

# Introduction to Machine Learning

Ranga Raju Vatsavai, Ph.D.

Chancellors Faculty Excellence Associate Professor in Geospatial Analytics  
Department of Computer Science, North Carolina State University (NCSU)

Feb. 25-27, 2019

## Probability - Basics

- Marginal Prob.

- Prob. of an event occurring,  $p(A)$ . It is not conditioned on another event.
- Prob. that a card drawn from deck is black: 0.5



- Joint Prob.

- Prob. of event A and B occurring,  $p(A \text{ and } B)$ . Its prob. of intersection of two or more events.
- Prob. of card is black and 5 is: 1/26

- Conditional Probability

- Prob. of event A occurring, given that event B occurs,  $p(A|B)$ .
- Say, you have drawn a black card, what's the prob. that's its 5.  $p(5|\text{black}) = 1/13$

## Probability - Basics

- Manipulation between conditional, joint, and marginal probabilities

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

- Event A is **independent** of B, if  $P(A) = P(A|B)$   
– Important consequence is multiplication rule

$$P(A \text{ and } B) = P(A|B)P(B) = P(A)P(B)$$

## Bayes Classifier

- A probabilistic framework for solving classification problems

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

- Conditional Probability:

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

- Bayes theorem:  $P(Y|X) = \frac{P(X|Y)P(Y)}{P(X)}$

## Example of Bayes Theorem

- Given:
  - A doctor knows that meningitis causes stiff neck 50% of the time
  - Prior probability of any patient having meningitis is 1/50,000
  - Prior probability of any patient having stiff neck is 1/20
- If a patient has stiff neck, what's the probability he/she has meningitis?

$$P(M|S) = \frac{P(S|M)P(M)}{P(S)} = \frac{0.5 \times 1/50000}{1/20} = 0.0002$$

## Using Bayes Theorem for Classification

- Consider each attribute and class label as random variables
- Given a record with attributes  $(X_1, X_2, \dots, X_d)$ 
  - Goal is to predict class Y
  - Specifically, we want to find the value of Y that maximizes  $P(Y|X_1, X_2, \dots, X_d)$
- Can we estimate  $P(Y|X_1, X_2, \dots, X_d)$  directly from data?

## Example Data

### Given a Test Record:

$X = (\text{Refund} = \text{No}, \text{Divorced}, \text{Income} = 120\text{K})$

Tid	Refund	Marital Status	Taxable Income	Evade
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	80K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

- Can we estimate

$P(\text{Evade} = \text{Yes} | X)$  and  $P(\text{Evade} = \text{No} | X)$ ?

In the following we will replace

Evade = Yes by Yes, and

Evade = No by No

## Using Bayes Theorem for Classification

- Approach:
  - compute posterior probability  $P(Y | X_1, X_2, \dots, X_d)$  using the Bayes theorem
- $$P(Y | X_1, X_2, \dots, X_d) = \frac{P(X_1, X_2, \dots, X_d | Y)P(Y)}{P(X_1, X_2, \dots, X_d)}$$
- Maximum a-posteriori:* Choose Y that maximizes  $P(Y | X_1, X_2, \dots, X_d)$
- Equivalent to choosing value of Y that maximizes  $P(X_1, X_2, \dots, X_d | Y) P(Y)$
- How to estimate  $P(X_1, X_2, \dots, X_d | Y)$ ?

## Example Data

**Given a Test Record:**

$$X = (\text{Refund} = \text{No}, \text{Divorced}, \text{Income} = 120K)$$

Tid	Refund	Marital Status	Taxable Income	Evade
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

### Using Bayes Theorem:

- $P(\text{Yes} | X) = \frac{P(X | \text{Yes})P(\text{Yes})}{P(X)}$
- $P(\text{No} | X) = \frac{P(X | \text{No})P(\text{No})}{P(X)}$
- How to estimate  $P(X | \text{Yes})$  and  $P(X | \text{No})$ ?

## Naïve Bayes Classifier

- Assume independence among attributes  $X_i$  when class is given:
  - $P(X_1, X_2, \dots, X_d | Y_j) = P(X_1 | Y_j) P(X_2 | Y_j) \dots P(X_d | Y_j)$
  - Now we can estimate  $P(X_i | Y_j)$  for all  $X_i$  and  $Y_j$  combinations from the training data
  - New point is classified to  $Y_j$  if  $P(Y_j) \prod P(X_i | Y_j)$  is maximal.

## Conditional Independence

- $X$  and  $Y$  are conditionally independent given  $Z$  if  $P(X | YZ) = P(X | Z)$
- Example: Arm length and reading skills
  - Young child has shorter arm length and limited reading skills, compared to adults
  - If age is fixed, no apparent relationship between arm length and reading skills
  - Arm length and reading skills are conditionally independent given age

## Naïve Bayes on Example Data

**Given a Test Record:**

$$X = (\text{Refund} = \text{No}, \text{Divorced}, \text{Income} = 120K)$$

Tid	Refund	Marital Status	Taxable Income	Evade
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

- $P(X | \text{Yes}) =$ 
  - $P(\text{Refund} = \text{No} | \text{Yes}) \times$
  - $P(\text{Divorced} | \text{Yes}) \times$
  - $P(\text{Income} = 120K | \text{Yes})$
- $P(X | \text{No}) =$ 
  - $P(\text{Refund} = \text{No} | \text{No}) \times$
  - $P(\text{Divorced} | \text{No}) \times$
  - $P(\text{Income} = 120K | \text{No})$

## Estimate Probabilities from Data

Tid	Refund	Marital Status	Taxable Income	Evaade
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

- Class:  $P(Y) = N_c/N$ 
  - e.g.,  $P(\text{No}) = 7/10$ ,  $P(\text{Yes}) = 3/10$
- For categorical attributes:
$$P(X_i | Y_k) = |X_{ik}| / N_{c_k}$$
  - where  $|X_{ik}|$  is number of instances having attribute value  $X_i$  and belonging to class  $Y_k$
  - Examples:
$$P(\text{Status=Married} | \text{No}) = 4/7$$

$$P(\text{Refund}=Yes | \text{Yes})=0$$

## Estimate Probabilities from Data

- For continuous attributes:
  - Discretization:** Partition the range into K bins:
    - Replace continuous value with bin value
    - Attribute changed from continuous to ordinal
  - Probability density estimation:**
    - Assume attribute follows a normal distribution
    - Use data to estimate parameters of distribution (e.g., mean and standard deviation)
    - Once probability distribution is known, use it to estimate the conditional probability  $P(X_i | Y)$

## Estimate Probabilities from Data

Tid	Refund	Marital Status	Taxable Income	Evaade
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

- Normal distribution:
$$P(X_i | Y_j) = \frac{1}{\sqrt{2\pi\sigma_{ij}^2}} e^{-\frac{(X_i - \mu_j)^2}{2\sigma_{ij}^2}}$$
  - One for each  $(X_i, Y_j)$  pair
- For (Income, Class=No):
  - If Class=No
    - sample mean = 110
    - sample variance = 2975

$$P(\text{Income} = 120 | \text{No}) = \frac{1}{\sqrt{2\pi(2975)}} e^{-\frac{(120-110)^2}{2(2975)}} = 0.0072$$

## Example of Naïve Bayes Classifier

Given a Test Record:

$$X = (\text{Refund} = \text{No}, \text{Divorced}, \text{Income} = 120K)$$

Naïve Bayes Classifier:

$$\begin{aligned} P(\text{Refund} = \text{Yes} | \text{No}) &= 3/7 \\ P(\text{Refund} = \text{No} | \text{No}) &= 4/7 \\ P(\text{Refund} = \text{Yes} | \text{Yes}) &= 0 \\ P(\text{Refund} = \text{No} | \text{Yes}) &= 1 \\ P(\text{Marital Status} = \text{Single} | \text{No}) &= 2/7 \\ P(\text{Marital Status} = \text{Divorced} | \text{No}) &= 1/7 \\ P(\text{Marital Status} = \text{Married} | \text{No}) &= 4/7 \\ P(\text{Marital Status} = \text{Single} | \text{Yes}) &= 2/3 \\ P(\text{Marital Status} = \text{Divorced} | \text{Yes}) &= 1/3 \\ P(\text{Marital Status} = \text{Married} | \text{Yes}) &= 0 \end{aligned}$$

For Taxable Income:

$$\begin{aligned} \text{If class = No: sample mean} &= 110 \\ \text{sample variance} &= 2975 \end{aligned}$$

$$\begin{aligned} \text{If class = Yes: sample mean} &= 90 \\ \text{sample variance} &= 25 \end{aligned}$$

$$\text{Since } P(X|\text{No})P(\text{No}) > P(X|\text{Yes})P(\text{Yes})$$

$$\text{Therefore } P(\text{No}|X) > P(\text{Yes}|X)$$

$$\Rightarrow \text{Class} = \text{No}$$

$$\bullet P(X | \text{No}) = P(\text{Refund}=\text{No} | \text{No}) \times P(\text{Divorced} | \text{No}) \times P(\text{Income}=120K | \text{No}) = 4/7 \times 1/7 \times 0.0072 = 0.0006$$

$$\bullet P(X | \text{Yes}) = P(\text{Refund}=\text{No} | \text{Yes}) \times P(\text{Divorced} | \text{Yes}) \times P(\text{Income}=120K | \text{Yes}) = 1 \times 1/3 \times 1.2 \times 10^{-9} = 4 \times 10^{-10}$$

## Example of Naïve Bayes Classifier

### Given a Test Record:

$X = (\text{Refund} = \text{No}, \text{Divorced}, \text{Income} = 120\text{K})$

#### Naïve Bayes Classifier:

$P(\text{Refund} = \text{Yes} | \text{No}) = 3/7$   
 $P(\text{Refund} = \text{No} | \text{No}) = 4/7$   
 $P(\text{Refund} = \text{Yes} | \text{Yes}) = 0$   
 $P(\text{Refund} = \text{No} | \text{Yes}) = 1$   
 $P(\text{Marital Status} = \text{Single} | \text{No}) = 2/7$   
 $P(\text{Marital Status} = \text{Divorced} | \text{No}) = 1/7$   
 $P(\text{Marital Status} = \text{Married} | \text{No}) = 4/7$   
 $P(\text{Marital Status} = \text{Single} | \text{Yes}) = 2/3$   
 $P(\text{Marital Status} = \text{Divorced} | \text{Yes}) = 1/3$   
 $P(\text{Marital Status} = \text{Married} | \text{Yes}) = 0$

For Taxable Income:  
If class = No: sample mean = 110  
sample variance = 2975  
If class = Yes: sample mean = 90  
sample variance = 25

- $P(\text{Yes}) = 3/10$   
 $P(\text{No}) = 7/10$
- $P(\text{Yes} | \text{Divorced}) = 1/3 \times 3/10 / P(\text{Divorced})$   
 $P(\text{No} | \text{Divorced}) = 1/7 \times 7/10 / P(\text{Divorced})$
- $P(\text{Yes} | \text{Refund} = \text{No, Divorced}) = 1 \times 1/3 \times 3/10 / P(\text{Divorced, Refund} = \text{No})$   
 $P(\text{No} | \text{Refund} = \text{No, Divorced}) = 4/7 \times 1/7 \times 7/10 / P(\text{Divorced, Refund} = \text{No})$

## Issues with Naïve Bayes Classifier

#### Naïve Bayes Classifier:

$P(\text{Refund} = \text{Yes} | \text{No}) = 3/7$   
 $P(\text{Refund} = \text{No} | \text{No}) = 4/7$   
 $P(\text{Refund} = \text{Yes} | \text{Yes}) = 0$   
 $P(\text{Refund} = \text{No} | \text{Yes}) = 1$   
 $P(\text{Marital Status} = \text{Single} | \text{No}) = 2/7$   
 $P(\text{Marital Status} = \text{Divorced} | \text{No}) = 1/7$   
 $P(\text{Marital Status} = \text{Married} | \text{No}) = 4/7$   
 $P(\text{Marital Status} = \text{Single} | \text{Yes}) = 2/3$   
 $P(\text{Marital Status} = \text{Divorced} | \text{Yes}) = 1/3$   
 $P(\text{Marital Status} = \text{Married} | \text{Yes}) = 0$

For Taxable Income:  
If class = No: sample mean = 110  
sample variance = 2975  
If class = Yes: sample mean = 90  
sample variance = 25

- $P(\text{Yes}) = 3/10$   
 $P(\text{No}) = 7/10$
- $P(\text{Yes} | \text{Married}) = 0 \times 3/10 / P(\text{Married})$   
 $P(\text{No} | \text{Married}) = 4/7 \times 7/10 / P(\text{Married})$

## Issues with Naïve Bayes Classifier

Consider the table with Tid = 7 deleted

Tid	Refund	Marital Status	Taxable Income	Evaade
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

#### Naïve Bayes Classifier:

$P(\text{Refund} = \text{Yes} | \text{No}) = 2/6$   
 $P(\text{Refund} = \text{No} | \text{No}) = 4/6$   
 $P(\text{Refund} = \text{Yes} | \text{Yes}) = 0$   
 $P(\text{Refund} = \text{No} | \text{Yes}) = 1$   
 $P(\text{Marital Status} = \text{Single} | \text{No}) = 2/6$   
 $P(\text{Marital Status} = \text{Divorced} | \text{No}) = 0$   
 $P(\text{Marital Status} = \text{Married} | \text{No}) = 4/6$   
 $P(\text{Marital Status} = \text{Single} | \text{Yes}) = 2/3$   
 $P(\text{Marital Status} = \text{Divorced} | \text{Yes}) = 1/3$   
 $P(\text{Marital Status} = \text{Married} | \text{Yes}) = 0/3$

For Taxable Income:

If class = No: sample mean = 91  
sample variance = 685

If class = Yes: sample mean = 90  
sample variance = 25

Given X = (Refund = Yes, Divorced, 120K)

Naïve Bayes will not be able to classify

$P(X | \text{No}) = 2/6 \times 0 \times 0.0083 = 0$

X as Yes or No!

$P(X | \text{Yes}) = 0 \times 1/3 \times 1.2 \times 10^{-9} = 0$

## Issues with Naïve Bayes Classifier

• If one of the conditional probabilities is zero, then the entire expression becomes zero

• Need to use other estimates of conditional probabilities than simple fractions

• Probability estimation:

c: number of classes

p: prior probability of the class

m: parameter

$$\text{Original : } P(A_i | C) = \frac{N_{ic}}{N_c}$$

$$\text{Laplace : } P(A_i | C) = \frac{N_{ic} + 1}{N_c + c}$$

$$\text{m - estimate : } P(A_i | C) = \frac{N_{ic} + mp}{N_c + m}$$

$N_{ic}$ : number of instances having attribute value  $A_i$  in class  $c$

$N_c$ : number of instances in the class

**Example of Naïve Bayes Classifier**

NC STATE UNIVERSITY

Name	Give Birth	Can Fly	Live in Water	Have Legs	Class
human	yes	no	no	yes	mammals
python	no	no	no	no	non-mammals
salmon	no	no	yes	no	non-mammals
penguin	yes	no	yes	no	mammals
frog	no	no	sometimes	yes	non-mammals
komodo	no	no	no	yes	non-mammals
bat	yes	yes	no	yes	mammals
pigeon	no	yes	no	yes	non-mammals
cat	yes	no	no	yes	mammals
leopard shark	yes	no	yes	no	non-mammals
rat	no	no	sometimes	yes	non-mammals
penguin	no	no	sometimes	yes	non-mammals
porcupine	yes	no	no	yes	mammals
eel	no	no	yes	no	non-mammals
salamander	no	no	sometimes	yes	non-mammals
gila monster	no	no	no	yes	non-mammals
platypus	no	no	no	yes	mammals
owl	no	yes	no	yes	non-mammals
zebra	yes	no	yes	no	mammals
seagull	no	yes	no	yes	non-mammals

A: attributes  
M: mammals  
N: non-mammals

$$P(A|M) = \frac{6}{7} \times \frac{6}{7} \times \frac{2}{7} \times \frac{2}{7} = 0.06$$

$$P(A|N) = \frac{1}{13} \times \frac{10}{13} \times \frac{3}{13} \times \frac{4}{13} = 0.0042$$

$$P(A|M)P(M) = 0.06 \times \frac{7}{20} = 0.021$$

$$P(A|N)P(N) = 0.004 \times \frac{13}{20} = 0.0027$$

$$P(A|M)P(M) > P(A|N)P(N)$$

$$\Rightarrow Mammals$$

**Naïve Bayes (Summary)**

NC STATE UNIVERSITY

- Robust to isolated noise points
- Handle missing values by ignoring the instance during probability estimate calculations
- Robust to irrelevant attributes
- Independence assumption may not hold for some attributes
  - Use other techniques such as Bayesian Belief Networks (BBN)

**Naïve Bayes**

NC STATE UNIVERSITY

- How does Naïve Bayes perform on the following dataset?

Conditional independence of attributes is violated

**Naïve Bayes**

NC STATE UNIVERSITY

- How does Naïve Bayes perform on the following dataset?

Naïve Bayes can construct oblique decision boundaries

## Maximum Likelihood Classification

- Uses Bayes decision rule:

$$x \in c_i \text{ iff } p(c_i | x) > p(c_j | x) \forall j \neq i$$

How to compute  $p(c_i | x)$ ?

Bayes' Theorem:

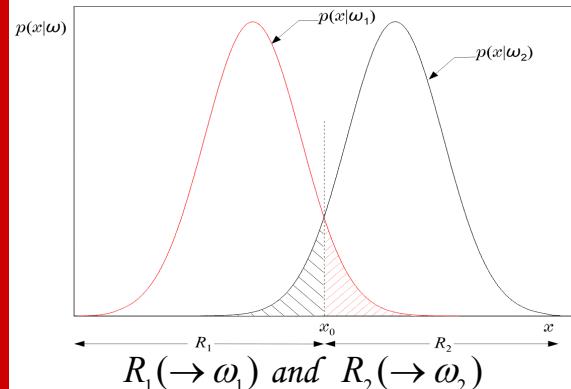
$$p(c_i | x) = \frac{p(x | c_i) \cdot p(c_i)}{p(x)}$$

Where  $p(x) = \text{prob. of finding a pixel from any class at location } x$

$$= \sum_{i=1}^M p(x | c_i) \cdot p(c_i)$$

**Note:**  $\omega, \gamma, c$  are used interchangeably for "class"

## Decision Rule



## Error and Optimality

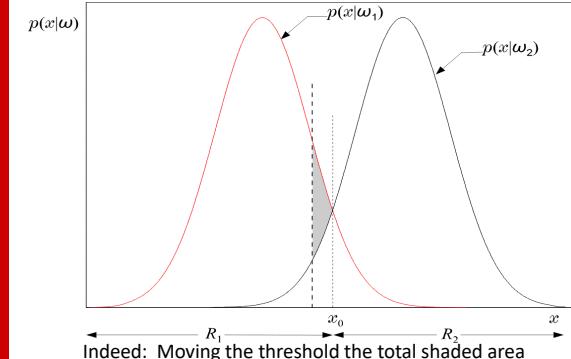
- Probability of error

– Total shaded area

$$P_e = \int_{-\infty}^{x_0} p(x | \omega_2) dx + \int_{x_0}^{+\infty} p(x | \omega_1) dx$$

- Bayesian classifier is **optimal** with respect to minimizing the classification error probability

## Optimality



## MLC

$p(x|\omega_i)$ : is the class conditional distribution and can be estimated from the training data (e.g., Gaussian)  
 $p(\omega_i)$ : prior probability

Modeling both of these terms is very important for accurate classification  
 Decision Rule:

$x \in \omega_i$  if  $p(x|\omega_i)p(\omega_i) > p(x|\omega_j)p(\omega_j) \forall i \neq j$

Simplified:

$$g_i(x) = \ln p(x|\omega_i) + \ln p(\omega_i)$$

$$x \in \omega_i \text{ if } g_i(x) > g_j(x) \forall i \neq j$$

## MLC

- Model Assumption

- Training samples were generated by multivariate Gaussian

$$p(x|\omega_i) = \frac{1}{\sqrt{(2\pi)^d |\Sigma_i|}} \exp\left[-\frac{1}{2}(x - \mu_i)^T \Sigma_i^{-1} (x - \mu_i)\right]$$

$$\text{MLE} \quad \hat{\mu}_i = \frac{1}{n} \sum_{k=1}^n x_k, \quad \hat{\Sigma}_i = \frac{1}{n} \sum_{k=1}^n (x_k - \hat{\mu}_i)(x_k - \hat{\mu}_i)^T$$

- For each feature vector  $x$   $p(\omega_i|x) = \frac{p(x|\omega_i)p(\omega_i)}{p(x)}$

- Apply, ML Decision

$$x \in \omega_i \text{ if } p(x|\omega_i)p(\omega_i) > p(x|\omega_j)p(\omega_j) \forall i \neq j$$

## Discriminant function

- $\ln(\cdot)$  is monotonic. Define:

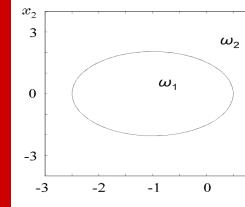
$$g_i(x) = \ln(p(x|\omega_i)P(\omega_i)) = \ln p(x|\omega_i) + \ln P(\omega_i)$$

$$g_i(x) = -\frac{1}{2}(x - \mu_i)^T \Sigma_i^{-1} (x - \mu_i) + \ln P(\omega_i) + -\frac{d}{2} \ln 2\pi - \frac{1}{2} \ln |\Sigma_i|$$

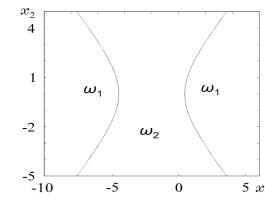
Note:  $\omega, y, c$  are used interchangeably for "class"

## Discriminant function

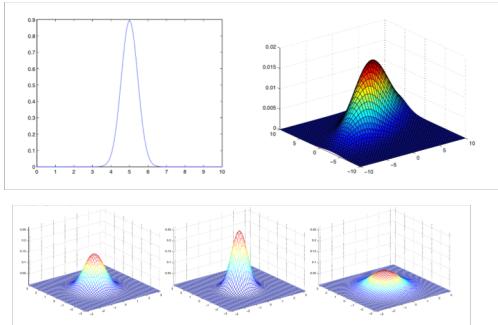
- Note that  $g_i(x)$  is quadratic, and produces quadrics, ellipsoids, parabolas, hyperbolas, and pairs of lines as decision boundaries



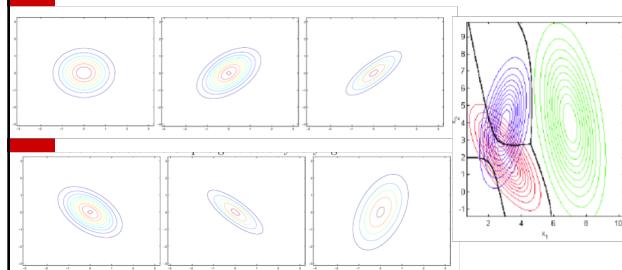
(a) Accurate Assessment of Covariance is important, as it influences decision boundaries



## Covariance Matrix



## Covariance Matrix



Accurate Assessment of Covariance is important,  
as it influences decision boundaries

Any Questions?