesis**
- **P(D) = Prior probability of training data D.**
- **P(h|D) = Probability of h given D.**
- **P(D|h) = Probability of D given h.**

- 
- **Bayesian Learning via**
    - Max Likelihood inference
    - Hierarchical Bayesian inference leading to the MAP hypothesis

# Bayesian Learning: Discrete input/output variables

---

- Here we consider the relationship between supervised learning, or function approximation problems, and Bayesian reasoning.
- We begin by considering how to design learning algorithms based on Bayes rule.
  - Consider a supervised learning problem in which we wish to approximate an unknown target function
    - $f : X \rightarrow Y$ , or equivalently  $P(Y|X)$ .
  - To begin, we will assume  $Y$  is a boolean-valued random variable, and  $X$  is a vector containing  $n$  boolean attributes.
  - In other words,  $X = \langle X_1, X_2, \dots, X_n \rangle$ , where  $X_i$  is the boolean random variable denoting the  $i$ th attribute of  $X$ .

# Probability Basics → Naïve Bayes Models

- **Introduction**
- **Probability Basics**
  - Probability Axioms
  - Conditional probabilities
  - Product Rule, Chain Rule, Bayes Rule
- **Bayes Nets And Naïve Bayes**
  - Learning
  - Independence
  - Conditional independence
  - Naïve Bayes derivation (discrete case plus smoothing)
- **Naïve Bayes Flavors**
  - Continuous input variables
  - Discrete input variables (2 flavors: Bernoulli, multinomial)
- **Case Study: Spam detector in Naïve Bayes**
- **Summary**

# Supervised Learning Via Bayes Rule

---

- Consider a supervised learning problem in which we wish to approximate an unknown target function  $f : X \rightarrow Y$ , or equivalently  $P(Y_j | X)$ .
- For pedagogical reasons assume discrete binary variables for both and input ( $X$ ) and output ( $Y$ )
  - To begin, we will assume  $Y$  is a boolean-valued random variable, and  $X$  is a vector containing  $n$  boolean attributes. In other words,  $X = hX_1; X_2 : \dots ; X_n$ , where  $X_i$  is the boolean random variable denoting the  $i$ th attribute of  $X$ .
- Applying Bayes rule, we see that  $P(Y = y_i | X)$  can be represented as

$$P(Y = y_i | X = x_k) = \frac{P(X = x_k | Y = y_i)P(Y = y_i)}{\sum_j P(X = x_k | Y = y_j)P(Y = y_j)}$$

- where  $y_m$  represents the  $m^{th}$  possible value for  $Y$ , and where the summation in the denominator is over all legal values of the random variable  $Y$

# Goal is to learn Class conditionals $\Pr(X|Y)$

---

Applying Bayes rule, we see that  $P(Y = y_i|X)$  can be represented as

$$P(Y = y_i|X = x_k) = \frac{P(X = x_k|Y = y_i)P(Y = y_i)}{\sum_j P(X = x_k|Y = y_j)P(Y = y_j)}$$

- Then use to classify X

# Learning Classifiers based on Bayes Rule

Posterior

Likelihood

Prior

$$P(A | B) = \frac{P(B | A)P(A)}{P(B)}$$

Assume

Y is a boolean-valued random variable,  
X is a vector containing n boolean attributes.  
In other words,  $X = \langle X_1, X_2, \dots, X_n \rangle$ , where  $X_i$  is the  
boolean random variable denoting the  $i$ th attribute of  $X$ .

Applying Bayes rule, we see that  $P(Y = y_i | X)$  can be represented as

$$P(Y = y_i | X = x_k) = \frac{P(X = x_k | Y = y_i)P(Y = y_i)}{\sum_j P(X = x_k | Y = y_j)P(Y = y_j)}$$

LHS

RHS

where  $y_m$  denotes the  $m$ th possible value for  $Y$ ,  $x_k$  denotes the  $k$ th possible vector value for  $X$ , and where the summation in the denominator is over all legal values of the random variable  $Y$ .

One way to learn  $P(Y|X)$  is to use the training data to estimate  $P(X|Y)$  and  $P(Y)$ . We can then use these estimates, together with Bayes rule above, to determine  $P(Y|X = x_k)$  for any new instance  $x_k$ .

Confidence interval for a binomial distribution

$$\hat{p} \pm z_{1-\frac{\alpha}{2}} \sqrt{\frac{1}{n} \hat{p}(1 - \hat{p})}$$

# Estimating $P(Y)$ , $P(X|Y)$ from data

- If we are going to train a Bayes classifier by estimating  $P(X|Y)$  and  $P(Y)$ , then it is reasonable to ask how much training data will be required to obtain reliable estimates of these distributions.
- Let us assume training examples are generated by drawing instances at random from an unknown underlying distribution  $P(X)$ , then allowing a teacher to label this example with its  $Y$  value.
- A 1,000 independently drawn training examples will usually suffice to obtain a maximum likelihood estimate of  $P(Y)$  that is within a few percent of its correct value when  $Y$  is a boolean variable.

Note for a simple proposition the Standard error:

$SE(P(Y)) = \text{SQRT}(0.1 * 0.9 / 1000) = 0.01$  for business labeled docs; assume 100 out 1000

- However, accurately estimating  $P(X|Y)$  typically requires many more examples

Assume  $n$  binary variables:  $2^n$  propositions X 1000 training examples

See example in a moment

Impractical: so what to do?

# Learning a (FULL) Bayesian Model

---

- **One way to learn  $P(Y | X)$** 
  1. Estimate the joint probability distribution  $P(X | Y)$  and  $P(Y)$  from the training data.
  2. Use these estimates, together with Bayes rule above, to determine  $P(Y | X = x_k)$  for any new instance  $x_k$ .
- **A maximum likelihood estimate of  $P(Y)$  can be accomplished with just a few hundred examples**
  - $\#ExamplesWithLabel/\#TotalNumberOfExamples$ . E.g., 45 examples are in class 1 and 55 are in class2 then  $Pr(Y=Class1) = 45/100=0.45$ .
- **However, estimating the joint probability distribution  $P(X | Y)$  requires an exponential amount training examples (even with this assumption of binary input and output variable)**
  - Assume  $n$  input attributes  $X_i$  take 2 discrete values and  $Y$  has 2 possible class values;  $2^n * 2$  possible states of the world (parameters)

# Estimating the Joint Prob. Directly?

---

- **$2^n * 2$  possible states of the world (parameter estimates)**
  - Assume  $n$  input attributes  $X_i$  take 2 discrete values and  $Y$  has 2 possible class values;  $2^n * 2$  possible states of the world (parameters) that we need to estimate from data
  - $\theta_{ij} = P(X = x_i | Y = y_j)$ ;  $2^n$  possible states of the input world 2 possible output states
- **Can reduce  $2^n - 1 * 2$  parameters by exploiting the sum-to-1**
  - Class conditional multinomial needs to sum to 1 so we exploit this and can infer one of the class conditional probs from the rest
  - Each state is a complex combination of feature values
- **Require 200k examples (or more) for reliable estimates of 10 binary variable problem**
  - So for 10 input variables and one output variable we have to estimate 2046 states. To estimate probabilities requiring at least 204,600 examples for reliable estimates.
- **So not very realistic, even in these WWW times**

# Learning a Bayesian Model

- $2^n * 2$  possible states of the world
  - Assume  $n$  input attributes  $X_i$  take 2 discrete values and  $Y$  has 2 possible class values;  $2^n * 2$  possible states of the world (parameters) that we need to estimate from data
    - $\theta_{ij} = P(X = x_i | Y = y_j)$ ;  $2^n$  possible states of the input world  $2$  possible output states

Index $X=X^i$	$X_1$ $X=x_1^i)$	$X_2$ $X=x_2^i)$	$Y$ $Y=y_j$	$\theta_{ij}=P(X=x^i y=y_j)$
1	1	1	1	$\theta_{11}$
2	0	1	1	$\theta_{21}$
3	1	0	1	$\theta_{31}$
4	0	0	1	$\theta_{41}$
5	1	1	0	$\theta_{50}$
6	0	1	0	$\theta_{60}$
7	1	0	0	$\theta_{70}$
8	0	0	0	$\theta_{80}$

Sum to 1

# Cant Estimate Joint Probability

---

- **Look for ways to combat the intractable data needs for learning a Bayesian Classifier**
  - Leverage the chain rule and other assumptions (Markov, Naïve Bayes); this leads to Bayesian Networks
  - Make a conditional independence assumption; this leads to a Naïve Bayes classifier
    - Reduces the number of parameters from  $2^n - 1 * 2$  parameters to  $2n$

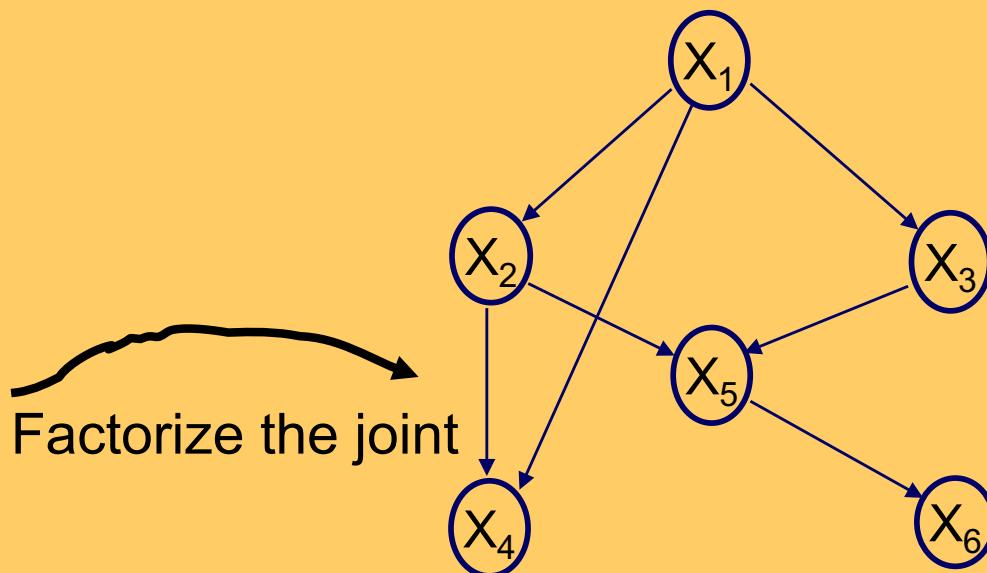
# Bayesian Networks

GOAL:  $2^N$  Possibilities  $\rightarrow 2^N$

P: Joint Probability Distribution

#	X1	X2	X3	X4	X5	X6	Pr(X1,X2, X3, X4, X5, X6)
1	1						
2	1						
..	1						
4..	1						
5	1						
6	1						
7	1						
64							

G: Directed Acyclic Graph



$$p(x_1, x_2, x_3, x_4, x_5, x_6)$$

1. Partial Order

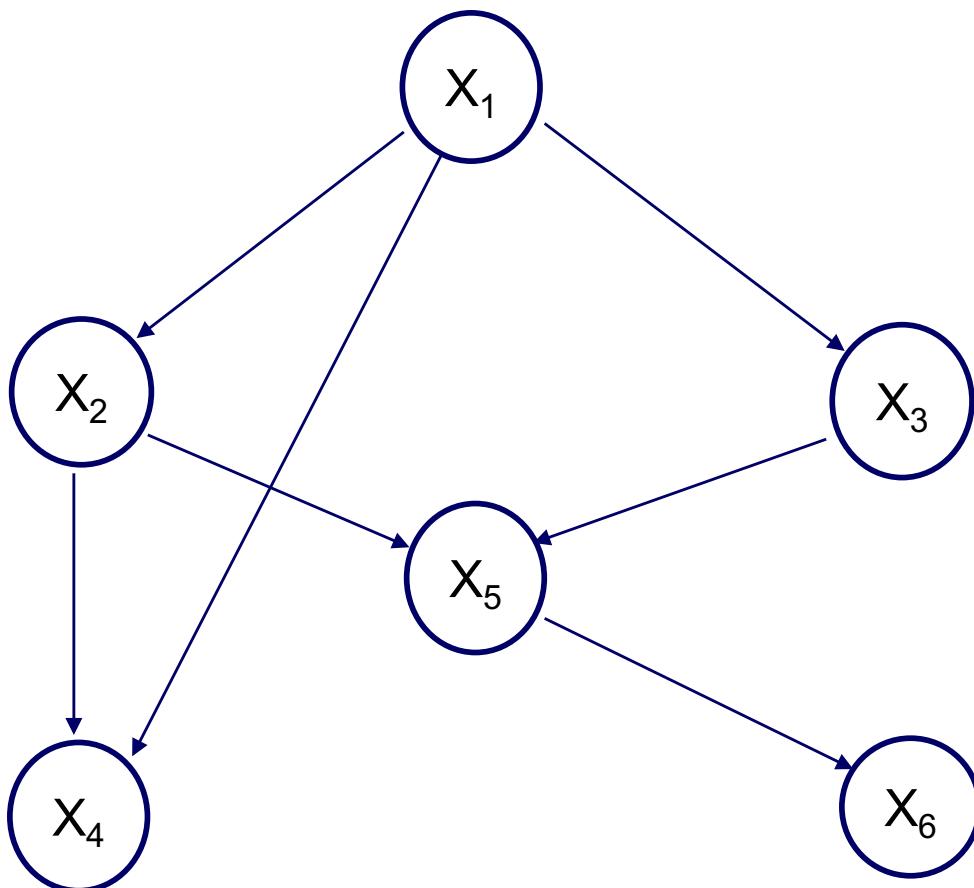
$$= p(x_1) p(x_2 | x_1) p(x_3 | x_1, x_2) p(x_4 | x_1, x_2, x_3) p(x_5 | x_4, x_3, x_2, x_1) p(x_6 | x_5, x_4, x_3, x_2, x_1)$$

$$= p(x_1) p(x_2 | x_1) p(x_3 | x_1) p(x_4 | x_2, x_1) p(x_5 | x_3, x_2) p(x_6 | x_5) \quad 3. \text{ Markov Property}$$

$$= p(x_1) p(x_2) p(x_3) p(x_4) p(x_5) p(x_6) \quad 4: \text{Independence (see next section)}$$

$$P(y|x_6, x_5, x_4, x_3, x_2, x_1) = p(x_1|y) p(x_2|y) p(x_3|y) p(x_4|y) p(x_5|y) p(x_6|y) \quad 5: \text{Naïve Bayes via Cond. Independence}$$

# 1, 2, 3 Example of Factorization



$$p(x_1, x_2, x_3, x_4, x_5, x_6)$$

1. Partial Order

$$= p(x_1) p(x_2 | x_1) p(x_3 | x_1, x_2) p(x_4 | x_1, x_2, x_3) p(x_5 | x_4, x_3, x_2, x_1) p(x_6 | x_5, x_4, x_3, x_2, x_1)$$

$$= p(x_1) p(x_2 | x_1) p(x_3 | x_1) p(x_4 | x_2, x_1) p(x_5 | x_3, x_2) p(x_6 | x_5)$$

2. By Chain Rule

3. Markov Property

# Factorization

---

- Given a DAG G, topologically sort the variables:  
 $X_1, \dots, X_n$  (s.t. if  $i < j$ , then  $X_j$  is not an ancestor of  $X_i$ )
- For any joint distribution P that is Markov to G, factorize it as follows:  
$$\begin{aligned} P(X_1, \dots, X_n) \\ = P(X_1) P(X_2|X_1) \dots P(X_n|X_1, \dots, X_{n-1}) &\quad \text{By Chain Rule} \\ = \prod_i P(X_i | \text{Pa}(X_i)) &\quad \text{Exploit Local Markov Property} \end{aligned}$$
- Markov property: every **variable** is independent of its **non-descendants** given its **parents**.
- Markov Condition requires that every conditional independence in the graph is in the joint probability distribution.

# Probability Basics → Naïve Bayes Models

- **Introduction**
- **Probability Basics**
  - Probability Axioms
  - Conditional probabilities
  - Product Rule, Chain Rule, Bayes Rule
- **Bayes Nets And Naïve Bayes**
  - Learning
  - Independence
  - Conditional independence
  - Naïve Bayes derivation (discrete case plus smoothing)
- **Naïve Bayes Flavors**
  - Continuous input variables
  - Discrete input variables (2 flavors: Bernoulli, multinomial)
- **Case Study: Spam detector in Naïve Bayes**
- **Summary**

# Independence

Event

- In probability theory, two events are independent, statistically independent, or stochastically independent if the occurrence of one does not affect the probability of the other.

- $P(A|B) = Pr(A)$
  - As a result  $P(A \text{ and } B) = Pr(A) \times Pr(B);$

Variables

- Similarly, two random variables are independent if the realization of one does not affect the probability distribution of the other.

- The concept of independence extends to dealing with collections of more than two events or random variables, in which case the events are pairwise independent if each pair are independent of each other, and the events are mutually independent if each event is independent of each other combination of events.

[https://en.wikipedia.org/wiki/Independence\\_\(probability\\_theory\)](https://en.wikipedia.org/wiki/Independence_(probability_theory))

# $P(A \text{ and } B) = Pr(A) \times Pr(B)$

---

## Two events [\[edit\]](#)

Two events  $A$  and  $B$  are **independent** (often written as  $A \perp B$  or  $A \perp\!\!\!\perp B$ ) if and only if their **joint probability** equals the product of their probabilities:

$$P(A \cap B) = P(A)P(B).$$

Why this defines independence is made clear by rewriting with **conditional probabilities**:

$$P(A \cap B) = P(A)P(B) \Leftrightarrow P(A) = \frac{P(A)P(B)}{P(B)} = \frac{P(A \cap B)}{P(B)} = P(A | B)$$

and similarly

$$P(A \cap B) = P(A)P(B) \Leftrightarrow P(B) = P(B | A).$$

Thus, the occurrence of  $B$  does not affect the probability of  $A$ , and vice versa. Although the derived expressions may seem more intuitive, they are not the preferred definition, as the conditional probabilities may be undefined if  $P(A)$  or  $P(B)$  are 0. Furthermore, the preferred definition makes clear by symmetry that when  $A$  is independent of  $B$ ,  $B$  is also independent of  $A$ .

## More than two events [\[edit\]](#)

A finite set of events  $\{A_i\}$  is **pairwise independent** if and only if every pair of events is independent<sup>[2]</sup>—that is, if and only if for all distinct pairs of indices  $m, k$ ,

$$P(A_m \cap A_k) = P(A_m)P(A_k).$$

A finite set of events is **mutually independent** if and only if every event is independent of any intersection of the other events<sup>[2]</sup>—that is, if and only if for every  $n$ -element subset  $\{A_i\}$ ,

$$P\left(\bigcap_{i=1}^n A_i\right) = \prod_{i=1}^n P(A_i).$$

This is called the *multiplication rule* for independent events. Note that it is not a single condition involving only the product of all the probabilities of all single events (see [below](#) for a counterexample); it must hold true for all subset of events.

For more than two events, a mutually independent set of events is (by definition) pairwise independent; but the converse is not necessarily true (see [below](#) for a counterexample).

# Two events are independent if...

---

**Definition 1.3** Two events  $E$  and  $F$  are independent if one of the following hold:

1.  $P(E|F) = P(E)$  and  $P(E) \neq 0, P(F) \neq 0$ .
2.  $P(E) = 0$  or  $P(F) = 0$ .

Notice that the definition states that the two events are independent even though it is based on the conditional probability of  $E$  given  $F$ . The reason is that independence is symmetric. That is, if  $P(E) \neq 0$  and  $P(F) \neq 0$ , then  $P(E|F) = P(E)$  if and only if  $P(F|E) = P(F)$ . It is straightforward to prove that  $E$  and  $F$  are independent if and only if  $P(E \cap F) = P(E)P(F)$ .

**Independence is symmetric**

# Independence Example

## $\Pr(\text{Queen}) = \Pr(\text{Queen}|\text{Spade})$

---

**Example 1.6** Again consider drawing the top card from a deck of cards, let Queen be the set of the 4 queens, RoyalCard be the set of the 12 royal cards, and Spade be the set of the 13 spades. Then

$$P(\text{Queen}) = \frac{1}{13}$$

$$P(\text{Queen}|\text{RoyalCard}) = \frac{P(\text{Queen} \cap \text{RoyalCard})}{P(\text{RoyalCard})} = \frac{1/13}{3/13} = \frac{1}{3}$$

$$P(\text{Queen}|\text{Spade}) = \frac{P(\text{Queen} \cap \text{Spade})}{P(\text{Spade})} = \frac{1/52}{1/4} = \frac{1}{13}.$$

Notice in the previous example that  $P(\text{Queen}|\text{Spade}) = P(\text{Queen})$ . This means that finding out the card is a spade does not make it more or less probable that it is a queen. That is, the knowledge of whether it is a spade is irrelevant to whether it is a queen. We say that the two events are independent in this case, which is formalized in the following definition.

[Learning Bayesian Networks, Prentice-Hall, Richard E. Neapolitan]

# Independence more generally on sets: Set A and B are independent

---

**Definition 1.6** Suppose we have a probability space  $(\Omega, P)$ , and two sets A and B containing random variables defined on  $\Omega$ . Then the sets A and B are said to be independent if, for all values of the variables in the sets a and b, the events  $A = a$  and  $B = b$  are independent. That is, either  $P(a) = 0$  or  $P(b) = 0$  or

$$P(a|b) = P(a).$$

When this is the case, we write

$$I_P(A, B),$$

where  $I_P$  stands for independent in  $P$ .

# Independence Example for Random Variables R, T, S

**Example 1.18** Let  $\Omega$  be the set of all cards in an ordinary deck, and let  $P$  assign  $1/52$  to each card. Define random variables as follows:

Variable	Value	Outcomes Mapped to this Value
$R$	$r_1$	All royal cards
	$r_2$	All nonroyal cards
$T$	$t_1$	All tens and jacks
	$t_2$	All cards that are neither tens nor jacks
$S$	$s_1$	All spades
	$s_2$	All nonspades

Then we maintain the sets  $\{R, T\}$  and  $\{S\}$  are independent. That is,

$$I_P(\{R, T\}, \{S\}).$$

To show this, we need show for all values of  $r$ ,  $t$ , and  $s$  that

$$P(r, t | s) = P(r, t).$$

(Note that it we do not show brackets to denote sets in our probabilistic expression because in such an expression a set represents the members of the set. See the discussion following Example 1.14.) The following table shows this is the case:

# Independence Example for Random Variables R, T, S

---

$s$	$r$	$t$	$P(r, t s)$	$P(r, t)$
$s1$	$r1$	$t1$	$1/13$	$4/52 = 1/13$
$s1$	$r1$	$t2$	$2/13$	$8/52 = 2/13$
$s1$	$r2$	$t1$	$1/13$	$4/52 = 1/13$
$s1$	$r2$	$t2$	$9/13$	$36/52 = 9/13$
$s2$	$r1$	$t1$	$3/39 = 1/13$	$4/52 = 1/13$
$s2$	$r1$	$t2$	$6/39 = 2/13$	$8/52 = 2/13$
$s2$	$r2$	$t1$	$3/39 = 1/13$	$4/52 = 1/13$
$s2$	$r2$	$t2$	$27/39 = 9/13$	$36/52 = 9/13$

**Definition 1.7** Suppose we have a probability space  $(\Omega, P)$ , and three sets A, B, and C containing random variable defined on  $\Omega$ . Then the sets A and B are said to be conditionally independent given the set C if, for all values of the variables in the sets a, b, and c, whenever  $P(c) \neq 0$ , the events  $A = a$  and  $B = b$  are conditionally independent given the event  $C = c$ . That is, either  $P(a|c) = 0$  or  $P(b|c) = 0$  or

$$P(a|b, c) = P(a|c).$$

When this is the case, we write

$$I_P(A, B|C).$$

# Probability Basics → Naïve Bayes Models

- **Introduction**
- **Probability Basics**
  - Probability Axioms
  - Conditional probabilities
  - Product Rule, Chain Rule, Bayes Rule
- **Bayes Nets And Naïve Bayes**
  - Learning
  - Independence
  - Conditional independence
  - Naïve Bayes derivation (discrete case plus smoothing)
- **Naïve Bayes Flavors**
  - Continuous input variables
  - Discrete input variables (2 flavors: Bernoulli, multinomial)
- **Case Study: Spam detector in Naïve Bayes**
- **Summary**

# Conditional Independence

*Definition:* Given random variables  $X, Y$  and  $Z$ , we say  $X$  is **conditionally independent** of  $Y$  given  $Z$ , if and only if the probability distribution governing  $X$  is independent of the value of  $Y$  given  $Z$ ; that is

Knowing  $Y$  given  $Z$  tells us nothing about  $X$

$$(\forall i, j, k) P(X = x_i | Y = y_j, Z = z_k) = P(X = x_i | Z = z_k)$$

As an example, consider three boolean random variables to describe the current weather: *Rain*, *Thunder* and *Lightning*. We might reasonably assert that *Thunder* is independent of *Rain* given *Lightning*. Because we know *Lightning* causes *Thunder*, once we know whether or not there is *Lightning*, no additional information about *Thunder* is provided by the value of *Rain*. Of course there is a clear dependence of *Thunder* on *Rain* in general, but there is no *conditional* dependence once we know the value of *Lightning*.

$\Pr(\text{Rain} | \text{Lightning}, \text{Thunder}) == \Pr(\text{Rain} | \text{Lightning})$

$$\mathbf{P(R|L, T)= P(R|L)}$$

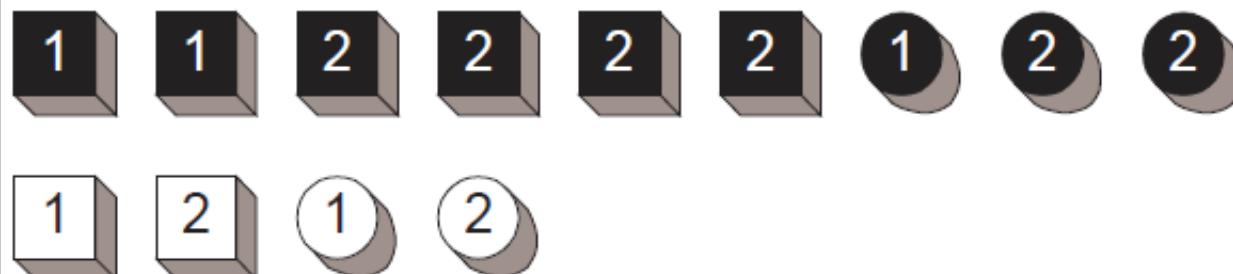
# Conditional independence notation

---

- **V is conditionally independent of set  $Vi$  given set  $Vj$** 
  - if  $p(V | Vi, Vj) = p(V | Vj)$
  - notation:  $I(V, Vi | Vj)$  or  $V \perp Vi | Vj$
- **Intuition**
  - if  $I(V, Vi | Vj)$  then knowing  $Vi$  &  $Vj$  tells nothing more about  $V$  than knowing  $Vj$  alone
  - if we know  $Vj$  we can ignore  $Vi$

$$P(Queen | Club) = \frac{P(Queen, Club)}{P(Club)} = \frac{\frac{1}{52}}{\frac{13}{52}} = \frac{1}{13}$$

$$P(Queen) = \frac{4}{52} = \frac{1}{13}$$



# Conditionally Independent

Figure 1.2: Containing a '1' and being a square are not independent, but they are conditionally independent given the object is black and given it is white.

**Example 1.19** Let  $\Omega$  be the set of all objects in Figure 1.2, and let  $P$  assign  $1/13$  to each object. Define random variables  $S$  (for shape),  $V$  (for value), and  $C$  (for color) as follows:

Variable	Value	Outcomes Mapped to this Value
$V$	$v1$	All objects containing a '1'
	$v2$	All objects containing a '2'
$S$	$s1$	All square objects
	$s2$	All round objects
$C$	$c1$	All black objects
	$c2$	All white objects

Then we maintain that  $\{V\}$  and  $\{S\}$  are conditionally independent given  $\{C\}$ . That is,

$$I_P(\{V\}, \{S\} | \{C\}).$$

To show this, we need show for all values of  $v$ ,  $s$ , and  $c$  that

$$P(v|s, c) = P(v|c).$$

The results in Example 1.8 show  $P(v1|s1, c1) = P(v1|c1)$  and  $P(v1|s1, c2) = P(v1|c2)$ . The table that follows shows the equality holds for the other values of the variables too:

$c$	$s$	$v$	$P(v s, c)$	$P(v c)$
$c1$	$s1$	$v1$	$2/6 = 1/3$	$3/9 = 1/3$
	$s1$	$v2$	$4/6 = 2/3$	$6/9 = 2/3$
$c1$	$s2$	$v1$	$1/3$	$3/9 = 1/3$
	$s2$	$v2$	$2/3$	$6/9 = 2/3$
$c2$	$s1$	$v1$	$1/2$	$2/4 = 1/2$
	$s1$	$v2$	$1/2$	$2/4 = 1/2$
$c2$	$s2$	$v1$	$1/2$	$2/4 = 1/2$
	$s2$	$v2$	$1/2$	$2/4 = 1/2$

[Learning Bayesian Networks,  
Prentice-Hall, Richard E.  
Neapolitan]

# Conditional Independence (Wiki)

---

## Conditional independence [ edit ]

Main article: [Conditional independence](#)

Intuitively, two random variables  $X$  and  $Y$  are conditionally independent given  $Z$  if, once  $Z$  is known, the value of  $Y$  does not add any additional information about  $X$ . For instance, two measurements  $X$  and  $Y$  of the same underlying quantity  $Z$  are not independent, but they are **conditionally independent given  $Z$**  (unless the errors in the two measurements are somehow connected).

The formal definition of conditional independence is based on the idea of conditional distributions. If  $X$ ,  $Y$ , and  $Z$  are discrete random variables, then we define  $X$  and  $Y$  to be *conditionally independent given  $Z$*  if

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

for all  $x$ ,  $y$  and  $z$  such that  $P(Z = z) > 0$ . On the other hand, if the random variables are **continuous** and have a joint probability density function  $p$ , then  $X$  and  $Y$  are **conditionally independent** given  $Z$  if

$$p_{XY|Z}(x, y | z) = p_{X|Z}(x | z) \cdot p_{Y|Z}(y | z)$$

for all real numbers  $x$ ,  $y$  and  $z$  such that  $p_Z(z) > 0$ .

If  $X$  and  $Y$  are conditionally independent given  $Z$ , then

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

for any  $x$ ,  $y$  and  $z$  with  $P(Z = z) > 0$ . That is, the conditional distribution for  $X$  given  $Y$  and  $Z$  is the same as that given  $Z$  alone. A similar equation holds for the conditional probability density functions in the continuous case.

Independence can be seen as a special kind of conditional independence, since probability can be seen as a kind of conditional probability given no events.

[https://en.wikipedia.org/wiki/Conditional\\_independence](https://en.wikipedia.org/wiki/Conditional_independence)

# Conditional independence notation

---

- **V is conditionally independent of set  $Vi$  given set  $Vj$** 
  - if  $p(V | Vi, Vj) = p(V | Vj)$
  - notation:  $I(V, Vi | Vj)$  or  $V \perp Vi | Vj$
- **Intuition**
  - if  $I(V, Vi | Vj)$  then knowing  $Vi$  &  $Vj$  tells nothing more about  $V$  than knowing  $Vj$  alone
  - if we know  $Vj$  we can ignore  $Vi$

$$P(Queen | Club) = \frac{P(Queen, Club)}{P(Club)} = \frac{\frac{1}{52}}{\frac{13}{52}} = \frac{1}{13}$$

$$P(Queen) = \frac{4}{52} = \frac{1}{13}$$

# Pairwise Independence

---

- Single  $V_i$  conditionally independent of single  $V_j$  given  $\mathbf{V}$ 
  - that is,  $I(V_i, V_j | \mathbf{V})$
- From definitions we have that
  - $p(V_i | V_j, \mathbf{V}) = p(V_i | \mathbf{V})$  and
  - $p(V_i, V_j | \mathbf{V}) = p(V_i | V_j, \mathbf{V}) p(V_j | \mathbf{V})$   
(Product Rule:  $P(A, B) = P(A | B)P(B)$  )
- Thus
  - $p(V_i, V_j | \mathbf{V}) = p(V_i | \mathbf{V}) p(V_j | \mathbf{V})$
  - $V_i$  and  $V_j$  is pairwise independent

# Independence: Pairwise and Conditional

---

- **Pairwise Independence**

- $P(A|C) = Pr(A)$
- $P(A,B|C)=P(A|C)P(B|C)$  Decompose joint probs using chain rule
  - Since  $P(A,B|C) = P(A|B,C)P(B|C)$  [Chain rule]
  - and  $P(A|B, C)= P(A|C)$

- **Conditional Independence**

- $P(A|B)=P(A)$
- $P(A|B, C)= P(A|C)$
- $P(A, B)=P(A)P(B)$

# Independence: Pairwise and Conditional

---

- **Pairwise Independence**

- $P(A|C) = \Pr(A)$
- $P(A,B|C) = P(A|C)P(B|C)$  Decompose joint probs using chain rule
  - Since  $P(A,B|C) = P(A|B,C)P(B|C)$  [Chain rule]
  - and  $P(A|B, C) = P(A|C)$

- **Conditional Independence**

- $P(A|B) = P(A)$
- $P(A|B, C) = P(A|C)$
- $P(A, B) = P(A)P(B)$  Only a skip away from Naïve Bayes

# Probability Basics → Naïve Bayes Models

- **Introduction**
- **Probability Basics**
  - Probability Axioms
  - Conditional probabilities
  - Product Rule, Chain Rule, Bayes Rule
- **Bayes Nets And Naïve Bayes**
  - Learning
  - Independence
  - Conditional independence
  - Naïve Bayes derivation (discrete case plus smoothing)
- **Naïve Bayes Flavors**
  - Continuous input variables
  - Discrete input variables (2 flavors: Bernoulli, multinomial)
- **Case Study: Spam detector in Naïve Bayes**
- **Summary**

# Bayes Classifier: Simplify RHS

---

$$P(Y = y_k | X_1, X_2, \dots, X_N) = \frac{P(Y=y_k)P(X_1, X_2, \dots, X_N | Y=y_k)}{\sum_j P(Y=y_j)P(X_1, X_2, \dots, X_N | Y=y_j)}$$

$$P(Y = y_k | X_1, X_2, \dots, X_N) = \frac{P(Y=y_k)P(X_1, X_2, \dots, X_N | Y=y_k)}{\sum_j \cancel{P(Y=y_j)} \cancel{P(X_1, X_2, \dots, X_N | Y=y_j)}}$$

- The denominator is common to all classes  $y_k$  so set to 1

$$Y \leftarrow \operatorname{argmax}_{y_k} P(Y = y_k) \cdot \underline{P(Y=y_k)P(X_1, X_2, \dots, X_N | Y=y_k)}$$

argmax\_yk means find the value of yk that maximises the expression

# Bayes Classifier: Simplify RHS

---

$$P(Y = y_k | X_1, X_2, \dots, X_N) = \frac{P(Y=y_k)P(X_1, X_2, \dots, X_N | Y=y_k)}{\sum_j P(Y=y_j)P(X_1, X_2, \dots, X_N | Y=y_j)}$$

$$P(Y = y_k | X_1, X_2, \dots, X_N) = \frac{P(Y=y_k)P(X_1, X_2, \dots, X_N | Y=y_k)}{\sum_j \cancel{P(Y=y_j)P(X_1, X_2, \dots, X_N | Y=y_j)}}$$

- The denominator is common to all classes  $y_k$  so set to 1

$$Y \leftarrow \operatorname{argmax}_{y_k} P(Y = y_k) \boxed{P(X_1, X_2, \dots, X_N | Y=y_k)}$$

**BUT  $P(X_1, X_2, \dots, X_N | Y=y_k)$  is still complex**

# Simplify the class conditional term using the product rule and conditional independence

$$Y \leftarrow \operatorname{argmax}_{y_k} P(Y = y_k) P(X_1, X_2, \dots, X_N | Y = y_k)$$

$$P(X_1, X_2, \dots, X_N | Y = y_k)$$

X is vector of  $\langle X_1, X_2 \rangle$

Using the Product Rule  
Simplify the Class Cond'l term

$$\begin{aligned} P(X|Y) &= P(X_1, X_2|Y) & P(A, B) = P(A|B)P(B) \\ &= P(X_1|X_2, Y)P(X_2|Y) & \text{Assume cond'l independence.} \\ &= P(X_1|Y)P(X_2|Y) & \text{Naïve Bayes} \end{aligned}$$

Independence

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

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

$$P(X_1 | X_2, Y) = P(X_1 | Y)$$

Where

- the second line follows from a general property of probabilities (product Rule),
- the third line follows directly from our above definition of conditional independence.

# Key Slide

## Derive NB Algorithm

$$\begin{aligned} P(Y = y_k | X_1, X_2, \dots, X_N) &= \frac{P(Y=y_k)P(X_1, X_2, \dots, X_N | Y=y_k)}{\sum_j P(Y=y_j)P(X_1, X_2, \dots, X_N | Y=y_j)} \\ &= \frac{P(Y=y_k)\prod_i P(X_i | Y=y_k)}{\sum_j P(Y=y_j)\prod_i P(X_i | Y=y_j)} \end{aligned}$$

The Naive Bayes algorithm is a classification algorithm based on Bayes rule, that assumes the attributes  $X_1 \dots X_n$  are all conditionally independent of one another, given  $Y$ . The value of this assumption is that it dramatically simplifies the representation of  $P(X|Y)$ , and the problem of estimating it from the training data. Consider, for example, the case where  $X = \langle X_1, X_2 \rangle$ . In this case

### Independence

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

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

$$P(X_1 | X_2, Y) = P(X_1 | Y)$$

$X$  is vector of  $\langle X_1, X_2 \rangle$

$$\begin{aligned} P(X|Y) &= P(X_1, X_2|Y) \quad \text{Use Product Rule} \quad P(A, B) = P(A | B)P(B) \\ &= P(X_1|X_2, Y)P(X_2|Y) \\ &= P(X_1|Y)P(X_2|Y) \quad \text{Naïve Bayes} \end{aligned}$$

Where

- the second line follows from a general property of probabilities (product Rule),
- the third line follows directly from our above definition of conditional independence.

# Bayesian Classifier → Naïve Bayes

---

$$Y \leftarrow \operatorname{argmax}_{y_k} P(Y = y_k) \underline{P(X_1, X_2, \dots, X_N | Y = y_k)}$$

Multivariate Joint Prob → Production of univariate class conditionals

$2^n \rightarrow n$ . Probabilities to be learned from data for each class

$$Y \leftarrow \operatorname{argmax}_{y_k} P(Y = y_k) \prod_i P(X_i | Y = y_k)$$

argmax\_yk means find the value of yk that maximises the expression

# Conditional independence dramatically reduces model complexity: to $n$ parameters from $2^n$

---

More generally, when  $X$  contains  $n$  attributes which are conditionally independent of one another given  $Y$ , we have

$$P(X_1 \dots X_n | Y) = \prod_{i=1}^n P(X_i | Y) \quad (1)$$

Notice that when  $Y$  and the  $X_i$  are boolean variables, we need only  $2n$  parameters to define  $P(X_i = x_{ik} | Y = y_j)$  for the necessary  $i, j, k$ . This is a dramatic reduction compared to the  $2(2^n - 1)$  parameters needed to characterize  $P(X|Y)$  if we make no conditional independence assumption.

- **High Bias**
- **Assume no interactions between the variables**

# Naïve Bayes: summary so far

---

- A generative, parametric model
- Computes the conditional a-posterior probabilities of a categorical class variable given independent predictor variables using the Bayes rule.
  - So far focused on discrete input variables
  - See below how to model continuous input variables

# Probability Basics → Naïve Bayes Models

- **Introduction**
- **Probability Basics**
  - Probability Axioms
  - Conditional probabilities
  - Product Rule, Chain Rule, Bayes Rule
- **Bayes Nets And Naïve Bayes**
  - Learning
  - Independence
  - Conditional independence
  - Naïve Bayes derivation (discrete case plus smoothing)
- **Naïve Bayes Flavors**
  - Discrete input variables (2 flavors: Bernoulli, multinomial)
  - Continuous input variables
- **Case Study: Spam detector in Naïve Bayes**
- **Summary**

# Probability Basics → Naïve Bayes Models

- **Introduction**
- **Probability Basics**
  - Probability Axioms
  - Conditional probabilities
  - Product Rule, Chain Rule, Bayes Rule
- **Bayes Nets And Naïve Bayes**
  - Learning
  - Independence
  - Conditional independence
  - Naïve Bayes derivation (discrete case plus smoothing)
- **Naïve Bayes Flavors**
  - Discrete input variables (2 flavors: Bernoulli, multinomial)
  - Continuous input variables
- **Case Study: Spam detector in Naïve Bayes**
- **Summary**

# Parameter estimation for Naïve Bayes

See IRBook [Manning et al. 2008] Chapter 13 for more info

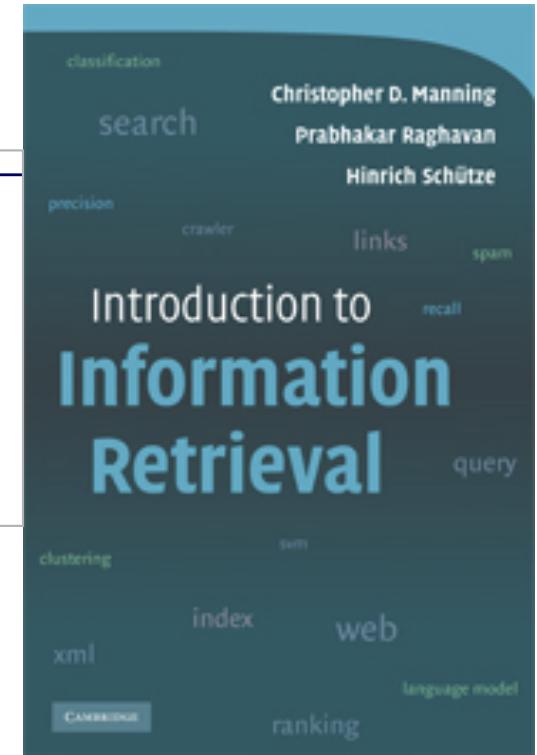
- **Multivariate Bernoulli model:**  $X_w = t$  if word present document

$$\hat{P}(X_w = t \mid c_j) = \text{fraction of documents of topic } c_j \text{ in which word } w \text{ appears}$$

- **Multinomial model:**

$$\hat{P}(X_i = w \mid c_j) = \text{fraction of times in which word } w \text{ appears across all documents of topic } c_j$$

- Can create a mega-document for topic  $j$  by concatenating all documents on this topic
- Use frequency of  $w$  in mega-document



# Naïve Bayes for text

PDF and HTML versions of Chapter 13 are available here  
[Christopher D. Manning, Prabhakar Raghavan and Hinrich Schütze, Introduction to Information Retrieval, Cambridge University Press. 2008.] <http://nlp.stanford.edu/IR-book/>

<http://www.cs.cmu.edu/~epxing/Class/10701-10s/Lecture/lecture5.pdf>

[http://www.cs.cmu.edu/~tom/10701\\_sp11/lectures.shtml](http://www.cs.cmu.edu/~tom/10701_sp11/lectures.shtml)

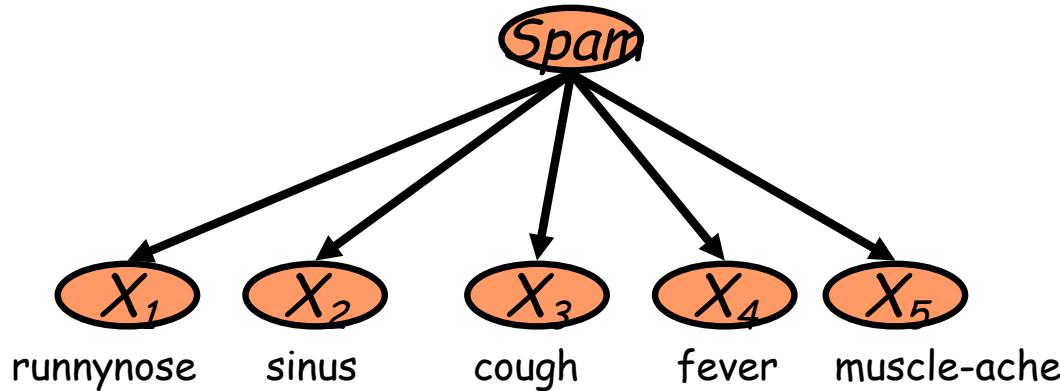
# Reading material

---

- **PDF and HTML versions of Chapter 13 are available here**
  - [Christopher D. Manning, Prabhakar Raghavan and Hinrich Schütze, Introduction to Information Retrieval, Cambridge University Press. 2008.] <http://nlp.stanford.edu/IR-book/>
  - <http://www.cs.cmu.edu/~epxing/Class/10701-10s/Lecture/lecture5.pdf>
  - [http://www.cs.cmu.edu/~tom/10701\\_sp11/lectures.shtml](http://www.cs.cmu.edu/~tom/10701_sp11/lectures.shtml)

# The Naïve Bayes Classifier

IRBook 13.3



- **Conditional Independence Assumption:** Features (term presence) are *independent* of each other given the class:

$$P(X_1, \dots, X_5 | C) = P(X_1 | C) \bullet P(X_2 | C) \bullet \dots \bullet P(X_5 | C)$$

# Learn NB Model as a dict with two columns

► Table 13.1 Data for parameter estimation examples.

	docID	words in document	in $c = China$ ?
training set	1	Chinese Beijing Chinese	yes
	2	Chinese Chinese Shanghai	yes
	3	Chinese Macao	yes
	4	Tokyo Japan Chinese	no
test set	5	Chinese Chinese Chinese Tokyo Japan	?

6 is the vocabulary size

Word	Pr(Word Class)	Pr(Word NotClass)
CLASSPRIORs	$\frac{3}{4} = 0.75$	$\frac{1}{4} = 0.25$
Tokyo	$0/(8) = 0$	$1/3$
Chinese	$5/8$	$1/3$
....		

# Use model to predict class of document

► Table 13.1 Data for parameter estimation examples.

	docID	words in document	in $c = China$ ?
training set	1	Chinese Beijing Chinese	yes
	2	Chinese Chinese Shanghai	yes
	3	Chinese Macao	yes
	4	Tokyo Japan Chinese	no
test set	5	Chinese Chinese Chinese Tokyo Japan	?

6 is the vocabulary size

Word	Pr(Word Class)	Pr(Word NotClass)
CLASSPRIORS	$\frac{3}{4} = 0.75$	$\frac{1}{4} = 0.25$
Tokyo	$0/(8) = 0$	$1/3$
Chinese	$5/8$	$1/3$
....		Classify Document={Tokyo, Chinese} $Pr(Tokyo Class) = 0/8$

$$Y \leftarrow argmax_{y_k} P(Y = y_k) \prod_i P(X_i | Y = y_k)$$

# Model file can be a dict with two columns

► Table 13.1 Data for parameter estimation examples.

	docID	words in document	in $c = China$ ?
training set	1	Chinese Beijing Chinese	yes
	2	Chinese Chinese Shanghai	yes
	3	Chinese Macao	yes
	4	Tokyo Japan Chinese	no
test set	5	Chinese Chinese Chinese Tokyo Japan	?

6 is the vocabulary size

Word	Pr(Word Class)	Pr(Word NotClass)
CLASSPRIORS	$\frac{3}{4} = 0.75$	$\frac{1}{4} = 0.25$
Tokyo	$0/(8) = 0$	$1/3$
Chinese	$5/8$	$1/3$
....		Classify Document={Tokyo, Chinese} $Pr(Tokyo Class) = 0/8$

$$Y \leftarrow argmax_{y_k} P(Y = y_k) \prod_i P(X_i | Y = y_k)$$

Run into trouble during classification if we have no estimate  $\theta_{ijk}$  for variable  $X_i$  with value j for class k. The whole expression go to 0 for that class k.

# Model file can be a dict with two columns

► Table 13.1 Data for parameter estimation examples.

	docID	words in document	in <i>c</i> = China?
training set	1	Chinese Beijing Chinese	yes
	2	Chinese Chinese Shanghai	yes
	3	Chinese Macao	yes
	4	Tokyo Japan Chinese	no
test set	5	Chinese Chinese Chinese Tokyo Japan	?

Plus 1 smoothing

+6 is the vocabulary size  
Prior 1/6 for each word  
For both classes

Word	Pr(Word Class)	Pr(Word NotClass)
CLASSPRIORS	$\frac{3}{4} = 0.75$	$\frac{1}{4} = 0.25$
Tokyo	$0+1/(8+6) = 1/14$	$1+1/(3+6) = 2/9$
Chinese	$5+1/(8+6) = 6/14$	$1+1/(3+6) = 2/9$
.....		

$$\text{Pr}(Tokyo|\text{Class}) = 0/8 \rightarrow 0+1/(8+6)$$

# Laplace Smoothing

► Table 13.1 Data for parameter estimation examples.

	docID	words in document	in $c = China$ ?
training set	1	Chinese Beijing Chinese	yes
	2	Chinese Chinese Shanghai	yes
	3	Chinese Macao	yes
	4	Tokyo Japan Chinese	no
test set	5	Chinese Chinese Chinese Tokyo Japan	?

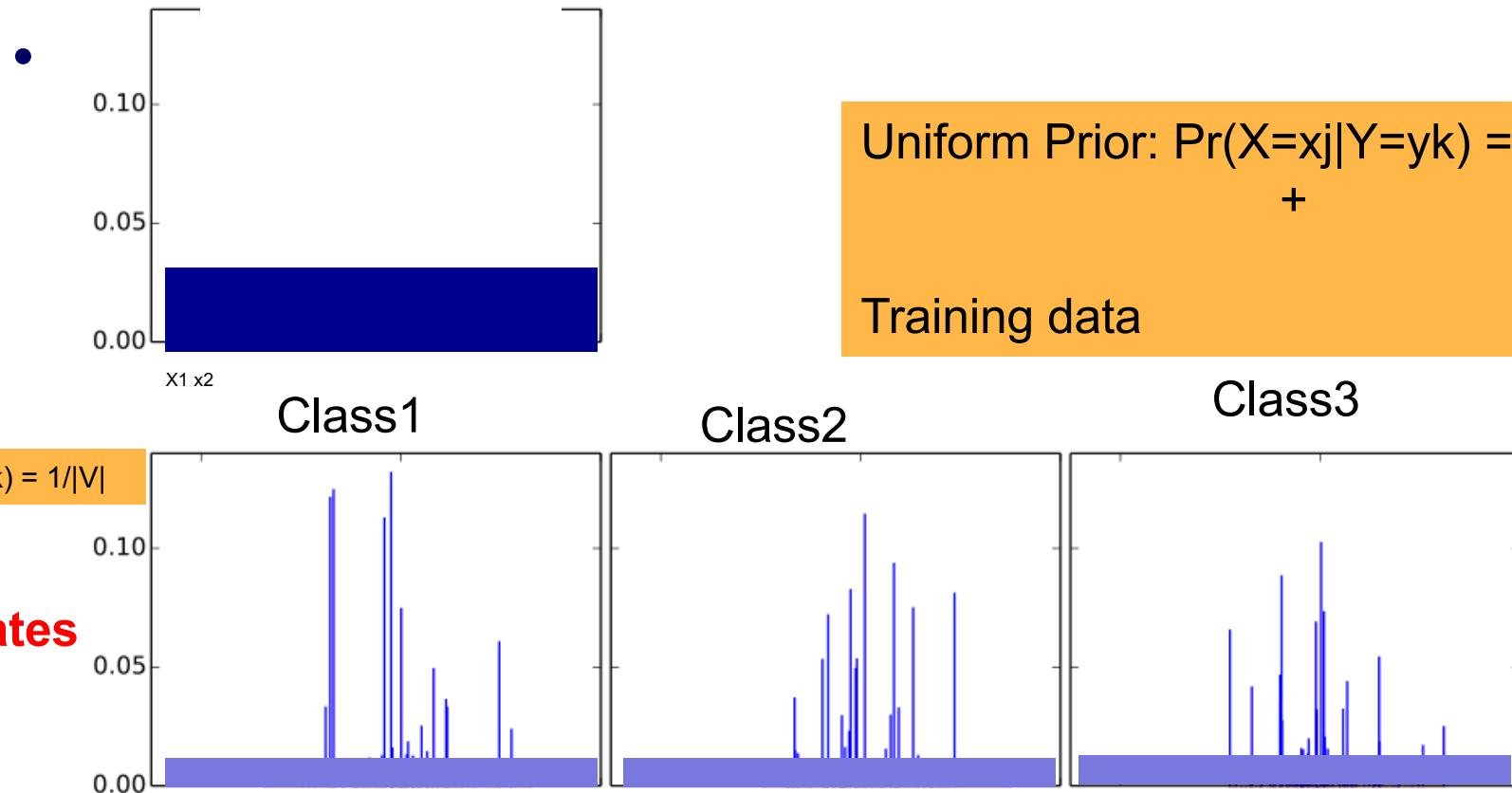
## Plus 1 smoothing

+6 is the vocabulary size  
Prior 1/6 for each word  
For both classes

Word	Pr(Word Class)	Pr(Word NotClass)	
CLASSPRIORS	$\frac{3}{4} = 0.75$	$\frac{1}{4} = 0.25$	
Tokyo	$0+1/(8+6) = 1/14$	$1+1/(3+6) = 2/9$	<b>MAP Estimates</b>
Chinese	$5+1/(8+6) = 6/14$	$1+1/(3+6) = 2/9$	Learning a NB Via Bayesian hierarchical model
.....			

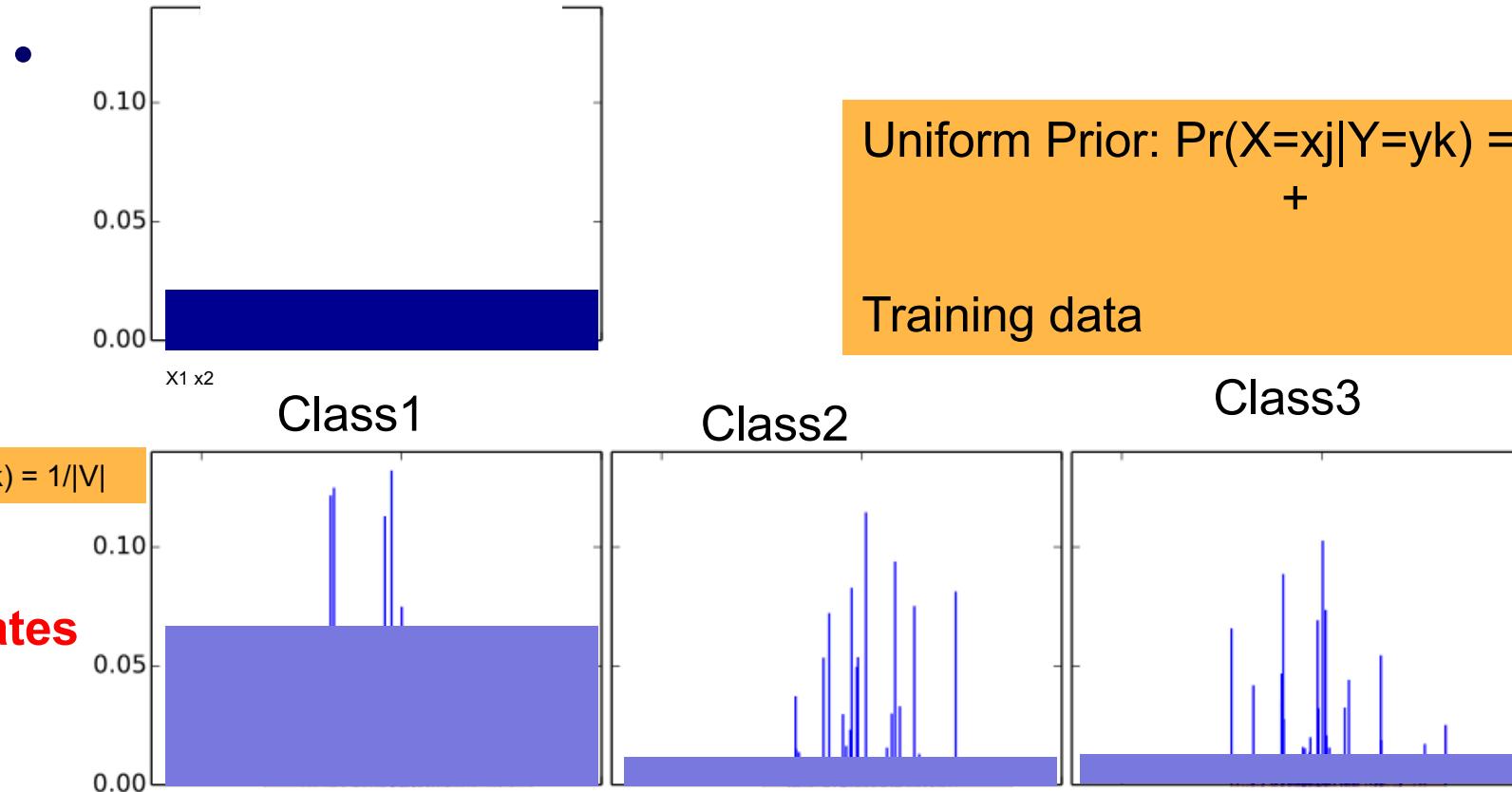
$$\text{Pr}(Tokyo|Class) = 0/8 \rightarrow 0+1/(8+6)$$

# MAP Estimates: Smoothed Probabilities



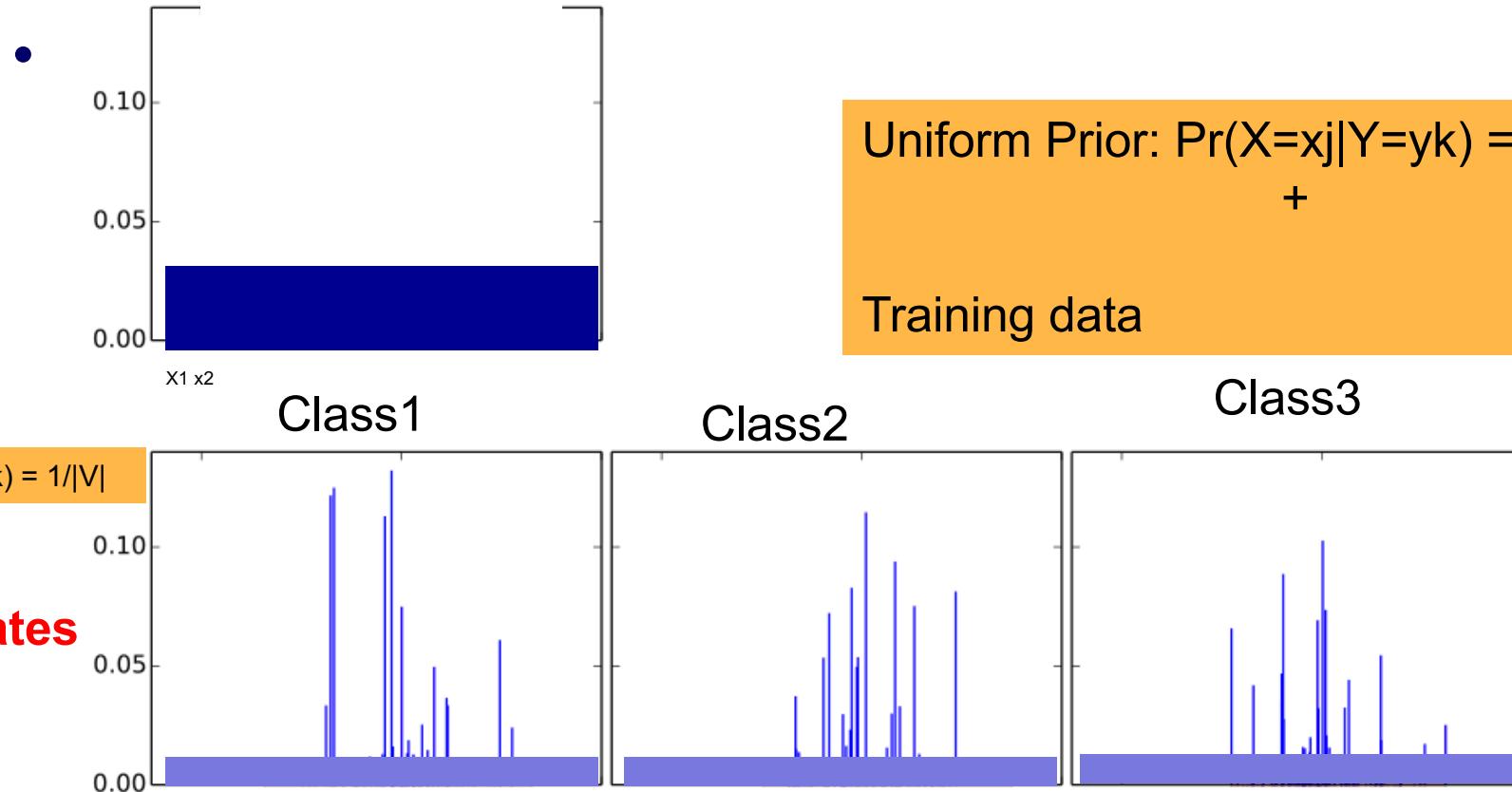
Huge Probability mass assigned to background model:  
Bias versus variance?

# MAP Estimates: Smoothed Probabilities



Huge Probability mass assigned to background model:  
Bias versus variance?

# MAP Estimates: Smoothed Probabilities



Huge Probability mass assigned to background model:  
Bias versus variance? High Bias!

# Class Prior Estimates

---

Maximum likelihood estimates for  $\pi_k$  are

$$\hat{\pi}_k = \hat{P}(Y = y_k) = \frac{\#D\{Y = y_k\}}{|D|} \quad (8)$$

where  $|D|$  denotes the number of elements in the training set  $D$ .

Alternatively, we can obtain a smoothed estimate, or equivalently a MAP estimate based on a Dirichlet prior over the  $\pi_k$  parameters assuming equal priors on each  $\pi_k$ , by using the following expression

$$\hat{\pi}_k = \hat{P}(Y = y_k) = \frac{\#D\{Y = y_k\} + l}{|D| + lK} \quad (9)$$

where  $K$  is the number of distinct values  $Y$  can take on, and  $l$  again determines the strength of the prior assumptions relative to the observed data  $D$ .

# Probability Basics → Naïve Bayes Models

- **Introduction**
- **Probability Basics**
  - Probability Axioms
  - Conditional probabilities
  - Product Rule, Chain Rule, Bayes Rule
- **Bayes Nets And Naïve Bayes**
  - Learning
  - Independence
  - Conditional independence
  - Naïve Bayes derivation (discrete case plus smoothing)
- **Naïve Bayes Flavors**
  - Discrete input variables (2 flavors: Bernoulli, multinomial)
  - Continuous input variables
- **Case Study: Spam detector in Naïve Bayes**
- **Summary**

# Naive Bayes for Continuous Inputs

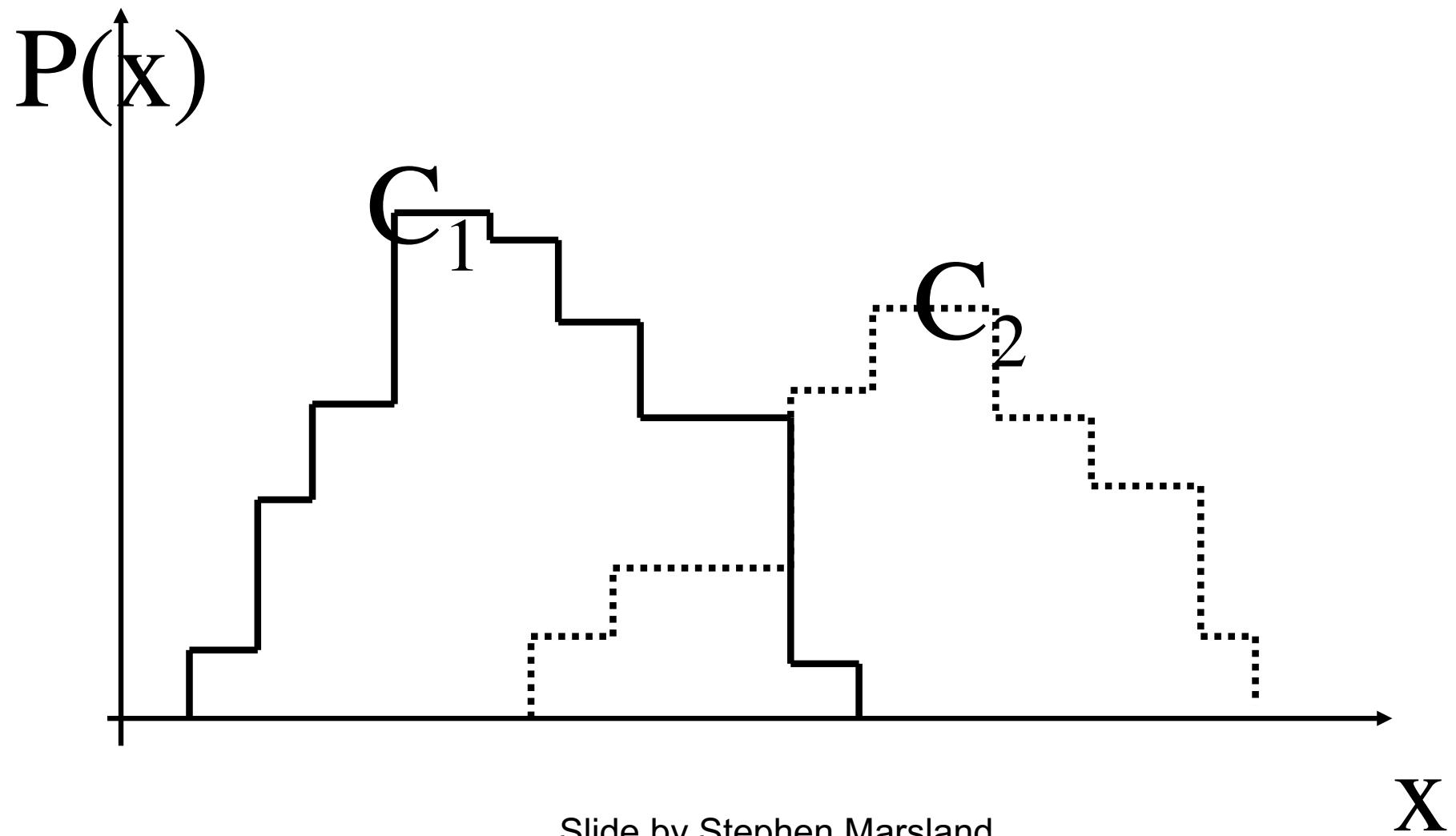
---

- When the  $X_i$  are continuous we must choose some other way to represent the distributions
  - $P(X_i|Y) = \text{Gaussian}(\mu, \sigma)$ .
- Mix discrete variables with continuous variables
  - One common approach is to assume that for each possible discrete value  $y_k$  of  $Y$ , the distribution of each continuous  $X_i$  is Gaussian, and is defined by a mean and standard deviation specific to  $X_i$  and  $y_k$ .
- In order to train such a Naïve Bayes classifier we must therefore estimate the mean and standard deviation of each of these Gaussians:

$$\mu_{ik} = E[X_i | Y = y_k]$$

$$\sigma_{ik}^2 = E[(X_i - \mu_{ik})^2 | Y = y_k]$$

# Feature Histograms



Slide by Stephen Marsland

# MLE-based Estimates for $\mu$ and $\sigma$

---

Continuous Inputs:  $\Pr(Y|X) \sim N(\mu, \sigma^2)$

- Again, we can use either maximum likelihood estimates (MLE) or maximum a posteriori (MAP) estimates for these parameters. The maximum likelihood estimator for  $\mu_{ik}$  is

$$\hat{\mu}_{ik} = \frac{1}{\sum_j \delta(Y^j = y_k)} \sum_j X_i^j \delta(Y^j = y_k) \quad (13)$$

where the superscript  $j$  refers to the  $j$ th training example, and where  $\delta(Y = y_k)$  is 1 if  $Y = y_k$  and 0 otherwise. Note the role of  $\delta$  here is to select only those training examples for which  $Y = y_k$ .

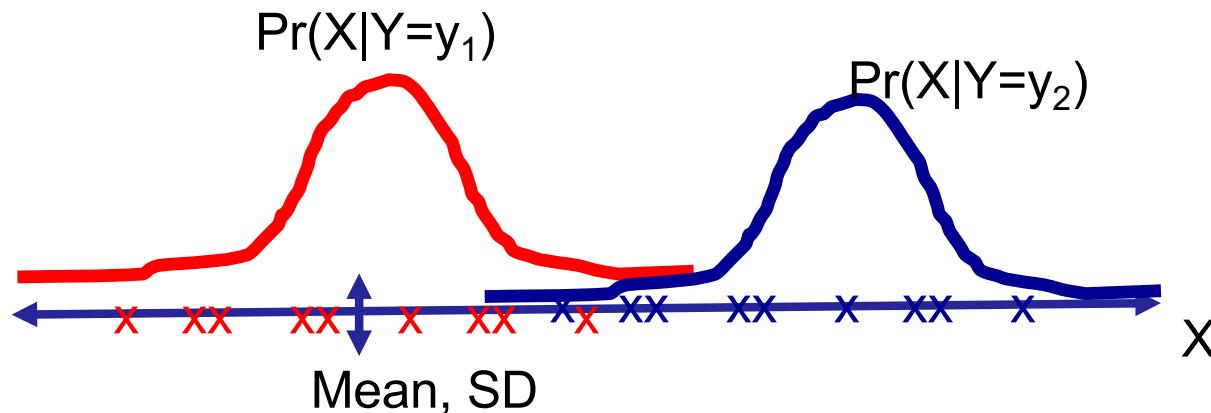
The maximum likelihood estimator for  $\sigma_{ik}^2$  is

$$\hat{\sigma}_{ik}^2 = \frac{1}{\sum_j \delta(Y^j = y_k)} \sum_j (X_i^j - \hat{\mu}_{ik})^2 \delta(Y^j = y_k) \quad (14)$$

# Estimate $\mu$ , $\sigma$ from data

$$\mu_{ik} = E[X_i | Y = y_k]$$

$$\sigma_{ik}^2 = E[(X_i - \mu_{ik})^2 | Y = y_k]$$



# Continuous Inputs: $\Pr(Y|X) \sim N(\mu, \sigma^2)$

---

If  $X$  is Normally distributed with mean  $\mu$  and standard deviation  $\sigma$ , we write

$$X \sim N(\mu, \sigma^2)$$

$\mu$  and  $\sigma$  are the **parameters** of the distribution.

---

The probability density of the Normal distribution is given by

$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} \exp^{-(x-\mu)^2/2\sigma^2}$$

For the purposes of this course we do not need to use this expression. It is included here for future reference.

# Naïve Bayes

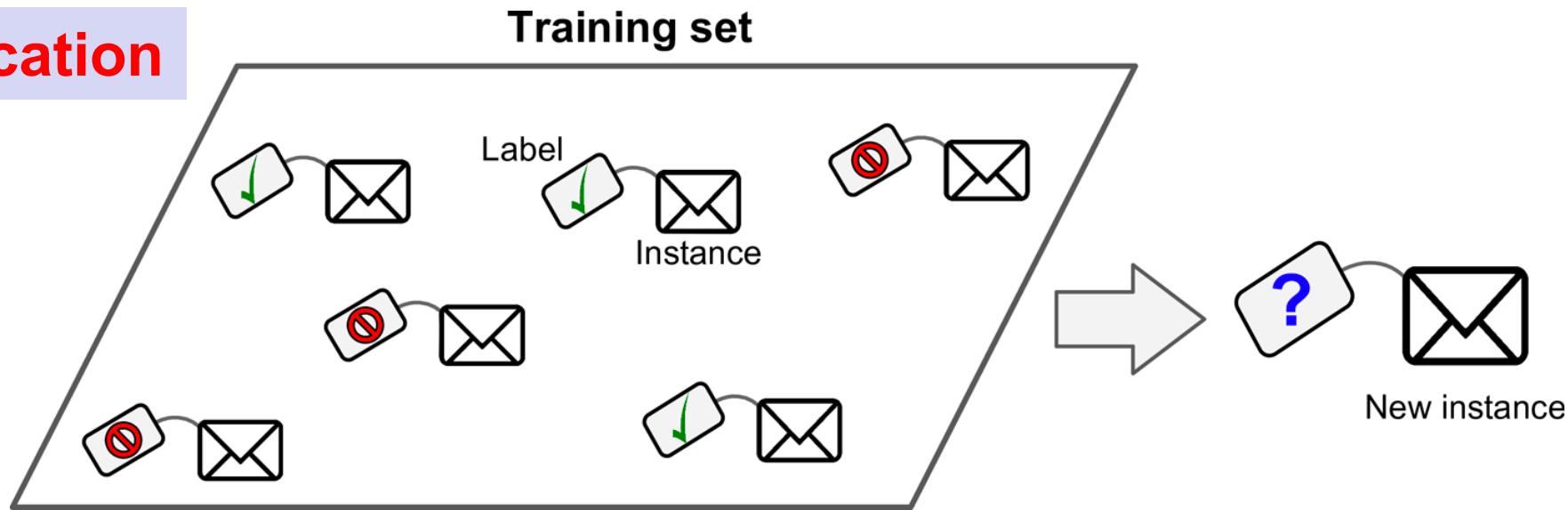
---

- **Combine discrete input variables with continuous input variables?**
  - Naïve Bayes, Decision trees
- **YES**
- **Whereas these requires feature transformations**
  - Logistic regression: one hot-encoding

# Probability Basics → Naïve Bayes Models

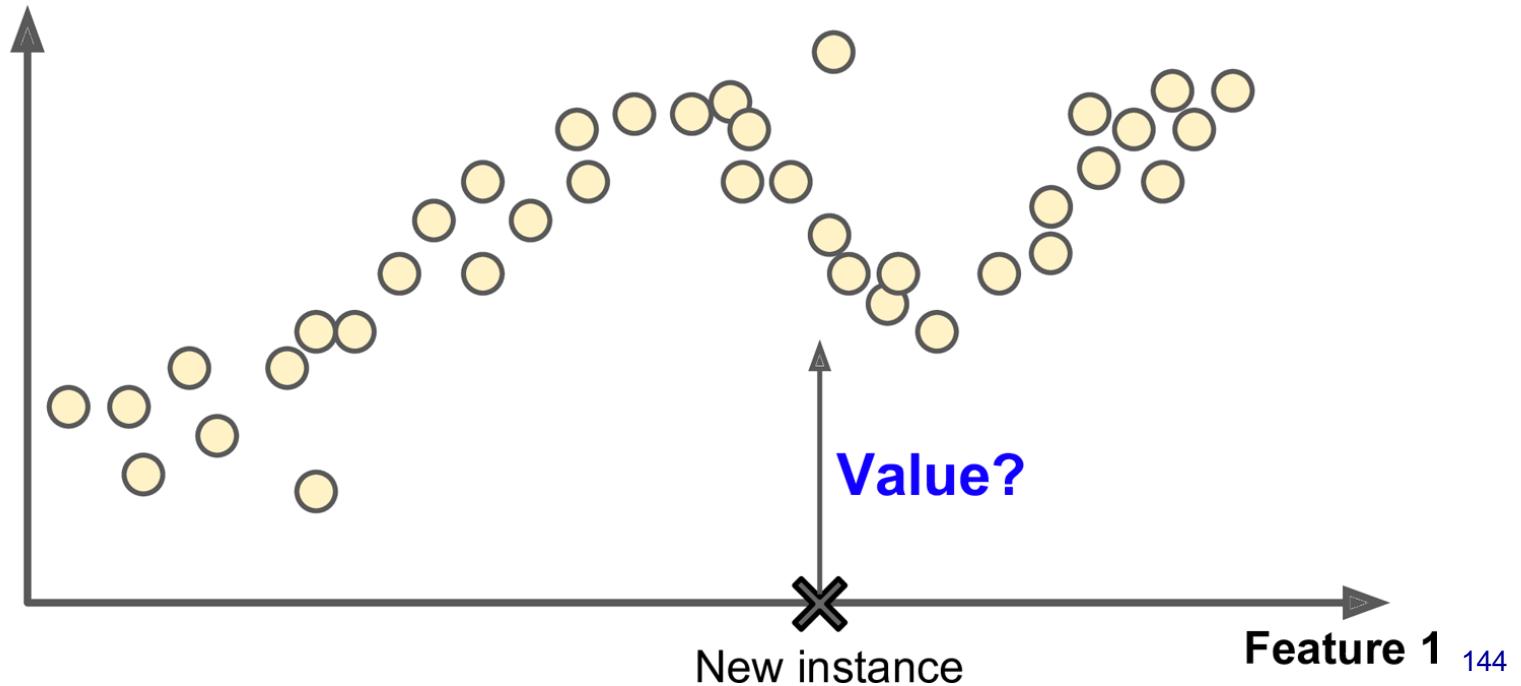
- **Introduction**
- **Probability Basics**
  - Probability Axioms
  - Conditional probabilities
  - Product Rule, Chain Rule, Bayes Rule
- **Bayes Nets And Naïve Bayes**
  - Learning
  - Independence
  - Conditional independence
  - Naïve Bayes derivation (discrete case plus smoothing)
- **Naïve Bayes Flavors**
  - Discrete input variables (2 flavors: Bernoulli, multinomial)
  - Continuous input variables
- **Case Study: Spam detector in Naïve Bayes**
- **Summary**

## Classification



Value

## Regression



# Email: SPAM versus Ham

The screenshot shows a Gmail inbox search results page with the query "in:spam". The results list several spam messages from various senders, each with a preview of the email content.

From	Subject
Ivan	Looking forward to see you in Xiamen,China for attending MMSTA2017
MR. BENJAMIN UDEVIS	Attention - CONTACT DHL COURIER COMPANY LIMITED Attention I am
Canadian-Pharmacy-24h	Free shipping on any order of \$40 or more! Hurry up to buy best qualit
Fr. John Fontana	Send us your All Souls petitions - If you can't see this message please us
UNITED BANK FOR AFRICA	VERY URGNET ATTENTION - BARRISTER MAURICE DEPARTMENT ATI
Credit one card	You may qualify for a platinumvisa - Follow these instructions to opt-out h
ALLPOINTS	Reservas em hotéis? Cadastre-se no programa ALLPOINTS e ganhe ag
BROWN Z.M	30/9/2017 acknowledge receipt - Good day, This is a letter of Intent for Inv
STEVE PATEMAN	Notification. - Ordering beneficiary, Be informed that we've been able to fru
Mr.DR WILLIAMS MORGAN	Attn, beneficiary - Attn, beneficiary This is to notify you that your fund has l
Peter Gaynor	ATTENTION ATTENTION ATTENTION!!! - Dear Friend, Greetings,I know y
Carlos Slim Helu	Donation for you!! - GREETINGS, My name is Carlos Slim Helu, A philant
Editor IJETTCS	Research Article are invited for Best Quality Journal IJETTCS (www.ije

# Email: SPAM versus Ham

The screenshot shows a Gmail inbox search results page for the query "in:spam". The search bar at the top contains "in:spam". Below the search bar, there is a red "COMPOSE" button. To the left, there is a sidebar with links to "Inbox (86,800)", "Important", "Sent Mail", and "Drafts (1,743)". The main area displays several spam messages:

Word	Pr(Word Class)	Pr(Word NotClass)
CLASSPRIORs	$\frac{3}{4} = 0.75$	$\frac{1}{4} = 0.25$
Tokyo	$0+1/(8+6) = 1/14$	$1+1/(3+6) = 2/9$
Chinese	$5+1/(8+6) = 6/14$	$1+1/(3+6) = 2/9$
....		

The first message is from "Ivan" with the subject "Looking forward to see you in Xiamen, China for attending MMSTA2017". The second message is from "MR. BENJAMIN UDEVIS" with the subject "Attention - CONTACT DHL COURIER COMPANY LIMITED Attention I am". The third message is from "Canadian-Pharmacy-24h" with the subject "Free shipping on any order of \$40 or more! Hurry up to buy best qualit". The fourth message is from "Fr. John Fontana" with the subject "Send us your All Souls petitions - If you can't see this message please us".

# Probability Basics → Naïve Bayes Models

- **Introduction**
- **Probability Basics**
  - Probability Axioms
  - Conditional probabilities
  - Product Rule, Chain Rule, Bayes Rule
- **Bayes Nets And Naïve Bayes**
  - Learning
  - Independence
  - Conditional independence
  - Naïve Bayes derivation (discrete case plus smoothing)
- **Naïve Bayes Flavors**
  - Discrete input variables (2 flavors: Bernoulli, multinomial)
  - Continuous input variables
- **Case Study: Spam detector in Naïve Bayes**
- **Summary**

# Summary

---

## ■ Bayesian prediction:

- requires solving density estimation problems.
- often difficult to estimate  $\Pr[\mathbf{x} \mid y]$  for  $\mathbf{x} \in \mathbb{R}^N$ .
- but, simple and easy to apply; widely used.

## ■ Naive Bayes:

- strong assumption.
- straightforward estimation problem.
- specific linear classifier.
- sometimes surprisingly good performance.

# **From complex class conditional joint probability distributions of the order $2^n$ to n**

- **From complex class conditional joint probability distributions of the order  $2^n$  to n**
- **$2^n$  Possibilities -> n combinations that we have to estimate class conditional probabilities for**

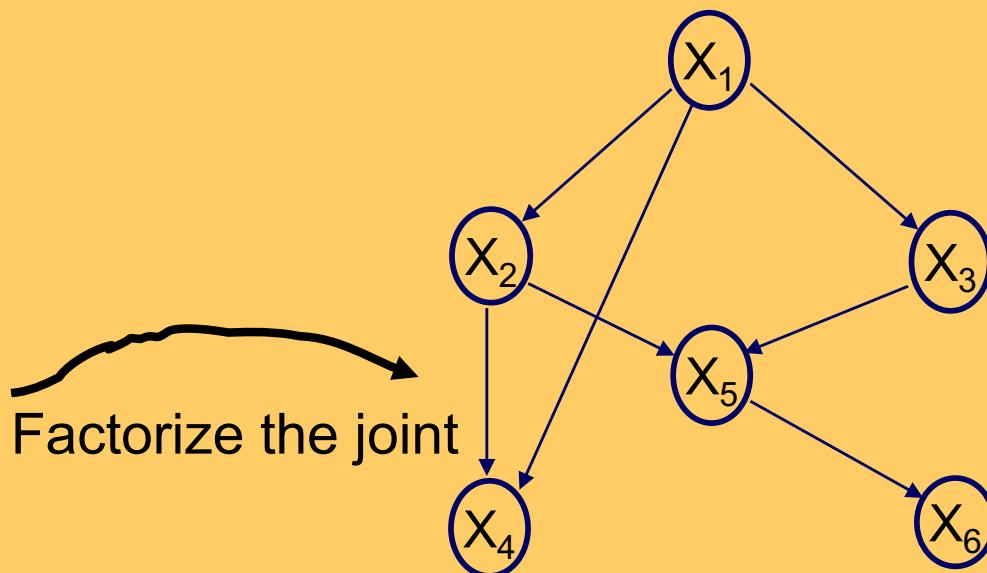
# Decompose Bayesian Networks

GOAL:  $2^n$  Possibilities  $\rightarrow n$

P: Joint Probability Distribution

#	X1	X2	X3	X4	X5	X6	Pr(X1,X2, X3, X4, X5, X6)
1	1						
2	1						
..	1						
4..	1						
5	1						
6	1						
7	1						
64							

G: Directed Acyclic Graph



$$p(x_1, x_2, x_3, x_4, x_5, x_6)$$

1. Partial Order

$$= p(x_1) p(x_2 | x_1) p(x_3 | x_1, x_2) p(x_4 | x_1, x_2, x_3) p(x_5 | x_4, x_3, x_2, x_1) p(x_6 | x_5, x_4, x_3, x_2, x_1)$$

$$= p(x_1) p(x_2 | x_1) p(x_3 | x_1) p(x_4 | x_2, x_1) p(x_5 | x_3, x_2) p(x_6 | x_5)$$

$$= p(x_1) p(x_2) p(x_3) p(x_4) p(x_5) p(x_6)$$

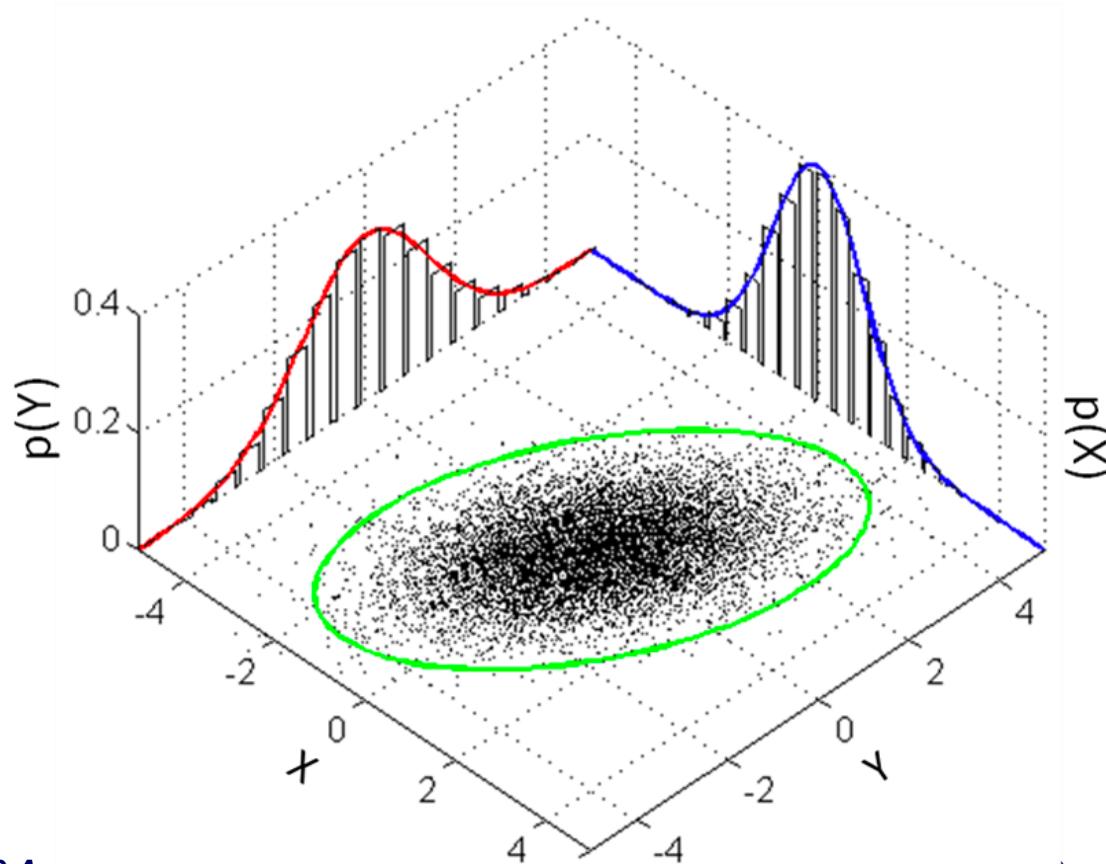
2. By Chain Rule

3. Markov Property

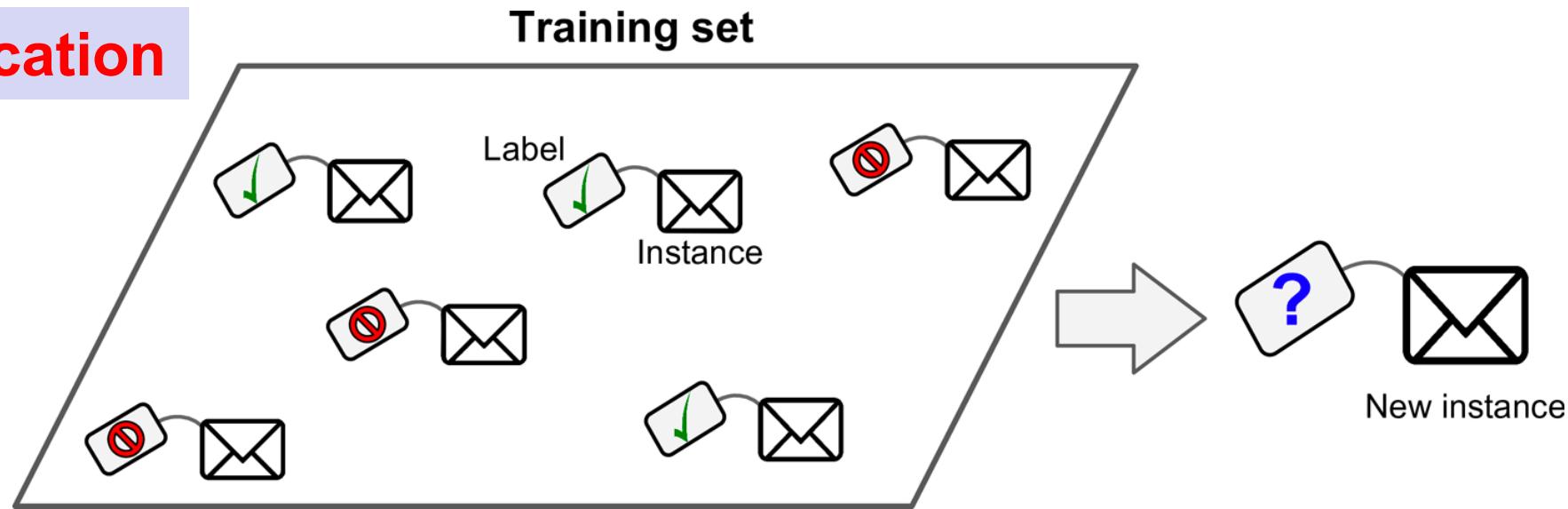
$$P(y|x_6, x_5, x_4, x_3, x_2, x_1) = p(x_1|y) p(x_2|y) p(x_3|y) p(x_4|y) p(x_5|y) p(x_6|y)$$

4: Independence (see next section)  
5: Naïve Bayes via  
Cond. Independence

- **Decompose large joint distributions of many variables into production many univariate distributions**

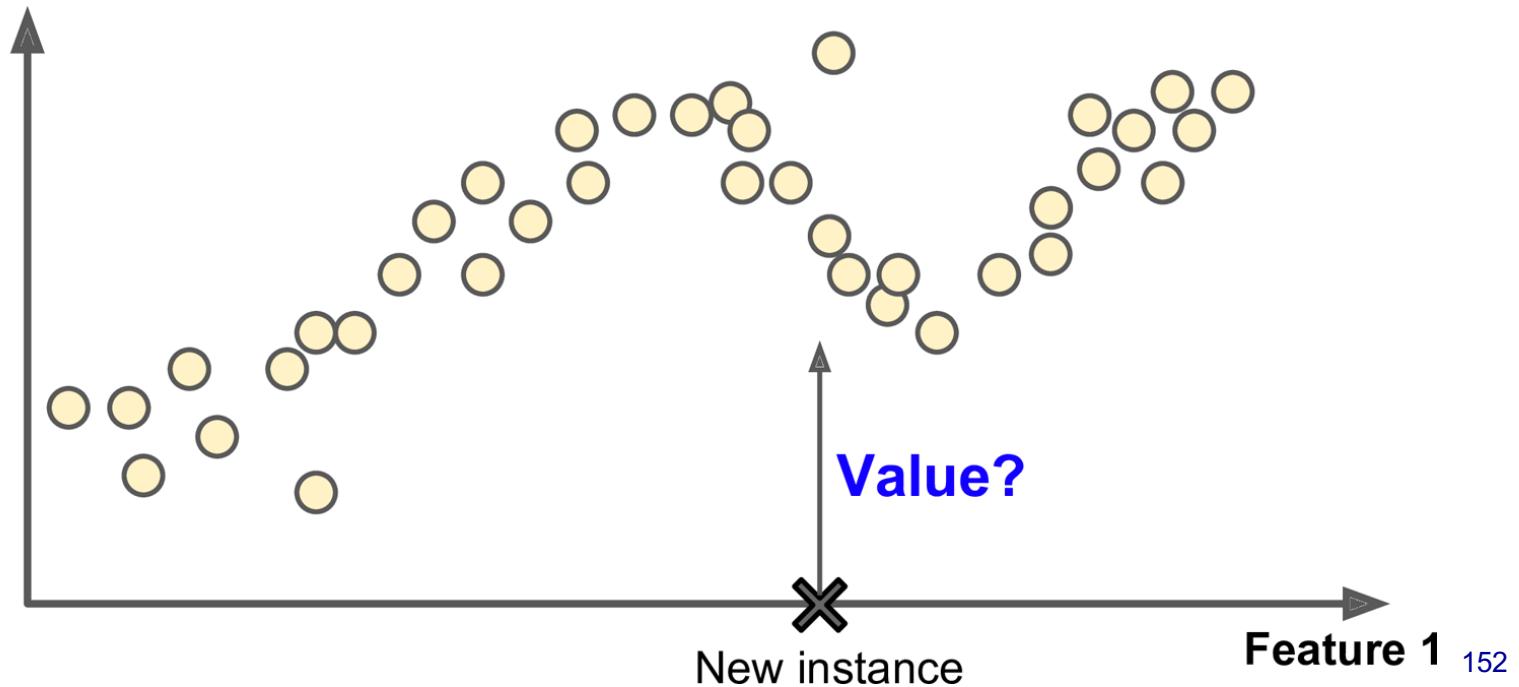


## Classification



Value

## Regression



# Comparison

assuming  $x$  in  $\{0, 1\}$

	Learning Objective	Training	Inference
Naïve Bayes	$\text{maximize} \sum_i \left[ \sum_j \log P(x_{ij}   y_i; \theta_j) + \log P(y_i; \theta_0) \right]$	$\theta_{kj} = \frac{\sum_i \delta(x_{ij} = 1 \wedge y_i = k) + r}{\sum_i \delta(y_i = k) + Kr}$	$\theta_1^T \mathbf{x} + \theta_0^T (1 - \mathbf{x}) > 0$ <p>where <math>\theta_{1j} = \log \frac{P(x_j = 1   y = 1)}{P(x_j = 1   y = 0)}</math>,</p> $\theta_{0j} = \log \frac{P(x_j = 0   y = 1)}{P(x_j = 0   y = 0)}$ 
Logistic Regression	$\text{maximize} \sum_i \log(P(y_i   \mathbf{x}, \boldsymbol{\theta})) + \lambda \ \boldsymbol{\theta}\ $ <p>where <math>P(y_i   \mathbf{x}, \boldsymbol{\theta}) = 1 / (1 + \exp(-y_i \boldsymbol{\theta}^T \mathbf{x}))</math></p>	Gradient ascent	$\boldsymbol{\theta}^T \mathbf{x} > 0$
Linear SVM	$\text{minimize} \lambda \sum_i \xi_i + \frac{1}{2} \ \boldsymbol{\theta}\ ^2$ <p>such that <math>y_i \boldsymbol{\theta}^T \mathbf{x} \geq 1 - \xi_i \quad \forall i</math></p>	Linear programming	$\boldsymbol{\theta}^T \mathbf{x} > 0$
Kernelized SVM	complicated to write	Quadratic programming	$\sum_i y_i \alpha_i K(\hat{\mathbf{x}}_i, \mathbf{x}) > 0$
Nearest Neighbor	most similar features $\rightarrow$ same label	Record data	$y_i$ <p>where <math>i = \operatorname{argmin}_i K(\hat{\mathbf{x}}_i, \mathbf{x})</math></p>

---



# End of lecture