



TEXAS A&M UNIVERSITY
Engineering

ECEN 758 Data Mining and Analysis: Lecture 12, Bayesian and Nearest Neighbor Classification

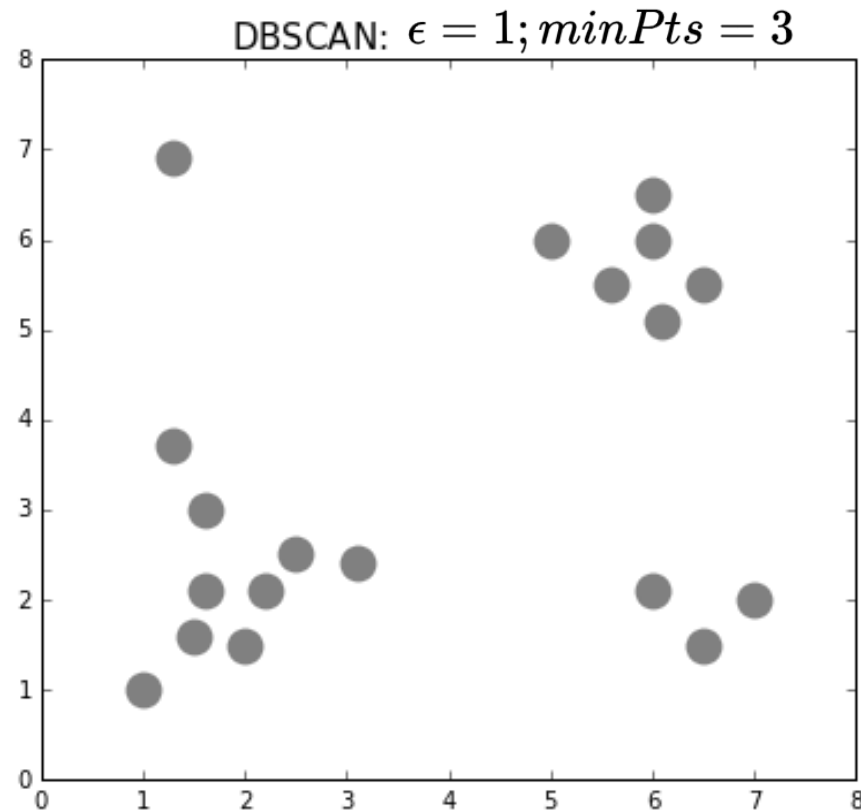
Joshua Peeples, Ph.D.

Assistant Professor

Department of Electrical and Computer Engineering

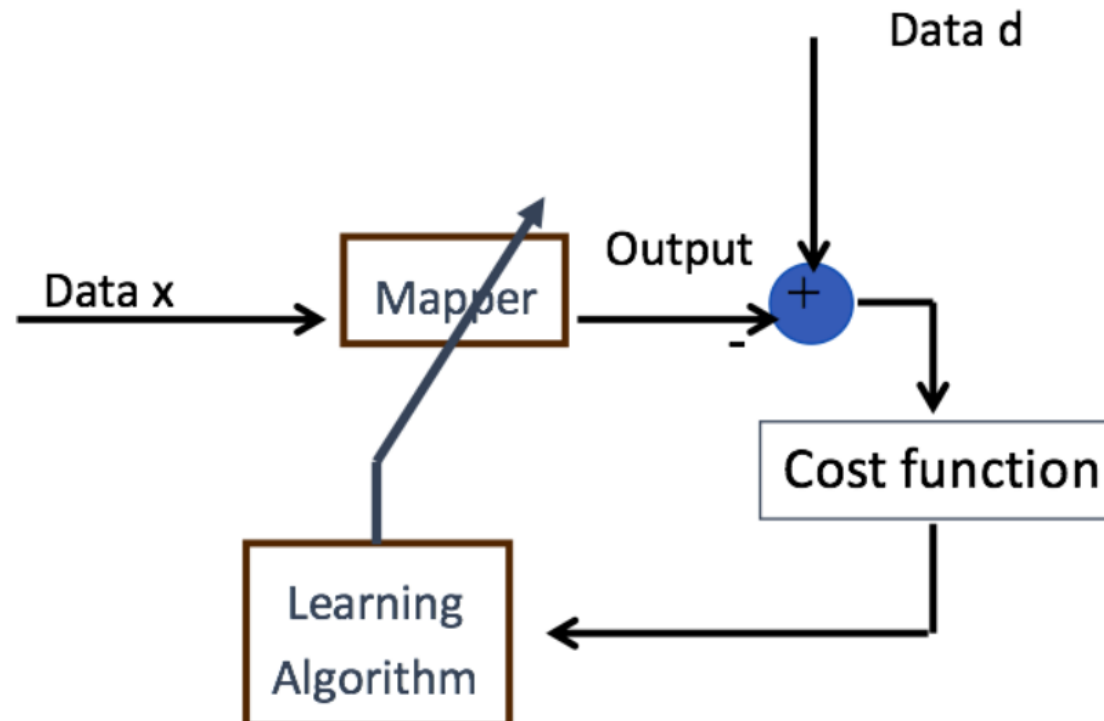
- Assignment #3 will be released next Wednesday (10/09)
- No class next Monday (10/07): Fall Break
- Guest lecture next Wednesday (10/09)
 - +1 on Midterm exam for attendance (must sign-in to attendance sheet)
 - +3 for one paragraph (5 – 7 sentences) summary of presentation (must attend lecture)
 - Section 700: Watch recording and submit 1 paragraph (+1) or 2 paragraphs (+3)
- Exam I in two weeks on Monday, 10/14

- Density-based Clustering



- Introduction to Classification I
 - Definition
 - Model types
 - Hyperparameters vs Parameters
- Bayesian and Nearest Neighbor Classification
- Reading: ZM Chapter 18

- In machine learning the model is derived from the data (observations)
- As a learning machine, the model can be modified over time, with additional data (observations), with the goal of improving outcomes



Many Sub-areas in Machine Learning



TEXAS A&M UNIVERSITY
Engineering

- Supervised Learning
- Unsupervised Learning
- Reinforcement Learning
- Semi-supervised Learning
- Self-supervised Learning
- Multiple Instance Learning
- Active Learning
- Transfer Learning
-

- Supervised learning
 - We “coach” the computer
 - Uses knowledge already learned
- Unsupervised learning
 - “We’re free!!”



Supervised Learning: Classification

Supervised Learning



TEXAS A&M UNIVERSITY
Engineering

Learning from experience

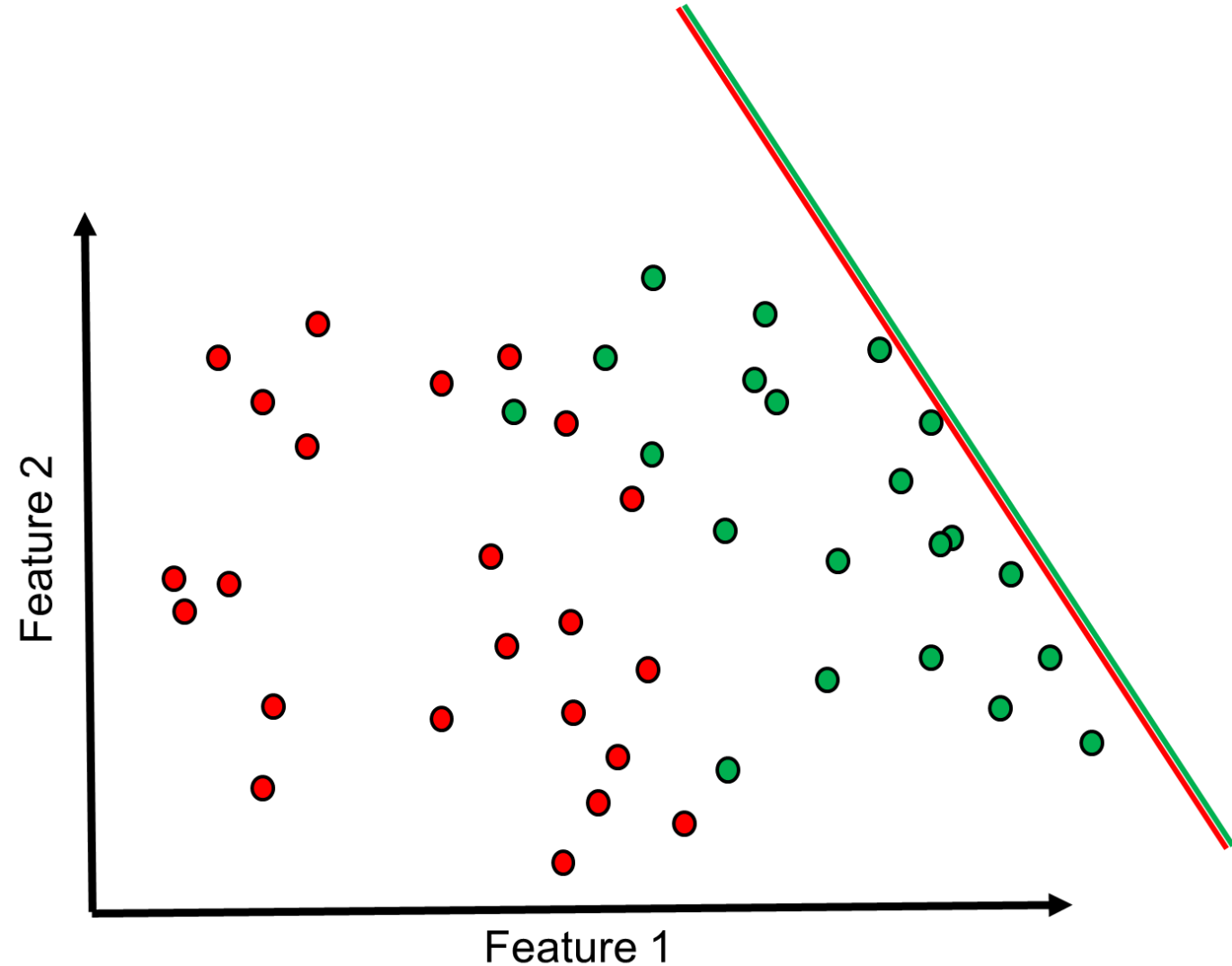


0: Macaw 1: Conure

Classification Example



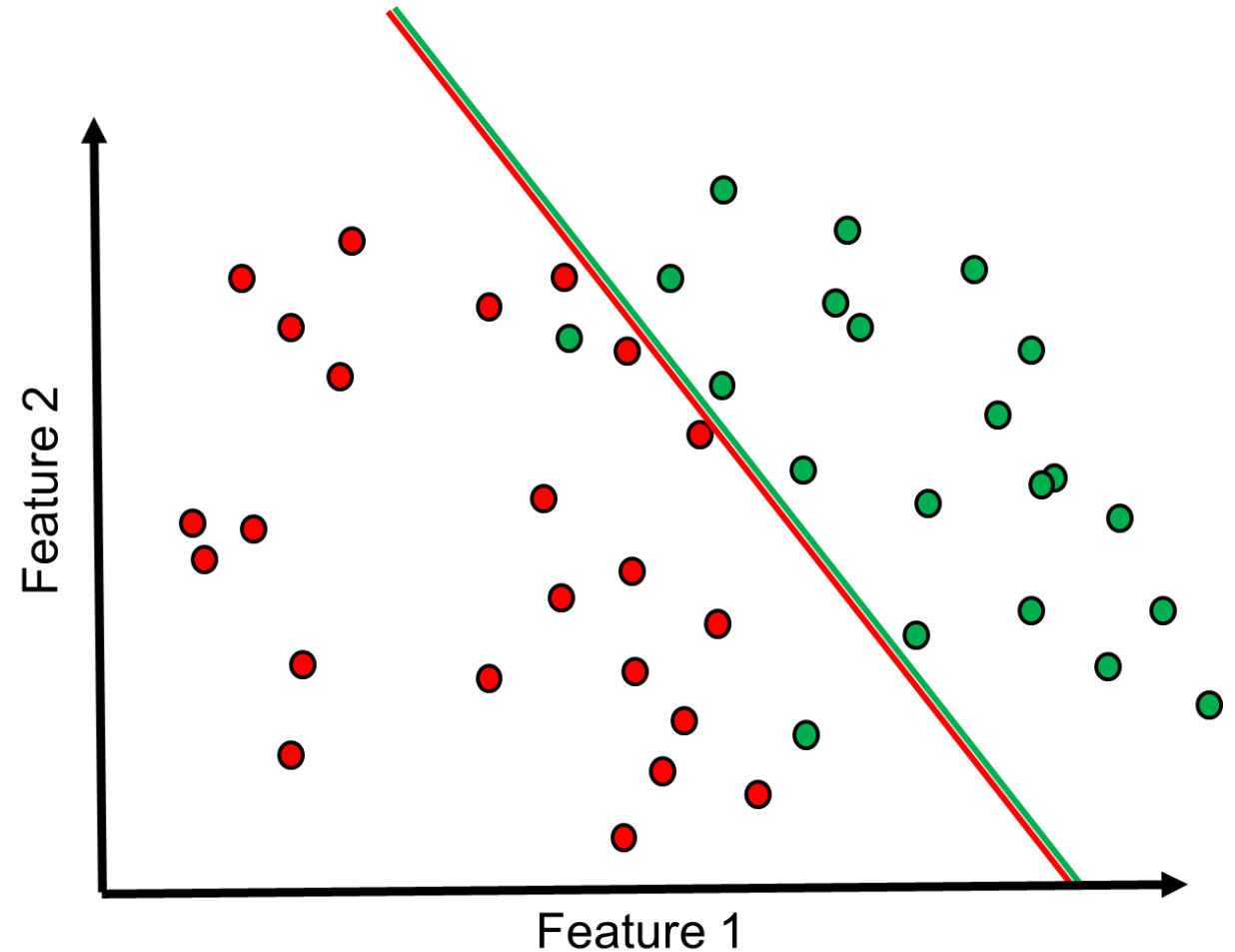
- Given a set of labeled (training) instances, learn a model
 - E.g., Find parameters for linear classifier
- Accurately predict new (test) samples



Classification Example



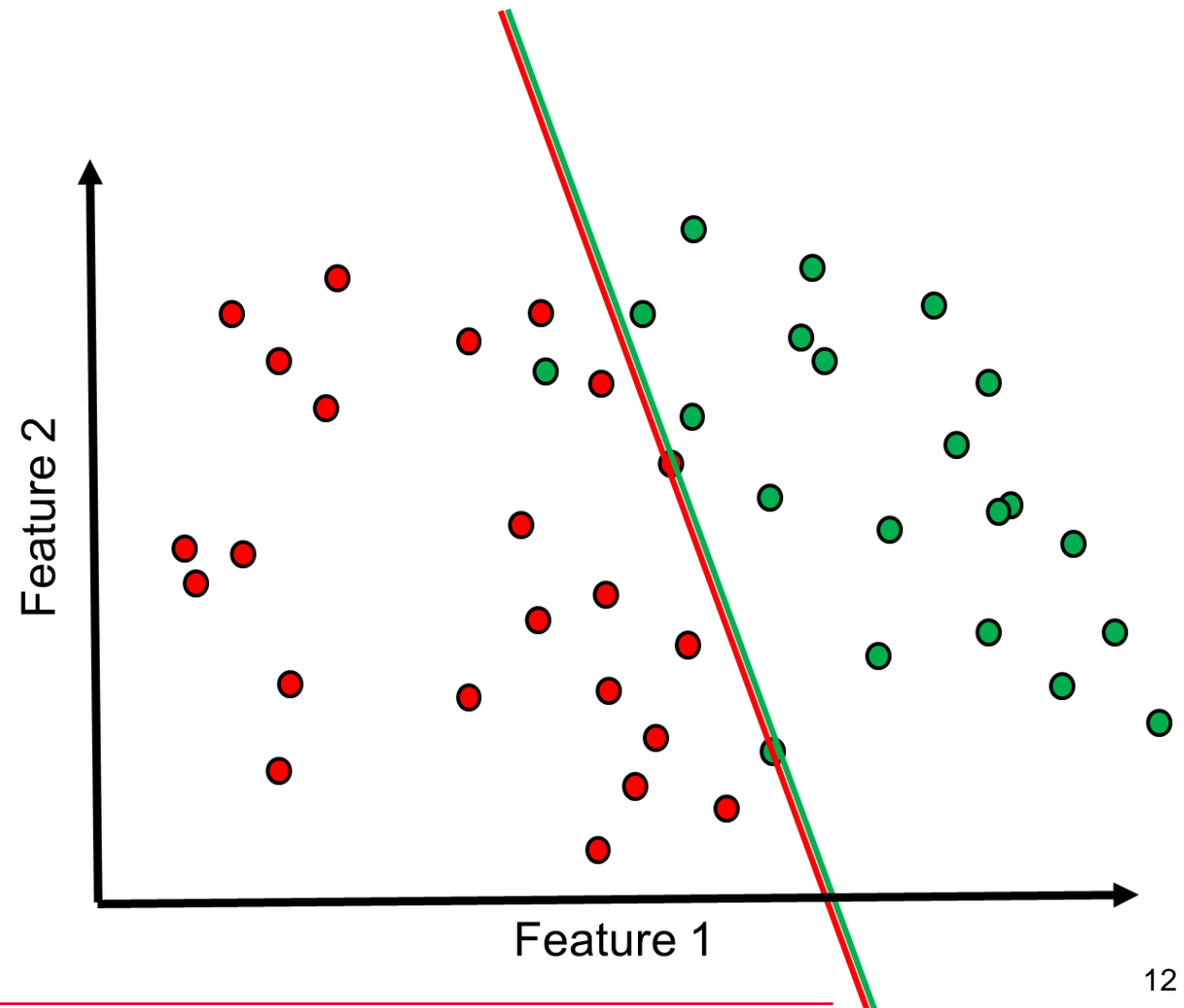
- Given a set of labeled (training) instances, learn a model
 - E.g., Find parameters for linear classifier
- Accurately predict new (test) samples



Classification Example



- Given a set of labeled (training) instances, learn a model
 - E.g., Find parameters for linear classifier
- Accurately predict new (test) samples



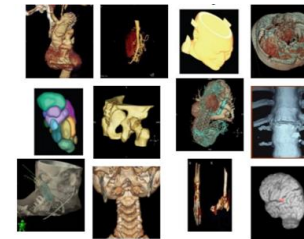


Classification Overview

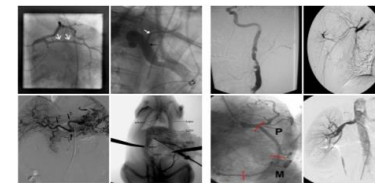
Classification Tasks



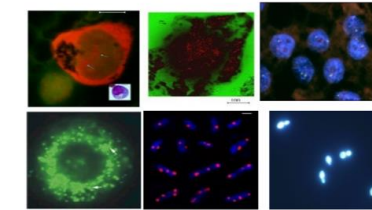
- Classification: given features X , predict label (class) y
- Examples:
 - Communication symbol recognition (input signal w/source & channel induced noise, classes: valid communication symbols)
 - Medical diagnosis (input: symptoms, classes: diseases)
 - Character recognition (input: handwritten characters, classes: {a..z})
 - Fraud detection (input: account activity, classes: fraud / no fraud)
 - ... many more



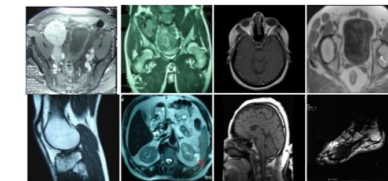
3D Reconstruction



Angiography



Fluorescence microscopy



Magnetic Resonance

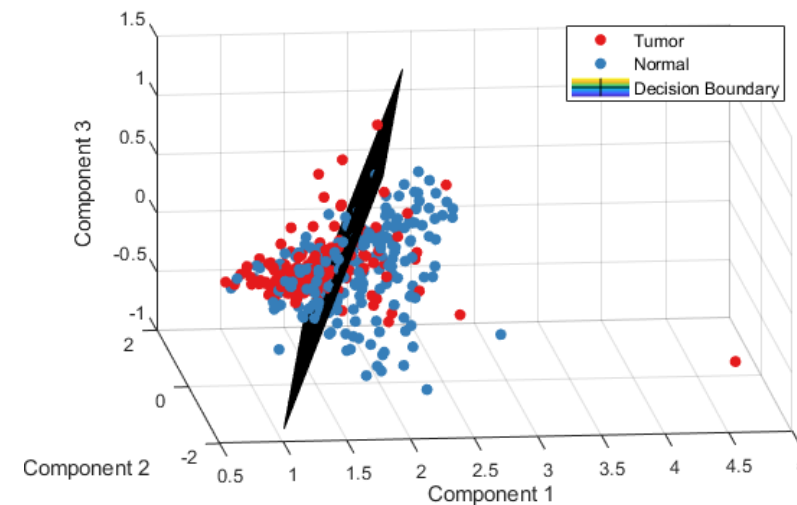
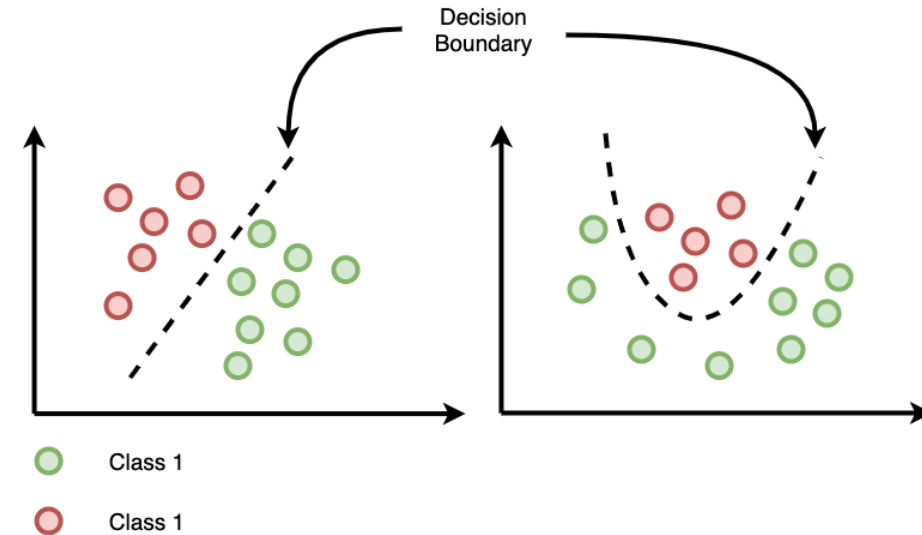
Src: Nowka 2015, IBM CLEF 2013



Image from: Madhu Ramiah- Medium

Decision Boundaries

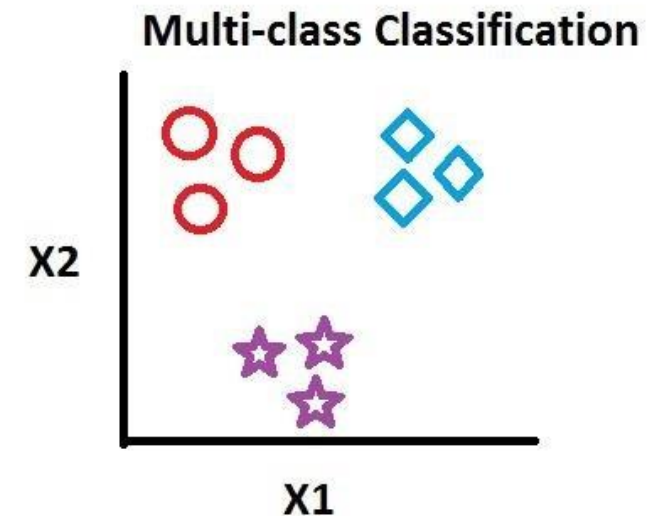
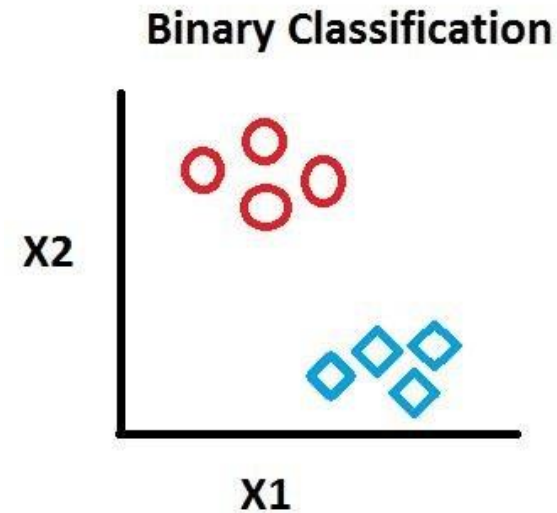
- Distinguish between classes
- In higher dimensions (>2), learning hyperplanes



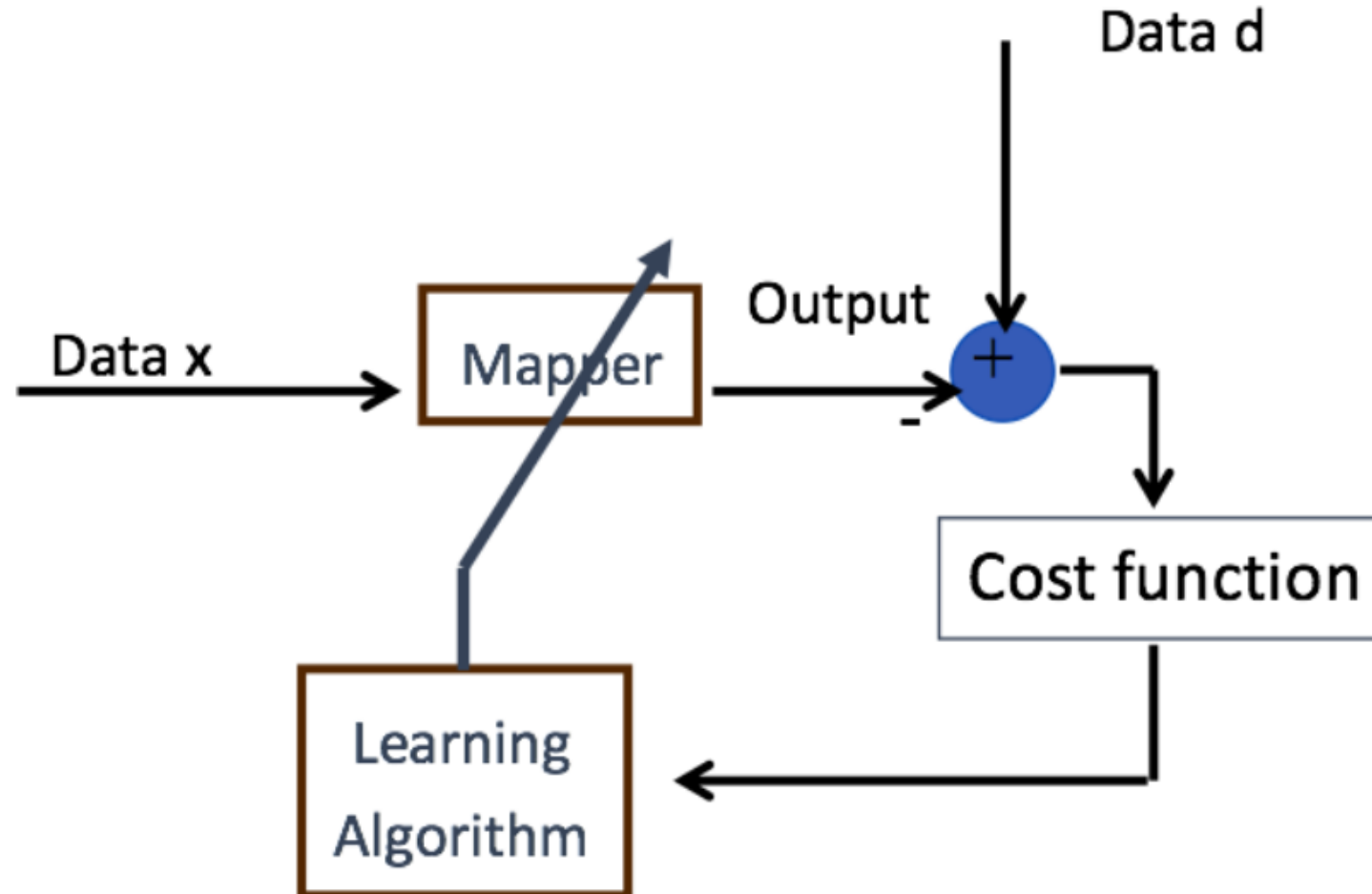
Types of Classification



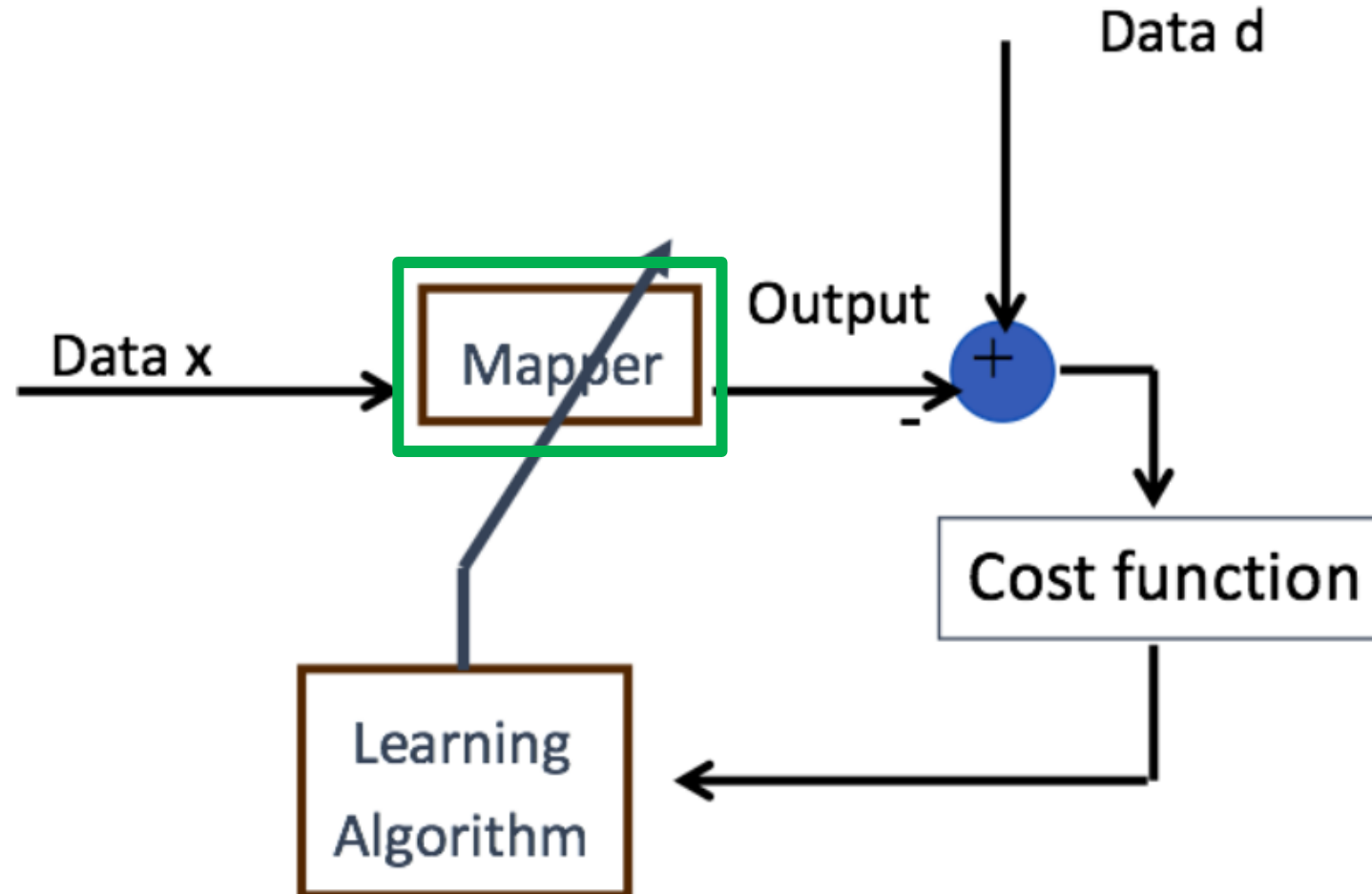
- **Binary**
- **Multi-class**
- Multi-label (binary)
- Multi-class, multi-output



Model Types



Model Types

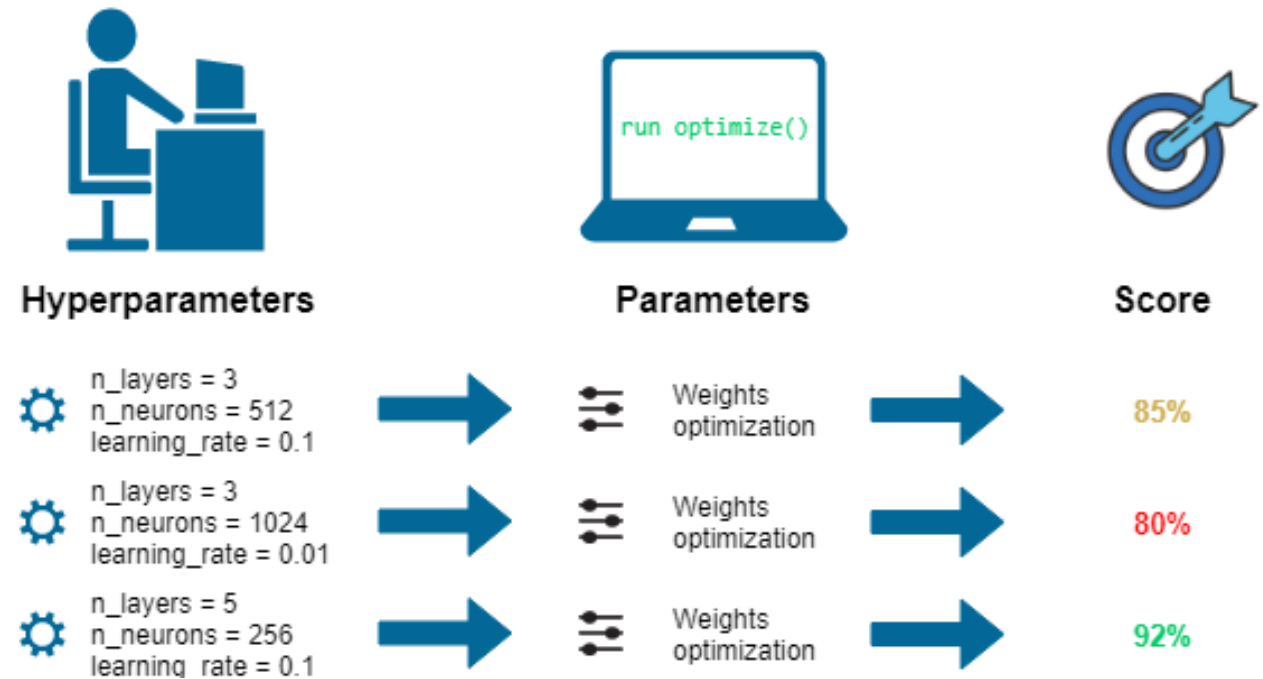


Parameters and Hyperparameters



TEXAS A&M UNIVERSITY
Engineering

- Parameters (θ) – (e.g. weights, means, covariances) are derived/modified (learned!) in the process of training the model
- Hyperparameters – (e.g. k in K-means) are set outside of the training/fit cycle and are used to direct characteristics of the resulting model

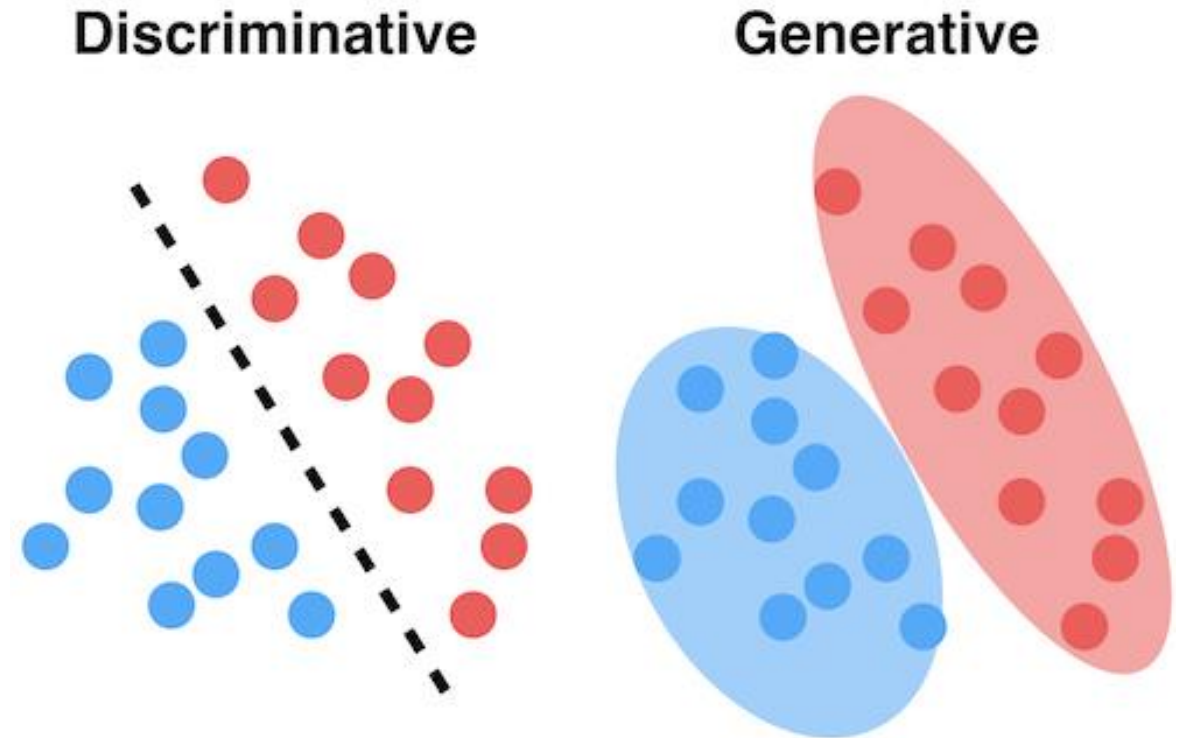


Discriminative vs Generative



TEXAS A&M UNIVERSITY
Engineering

- **Discriminative**
 - Focus on the decision boundary between classes
- **Generative**
 - Focus on distribution of classes



Parametric vs Non-parametric



- Parametric
 - Use known functional form
 - Makes assumptions of data
- Non-parametric
 - Free to learn any functional form
 - No strong assumptions of mapping function

Parametric

- ✓ Fast
- ✓ Simple
- ✓ Less data

- ▶ Limited complexity
- ▶ Strong assumptions
- ▶ Poor fit (if assumptions are not correct)

Nonparametric

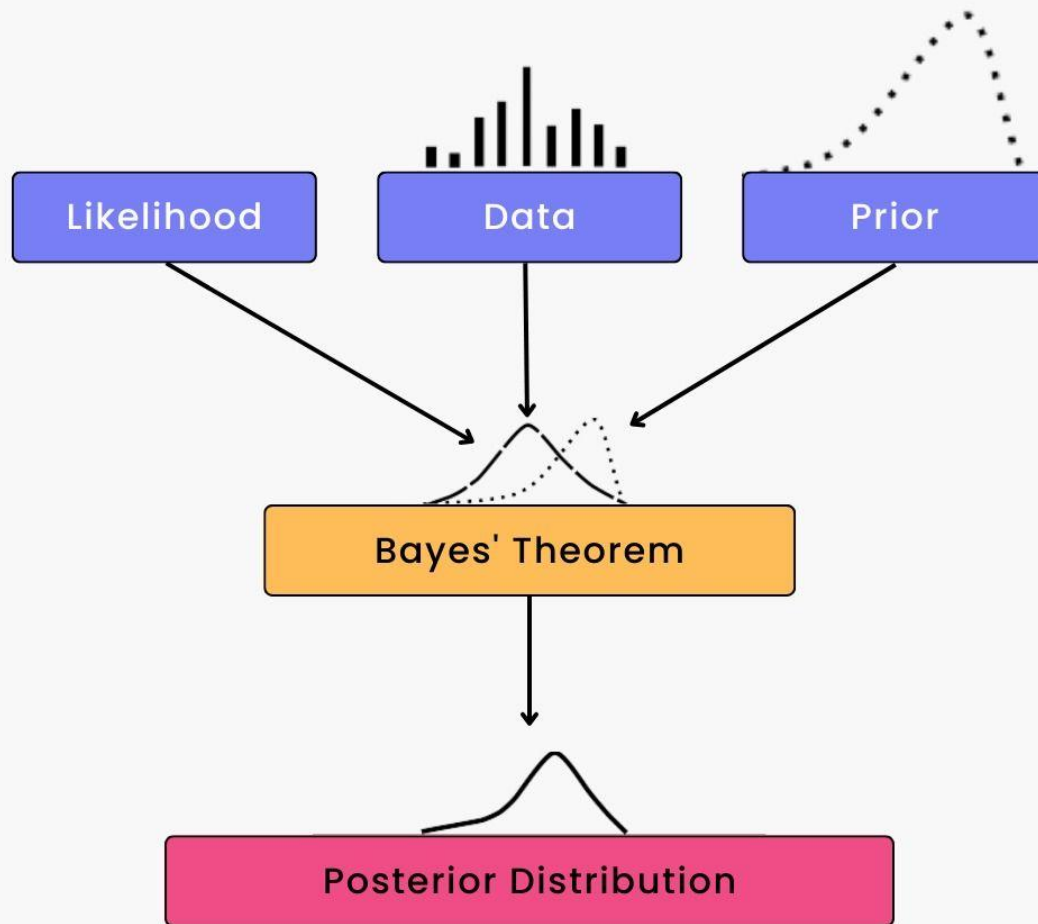
- ✓ Flexible
- ✓ Powerful
- ✓ Effective

- ▶ More data
- ▶ Computationally expensive
- ▶ Hard to interpret if models are too complex



Bayes Classifier

Bayes' Theorem



$$\begin{array}{c} \text{"Posterior"} \\ P(y|x) = \frac{\text{"Likelihood"} \quad \text{"Prior"} \\ P(x|y)}{\text{"Evidence"} \quad P(x)} P(y) \end{array}$$

- Given training dataset **D**
 - n points \mathbf{x}_i in d-dimensions
 - Each point has class label y_i
- Estimate posterior probability for each class c_i
- Predicted class has highest probability (\hat{y})

$$D = \begin{pmatrix} X_1 & X_2 & \cdots & X_d \\ x_{11} & x_{12} & \cdots & x_{1d} \\ x_{21} & x_{22} & \cdots & x_{2d} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \cdots & x_{nd} \end{pmatrix} y_i \in \{c_1, c_2, \dots, c_k\}$$

$$\hat{y} = \arg \max_{c_i} \{P(c_i | \mathbf{x})\}$$

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

- Given training dataset **D**
 - n points \mathbf{x}_i in d-dimensions
 - Each point has class label y_i
- $P(\mathbf{x})$ is fixed for given point

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

$$P(\mathbf{x}) = \sum_{j=1}^k P(\mathbf{x}|c_j) \cdot P(c_j)$$

$$\arg \max_{c_i} \{ P(\mathbf{x}|c_i) P(c_i) \}$$



Bayes Classifier: Parameter Estimation

- Can be given prior
 - E.g., (fair) coin toss
- Can find by counting observations
 - Number point data points in class (n_i) divided by total number of data points (n)

$$\mathbf{D}_i = \{\mathbf{x}_j \in \mathbf{D} \mid \mathbf{x}_j \text{ has class } y_j = c_i\}$$

$$\hat{P}(c_i) = \frac{n_i}{n}$$

Bayes Classifier: Likelihood



- Need to estimate joint probability of \mathbf{x} across all dimensions
- Assume each class is normally distributed

$$P(\mathbf{x}|c_i) = P(\mathbf{x} = (x_1, x_2, \dots, x_d) | c_i)$$

$$\hat{f}_i(\mathbf{x}) = \hat{f}(\mathbf{x} | \hat{\boldsymbol{\mu}}_i, \hat{\boldsymbol{\Sigma}}_i) = \frac{1}{(\sqrt{2\pi})^d \sqrt{|\hat{\boldsymbol{\Sigma}}_i|}} \exp \left\{ -\frac{(\mathbf{x} - \hat{\boldsymbol{\mu}}_i)^T \hat{\boldsymbol{\Sigma}}_i^{-1} (\mathbf{x} - \hat{\boldsymbol{\mu}}_i)}{2} \right\}$$

Bayes Classifier: Likelihood



- Find posterior probability, $P(c_i|\mathbf{x})$
- Predict class using highest posterior probability

$$P(c_i|\mathbf{x}) = \frac{\hat{f}_i(\mathbf{x})P(c_i)}{\sum_{j=1}^k \hat{f}_j(\mathbf{x})P(c_j)}$$

$$\hat{y} = \arg \max_{c_i} \left\{ \hat{f}_i(\mathbf{x})P(c_i) \right\}$$

$$\hat{f}_i(\mathbf{x}) = \hat{f}(\mathbf{x}|\hat{\boldsymbol{\mu}}_i, \hat{\boldsymbol{\Sigma}}_i) = \frac{1}{(\sqrt{2\pi})^d \sqrt{|\hat{\boldsymbol{\Sigma}}_i|}} \exp \left\{ -\frac{(\mathbf{x} - \hat{\boldsymbol{\mu}}_i)^T \hat{\boldsymbol{\Sigma}}_i^{-1} (\mathbf{x} - \hat{\boldsymbol{\mu}}_i)}{2} \right\}$$

Bayes Classifier Pseudocode



BayesClassifier ($D = \{(\mathbf{x}_j, y_j)\}_{j=1}^n$):

```
1 for  $i = 1, \dots, k$  do
2    $\mathbf{D}_i \leftarrow \{\mathbf{x}_j \mid y_j = c_i, j = 1, \dots, n\}$  // class-specific subsets
3    $n_i \leftarrow |\mathbf{D}_i|$  // cardinality
4    $\hat{P}(c_i) \leftarrow n_i/n$  // prior probability
5    $\hat{\boldsymbol{\mu}}_i \leftarrow \frac{1}{n_i} \sum_{\mathbf{x}_j \in \mathbf{D}_i} \mathbf{x}_j$  // mean
6    $\mathbf{Z}_i \leftarrow \mathbf{D}_i - \mathbf{1}_{n_i} \hat{\boldsymbol{\mu}}_i^T$  // centered data
7    $\hat{\boldsymbol{\Sigma}}_i \leftarrow \frac{1}{n_i} \mathbf{Z}_i^T \mathbf{Z}_i$  // covariance matrix
8 return  $\hat{P}(c_i), \hat{\boldsymbol{\mu}}_i, \hat{\boldsymbol{\Sigma}}_i$  for all  $i = 1, \dots, k$ 
```

Testing (\mathbf{x} and $\hat{P}(c_i), \hat{\boldsymbol{\mu}}_i, \hat{\boldsymbol{\Sigma}}_i$, for all $i \in [1, k]$):

```
9  $\hat{y} \leftarrow \arg \max_{c_i} \{f(\mathbf{x} | \hat{\boldsymbol{\mu}}_i, \hat{\boldsymbol{\Sigma}}_i) \cdot \hat{P}(c_i)\}$ 
10 return  $\hat{y}$ 
```

Bayes Classifier Example: Iris Data



- Class 1: Iris-setosa
- Class 2: Other classes
- Steps:
 - **Compute priors**
 - **Estimate parameters**
 - Compute posterior for new test points

$$\hat{P}(c_1) = \frac{n_1}{n} = \frac{50}{150} = 0.33$$

$$\hat{P}(c_2) = \frac{n_2}{n} = \frac{100}{150} = 0.67$$

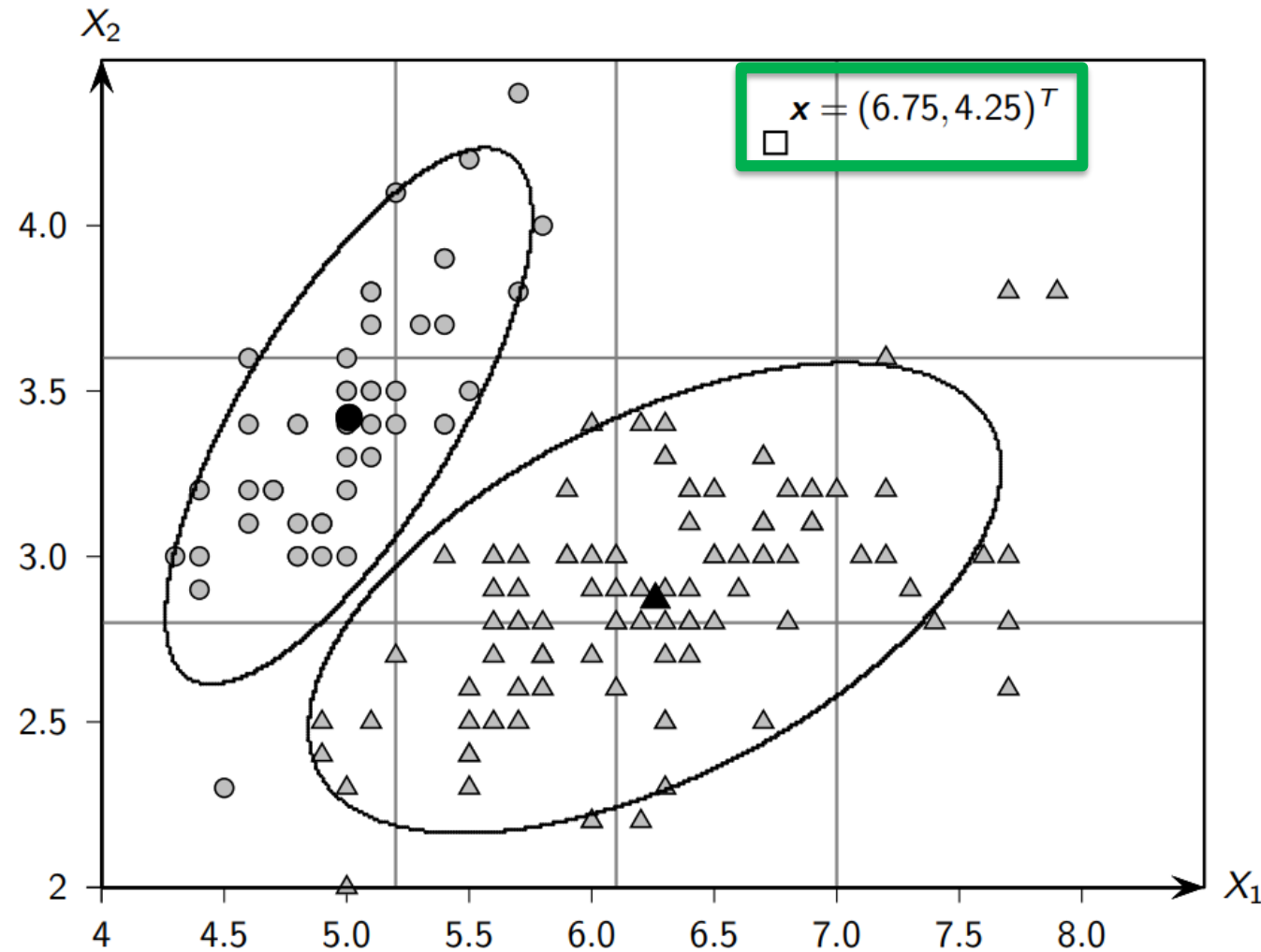
$$\hat{\mu}_1 = \begin{pmatrix} 5.006 \\ 3.418 \end{pmatrix}$$

$$\hat{\mu}_2 = \begin{pmatrix} 6.262 \\ 2.872 \end{pmatrix}$$

$$\hat{\Sigma}_1 = \begin{pmatrix} 0.1218 & 0.0983 \\ 0.0983 & 0.1423 \end{pmatrix}$$

$$\hat{\Sigma}_2 = \begin{pmatrix} 0.435 & 0.1209 \\ 0.1209 & 0.1096 \end{pmatrix}$$

Bayes Classifier Example: Iris Data



- Class 1: Iris-setosa
- Class 2: Other classes
- Steps:
 - Compute priors
 - Estimate parameters
 - **Compute posterior for new test points**

$$\hat{P}(c_2|\mathbf{x}) > \hat{P}(c_1|\mathbf{x})$$

$$\hat{P}(c_1|\mathbf{x}) \propto \hat{f}(\mathbf{x}|\hat{\mu}_1, \hat{\Sigma}_1) \hat{P}(c_1) = (4.951 \times 10^{-7}) \times 0.33 = 1.634 \times 10^{-7}$$

$$\hat{P}(c_2|\mathbf{x}) \propto \hat{f}(\mathbf{x}|\hat{\mu}_2, \hat{\Sigma}_2) \hat{P}(c_2) = (2.589 \times 10^{-5}) \times 0.67 = 1.735 \times 10^{-5}$$



Bayes Classifier: Categorical

Bayes Classifier: Categorical



- Compute joint probability mass function (PMF) from \mathbf{X}
- Compute joint PMF for each class by counting the number of times an attribute occurs

$$\text{dom}(X_j) = \{a_{j1}, a_{j2}, \dots, a_{jm_j}\}$$

$$P(\mathbf{x}|c_i) = f(\mathbf{v}|c_i) = f(\mathbf{X}_1 = \mathbf{e}_{1r_1}, \dots, \mathbf{X}_d = \mathbf{e}_{dr_d} | c_i)$$

$$\hat{f}(\mathbf{v}|c_i) = \frac{n_i(\mathbf{v})}{n_i}$$

- Avoid zero probabilities
- Use Laplace Smoothing
- Introduce smoothing parameter: α
- Larger values of α move data likelihood to uniform distribution
- m_j is domain for each attribute

$$\hat{f}(\mathbf{v}|c_i) = \frac{n_i(\mathbf{v}) + 1}{n_i + \prod_{j=1}^d m_j}$$

$$\hat{f}(\mathbf{v}_j|c_i) = \frac{n_i(\mathbf{v}) + \alpha}{n_i + \alpha * \prod_{j=1}^d m_j}$$

Bayes Classifier Example: Iris Data



Assume that the sepal length and sepal width attributes in the Iris dataset have been discretized as shown below.

Bins	Domain
[4.3, 5.2]	Very Short (a_{11})
(5.2, 6.1]	Short (a_{12})
(6.1, 7.0]	Long (a_{13})
(7.0, 7.9]	Very Long (a_{14})

(a) Discretized sepal length

Bins	Domain
[2.0, 2.8]	Short (a_{21})
(2.8, 3.6]	Medium (a_{22})
(3.6, 4.4]	Long (a_{23})

(b) Discretized sepal width

We have $|dom(X_1)| = m_1 = 4$ and $|dom(X_2)| = m_2 = 3$.

Bayes Classifier Example: Iris Data



	Class: c_1	X_2			\hat{f}_{X_1}
		Short (e_{21})	Medium (e_{22})	Long (e_{23})	
X_1	Very Short (e_{11})	1/50	33/50	5/50	39/50
	Short (e_{12})	0	3/50	8/50	13/50
	Long (e_{13})	0	0	0	0
	Very Long (e_{14})	0	0	0	0
\hat{f}_{X_2}		1/50	36/50	13/50	

	Class: c_2	X_2			\hat{f}_{X_1}
		Short (e_{21})	Medium (e_{22})	Long (e_{23})	
X_1	Very Short (e_{11})	6/100	0	0	6/100
	Short (e_{12})	24/100	15/100	0	39/100
	Long (e_{13})	13/100	30/100	0	43/100
	Very Long (e_{14})	3/100	7/100	2/100	12/100
\hat{f}_{X_2}		46/100	52/100	2/100	

Bayes Classifier Example: Iris Data



Consider a test point $\mathbf{x} = (5.3, 3.0)^T$ corresponding to the categorical point (Short, Medium), which is represented as $\mathbf{v} = (\mathbf{e}_{12}^T \quad \mathbf{e}_{22}^T)^T$.

The prior probabilities of the classes are $\hat{P}(c_1) = 0.33$ and $\hat{P}(c_2) = 0.67$.
The likelihood and posterior probability for each class is given as

$$\hat{P}(\mathbf{x}|c_1) = \hat{f}(\mathbf{v}|c_1) = 3/50 = 0.06$$

$$\hat{P}(\mathbf{x}|c_2) = \hat{f}(\mathbf{v}|c_2) = 15/100 = 0.15$$

$$\hat{P}(c_1|\mathbf{x}) \propto 0.06 \times 0.33 = 0.0198$$

$$\hat{P}(c_2|\mathbf{x}) \propto 0.15 \times 0.67 = 0.1005$$

In this case the predicted class is $\hat{y} = c_2$.

- May lack enough data to estimate joint pdf or pmf
 - Especially with many features
- Numeric attributes need to estimate covariances (d^2)
- Categorical attributes need to estimate all possible values
 - If attribute is binary, 2^d possibilities



Naive Bayes Classifier

- Assume attributes are independent

$$P(\mathbf{x}|c_i) = P(x_1, x_2, \dots, x_d|c_i) = \prod_{j=1}^d P(x_j|c_i)$$

- Likelihood for class c_i for dimension x_j

$$P(x_j|c_i) \propto f(x_j|\hat{\mu}_{ij}, \hat{\sigma}_{ij}^2) = \frac{1}{\sqrt{2\pi}\hat{\sigma}_{ij}} \exp \left\{ -\frac{(x_j - \hat{\mu}_{ij})^2}{2\hat{\sigma}_{ij}^2} \right\}$$

Naive Bayes Classifier



- Assumption leads to diagonal covariance
- Each class has mean vector and covariance matrix
 - 2d parameters to estimate per dimension x_j

$$\Sigma_i = \begin{pmatrix} \sigma_{i1}^2 & 0 & \dots & 0 \\ 0 & \sigma_{i2}^2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \sigma_{id}^2 \end{pmatrix}$$

Naive Bayes Classifier Pseudocode



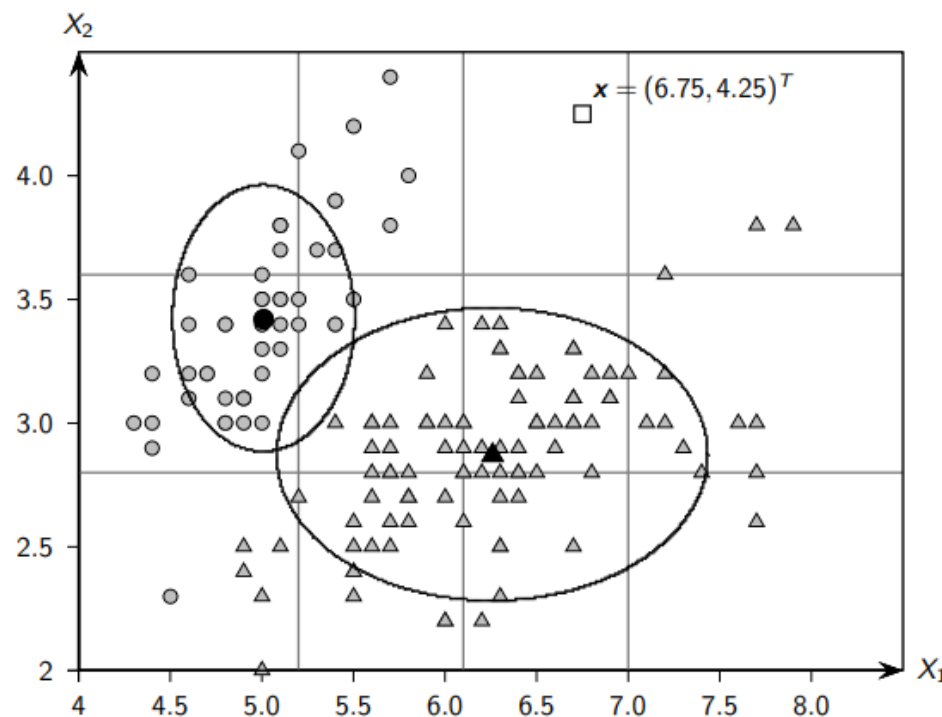
NaiveBayes ($D = \{(x_j, y_j)\}_{j=1}^n$):

```
1 for  $i = 1, \dots, k$  do
2    $D_i \leftarrow \{x_j^T \mid y_j = c_i, j = 1, \dots, n\}$  // class-specific subsets
3    $n_i \leftarrow |D_i|$  // cardinality
4    $\hat{P}(c_i) \leftarrow n_i/n$  // prior probability
5    $\hat{\mu}_i \leftarrow \frac{1}{n_i} \sum_{x_j \in D_i} x_j$  // mean
6    $\bar{D}_i = D_i - 1 \cdot \hat{\mu}_i^T$  // centered data for class  $c_i$ 
7   for  $j = 1, \dots, d$  do // class-specific var for  $j$ th attribute
8      $\hat{\sigma}_{ij}^2 \leftarrow \frac{1}{n_i} (\bar{X}_j^i)^T (\bar{X}_j^i)$  // variance
9    $\hat{\sigma}_i \leftarrow (\hat{\sigma}_{i1}^2, \dots, \hat{\sigma}_{id}^2)^T$  // class-specific attribute variances
10 return  $\hat{P}(c_i), \hat{\mu}_i, \hat{\sigma}_i$  for all  $i = 1, \dots, k$ 
```

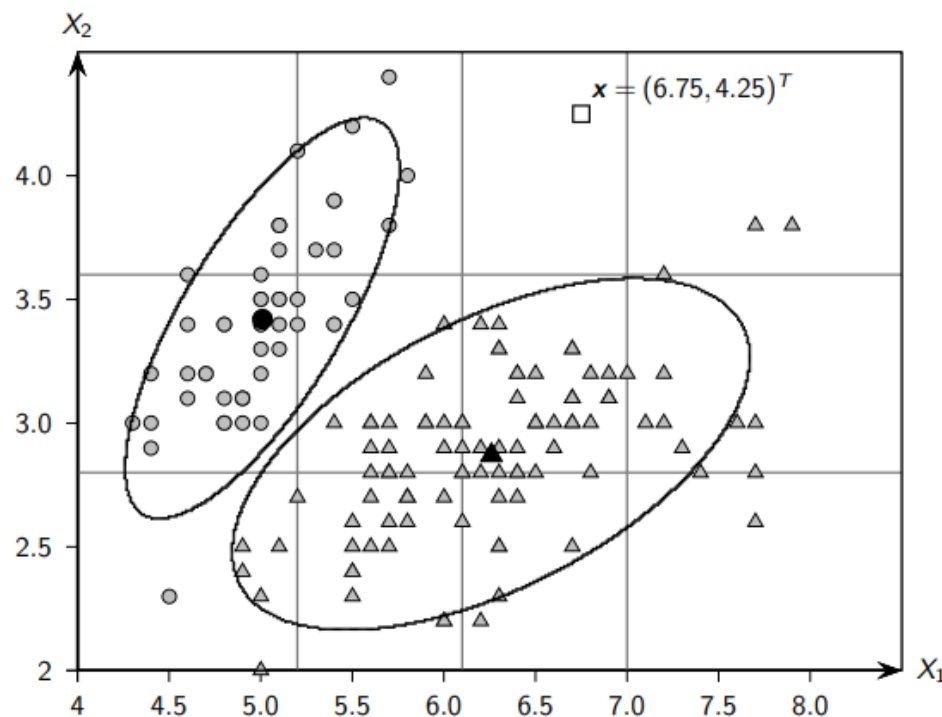
Testing (x and $\hat{P}(c_i), \hat{\mu}_i, \hat{\sigma}_i$, for all $i \in [1, k]$):

```
11  $\hat{y} \leftarrow \arg \max_{c_i} \left\{ \hat{P}(c_i) \prod_{j=1}^d f(x_j | \hat{\mu}_{ij}, \hat{\sigma}_{ij}^2) \right\}$ 
12 return  $\hat{y}$ 
```

Naive Bayes vs Full Bayes



(a) Naive Bayes



(b) Full Bayes



Naive Bayes Classifier: Categorical

Bayes Classifier: Categorical



- Compute joint probability mass function (PMF) from \mathbf{X}
- Compute joint PMF for each class by counting the number of times an attribute occurs

$$P(\mathbf{x}|c_i) = \prod_{j=1}^d P(x_j|c_i) = \prod_{j=1}^d f(\mathbf{x}_j = \mathbf{e}_{jr_j} | c_i)$$

$$\hat{f}(\mathbf{v}_j|c_i) = \frac{n_i(\mathbf{v}_j)}{n_i}$$

- Avoid zero probabilities
- Use Laplace Smoothing
- Introduce smoothing parameter: α
- Larger values of α move data likelihood to uniform distribution
- m_j is dimensionality

$$\hat{f}(\mathbf{v}_j|c_i) = \frac{n_i(\mathbf{v}_j) + 1}{n_i + m_j}$$

$$\hat{f}(\mathbf{v}_j|c_i) = \frac{n_i(\mathbf{v}_j) + \alpha}{n_i + \alpha * m_j}$$



Naive Bayes Classifier: Parameter Estimation

- Estimate posterior probability for each class c_i
- Predicted class has highest probability (\hat{y})
- Two approaches:
 - Maximum Likelihood Estimation (MLE)
 - Maximum a posteriori (MAP)

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

$$\arg \max_{c_i} \{ P(\mathbf{x}|c_i) P(c_i) \}$$

- Maximize posterior by maximizing data likelihood, $P(\mathbf{x}|c_i)$
- Similar to EM for GMM
 - Not “weighted” but only use samples from class to compute mean and variance for each feature and class

$$P(\mathbf{x}|c_i) = P(x_1, x_2, \dots, x_d | c_i) = \prod_{j=1}^d P(x_j | c_i)$$

$$P(x_j | c_i) \propto f(x_j | \hat{\mu}_{ij}, \hat{\sigma}_{ij}^2) = \frac{1}{\sqrt{2\pi}\hat{\sigma}_{ij}} \exp \left\{ -\frac{(x_j - \hat{\mu}_{ij})^2}{2\hat{\sigma}_{ij}^2} \right\}$$

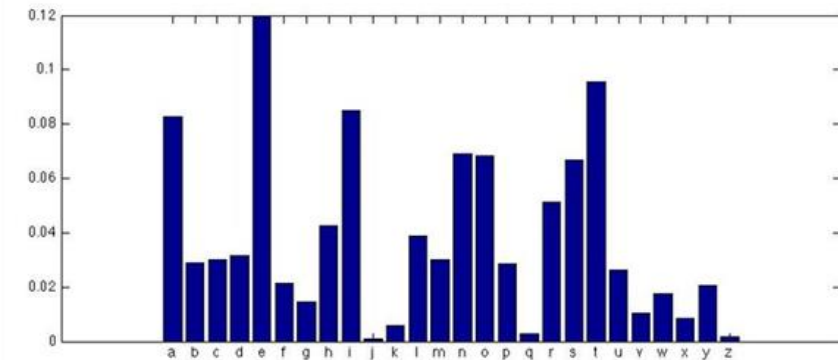
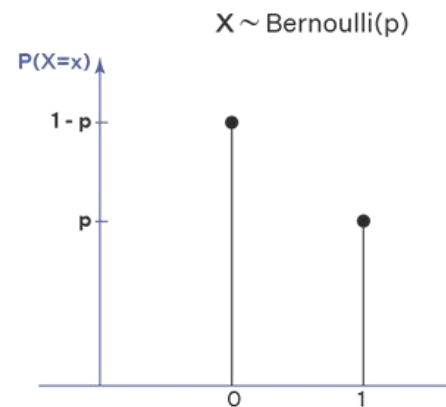
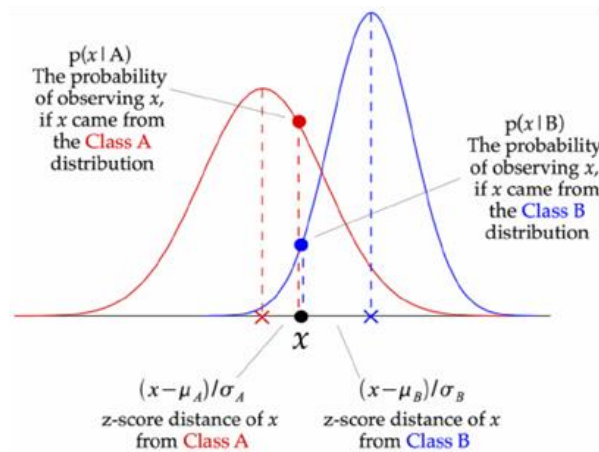
- Another option is to consider the most likely parameter value given the data, “maximum a posteriori” or MAP
- Now include prior with estimation

$$\arg \max_{\theta} P(\mathbf{X}|\theta)P(\theta)$$

Distributions Used with Naive Bayes



- Normal/Gaussian
- Bernoulli
- Multinomial



Distributions Used with Naive Bayes



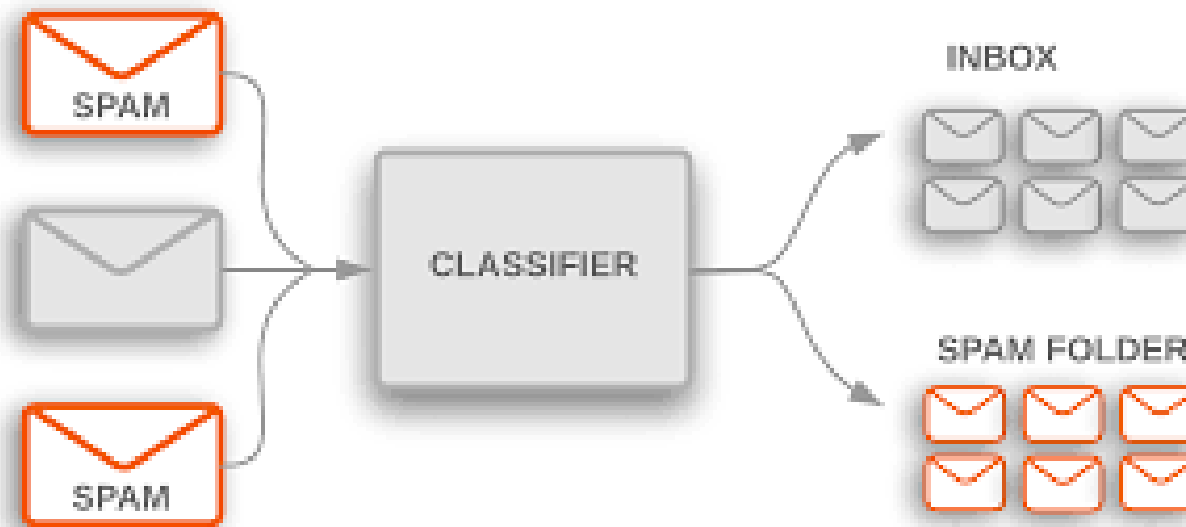
Classes	Features			
		Features are discrete Binary/Boolean functions	Features are discrete occurrence counts	Features are numerical values sampled from a continuous function with a Gaussian distribution
	Binary Labels {T,F}, {Pass,Fail}; {0,1}, {-1,1)	sklearn.naive_bayes.BernoulliNB() with 2 classes	sklearn.naive_bayes.MultinomialNB() with 2 classes	sklearn.naive_bayes.GaussianNB() with 2 classes
Multinomial Labels (3 or more classes) {0-6}, {airplane, bicycle, car}, {pass, fail type1, fail type 2, fail type 3}		sklearn.naive_bayes.BernoulliNB()	sklearn.naive_bayes.MultinomialNB()	sklearn.naive_bayes.GaussianNB()

* Use variant ComplementNB if unbalanced distribution and CategoricalNB if data is categorical {Monday, Tuesday, Wed....} rather than occurrences



Practical Applications of Naive Bayes









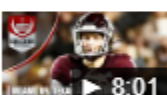








- [Jupyter Notebook](#)



Naive Bayes Application: Predicting Football Games



TEXAS A&M UNIVERSITY
Engineering

 SAMHOU	0	Final Sep 3	 Texas A&M	31		 App State	17	Final Sep 10	 Texas A&M	14	
 Miami FL	9	Final Sep 17	 Texas A&M	17		 Arkansas	21	Final Sat, Sep 24	 Texas A&M	23	
 Texas A&M	24	Final Sat, Oct 1	 23 MS State	42		 Texas A&M		Sat, Oct 8 7:00 PM	 1 Alabama		

Naive Bayes Application: Predicting Football Games



TEXAS A&M UNIVERSITY
Engineering

SAMHOU Texas A&M	0 31 ◀	Final Sep 3 ▶ 7:48	App State Texas A&M	17 ◀ 14	Final Sep 10 ▶ 9:39
Miami FL Texas A&M	9 17 ◀	Final Sep 17 ▶ 8:01	Arkansas Texas A&M	21 23 ◀	Final Sat, Sep 24 ▶ 8:27
Texas A&M 23 MS State	24 42 ◀	Final Sat, Oct 1 ▶ 11:01	Texas A&M 1 Alabama	Sat, Oct 8 7:00 PM	

$X = W, L, W, W, L$

Likelihood = $p^X * (1 - p)^{1-X}$

$$\text{MLE: } \theta_{ML} = \underset{\theta}{\operatorname{argmax}} L(x, \theta) \quad \theta_{MLE} : \hat{p} = \sum_{i=1}^D X_i / n = 0.6$$

$$\text{MAP: } \theta_{MAP} = \underset{\theta}{\operatorname{argmax}} L(x, \theta) * P(\theta) \quad \theta_{MAP} : \hat{p} = 0.615$$

Naive Bayes Application: Predicting Football Games



TEXAS A&M UNIVERSITY
Engineering

SAMHOU Texas A&M	0 31	Final Sep 3 7:48	App State Texas A&M	17 14	Final Sep 10 9:39
Miami FL Texas A&M	9 17	Final Sep 17 8:01	Arkansas Texas A&M	21 23	Final Sat, Sep 24 8:27
Texas A&M 23 MS State	24 42	Final Sat, Oct 1 11:01	Texas A&M 1 Alabama	Sat, Oct 8 7:00 PM	

$X = W, L, W, W, L$

MLE: $\theta_{ML} = \underset{\theta}{\operatorname{argmax}} L(X, \theta) \quad \theta_{MLE} : \hat{p} = \sum_{i=1}^D X_i / n = 0.6$

MAP: $\theta_{MAP} = \underset{\theta}{\operatorname{argmax}} L(X, \theta) * P(\theta) \quad \theta_{MAP} : \hat{p} = 0.615$

\hat{p}	prior	likelihood	posterior
0.538	0.181818	0.372044	0.067644
0.615	0.363636	0.380833	0.138485
0.666	0.090909	0.406964	0.036997
0.692	0.181818	0.426238	0.077498
0.844	0.090909	0.625548	0.056868
0.9	0.090909	0.739	0.067182

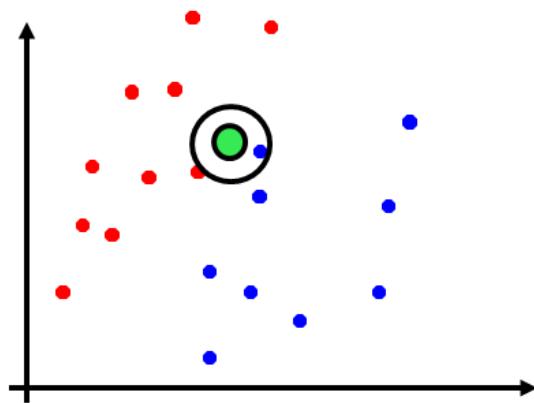


k-Nearest Neighbors Classifier

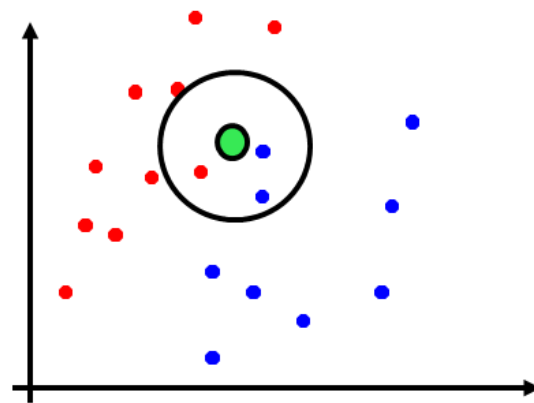
k-Nearest Neighbors Classification



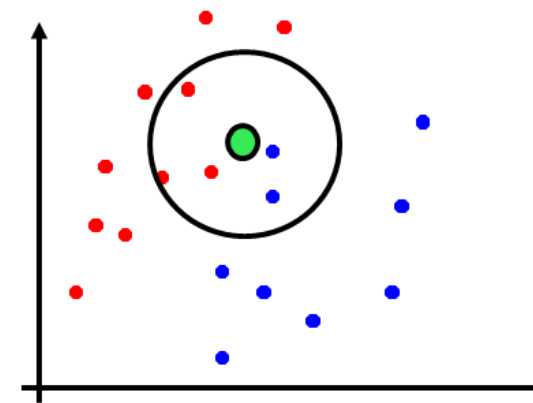
- Non-parametric and discriminative classifier
- Assign class based on the majority vote of the k -closest neighbors



$k=1$;
New sample =>



$k=3$;
New sample =>

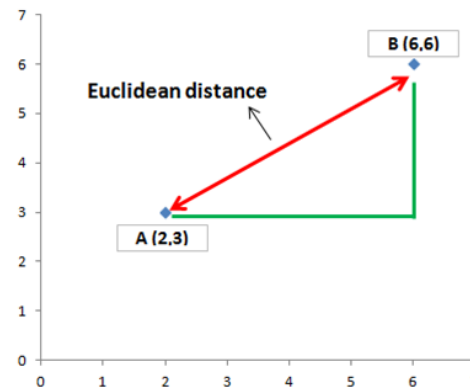


$k=5$;
New sample =>

k-Nearest Neighbors Classification



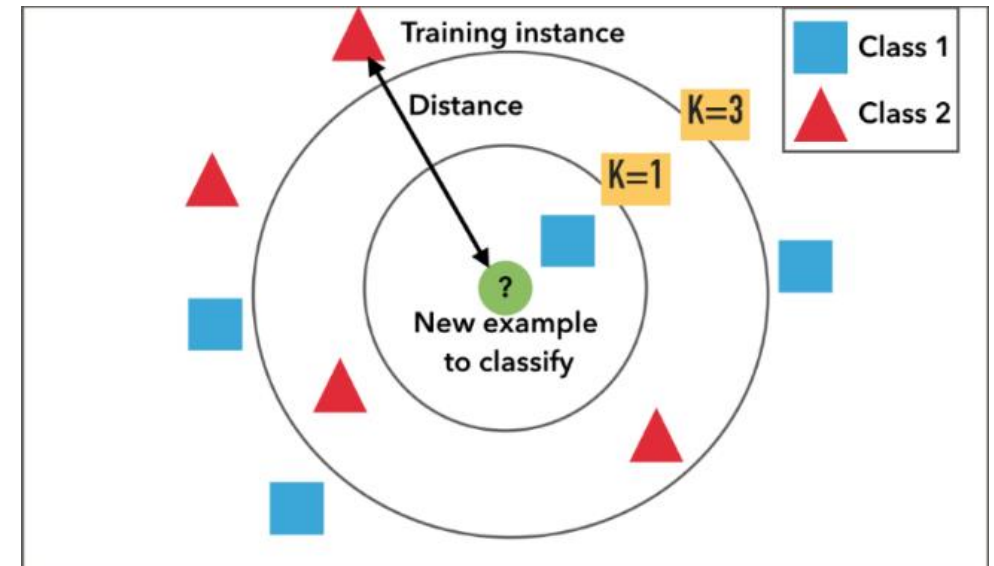
- Very simple nonparametric model
 - Select a distance function $L(x, x')$ (e.g., Euclidean)
 - Choose a hyperparameter K (usually odd)



$$\text{Euclidean distance } (a, b) = \sqrt{(a_1 - b_1)^2 + (a_2 - b_2)^2}$$

k-Nearest Neighbors Classification

- For each instance in the test set (new, unobserved) x ,
- Select in D the k examples that are nearest to x according to $L(\mathbf{x}, \mathbf{x}^i)$ and keep their index in set $\{s_1, \dots, s_K\}$, find:
 - For classification: $\hat{y} = \underset{c_i}{\operatorname{argmax}} K_i$
- In scikit-learn we use `KNeighborsClassifier()`

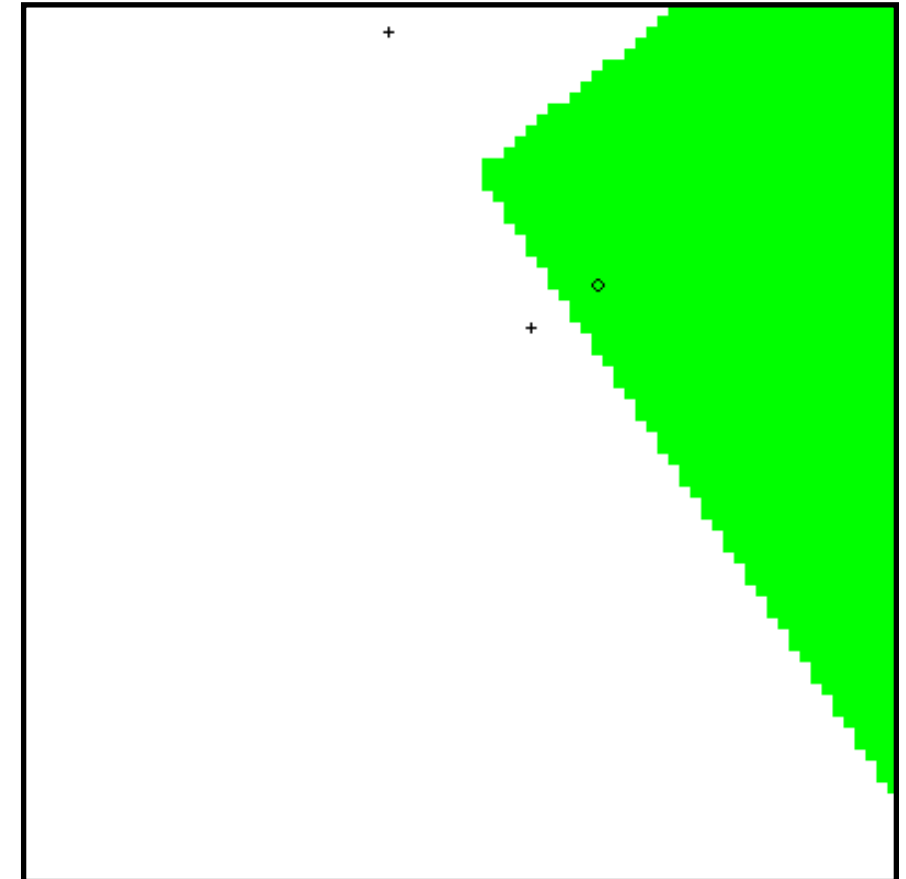


k-Nearest Neighbors Classification



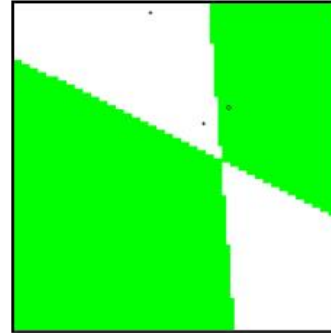
TEXAS A&M UNIVERSITY
Engineering

- Classification from similarity
 - Case-based reasoning
 - Predict an instance's label using similar instances
- Nearest-neighbor classification
 - 1-NN: copy the label of the most similar data point
 - K-NN: vote the k nearest neighbors
 - Key issue: how to define similarity and number of k



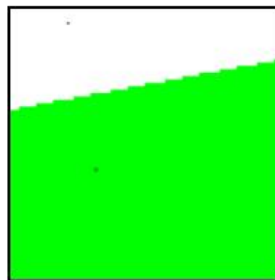
<http://www.cs.cmu.edu/~zhuxj/courseproject/knndemo/KNN.html>

Decision Boundaries for Varying Number Training Samples

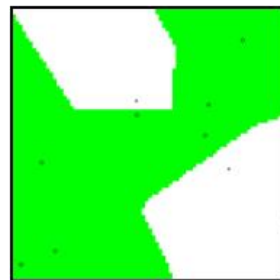


Truth

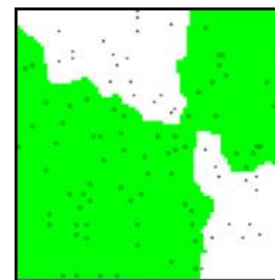
2 Observations



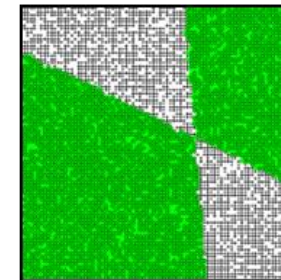
10
Observations



100
Observations



10000
Observations



<http://www.cs.cmu.edu/~zhuxj/courseproject/knndemo/KNN.html>

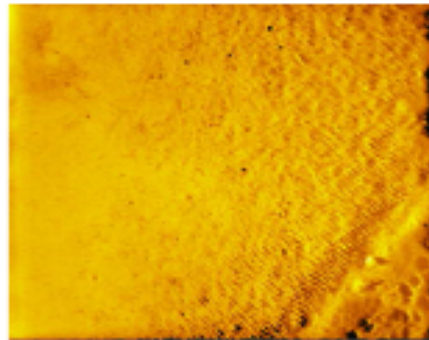


Practical Application of k-NN

Sonar Image Segmentation



TEXAS A&M UNIVERSITY
Engineering



(a) Original image



(b) PFLICM - Crater



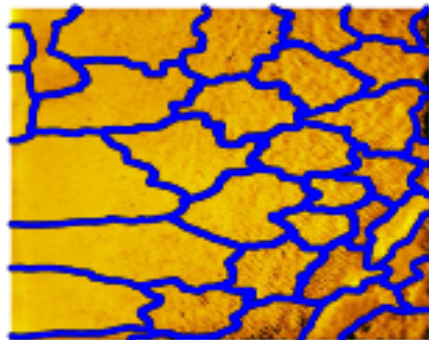
(c) PFLICM - Flat



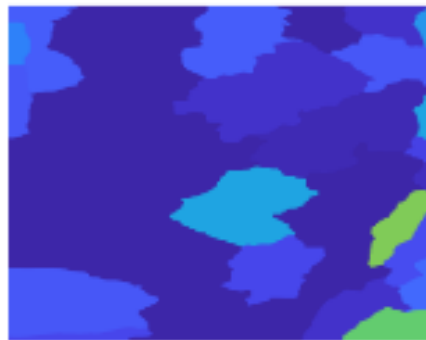
(d) PFLICM - Ripples



(e) PFLICM - Rocky



(f) Superpixels



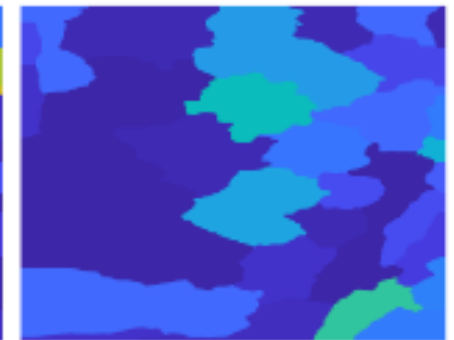
(g) PKNN - Crater



(h) PKNN - Flat



(i) PKNN - Ripples



(j) PKNN - Rocky

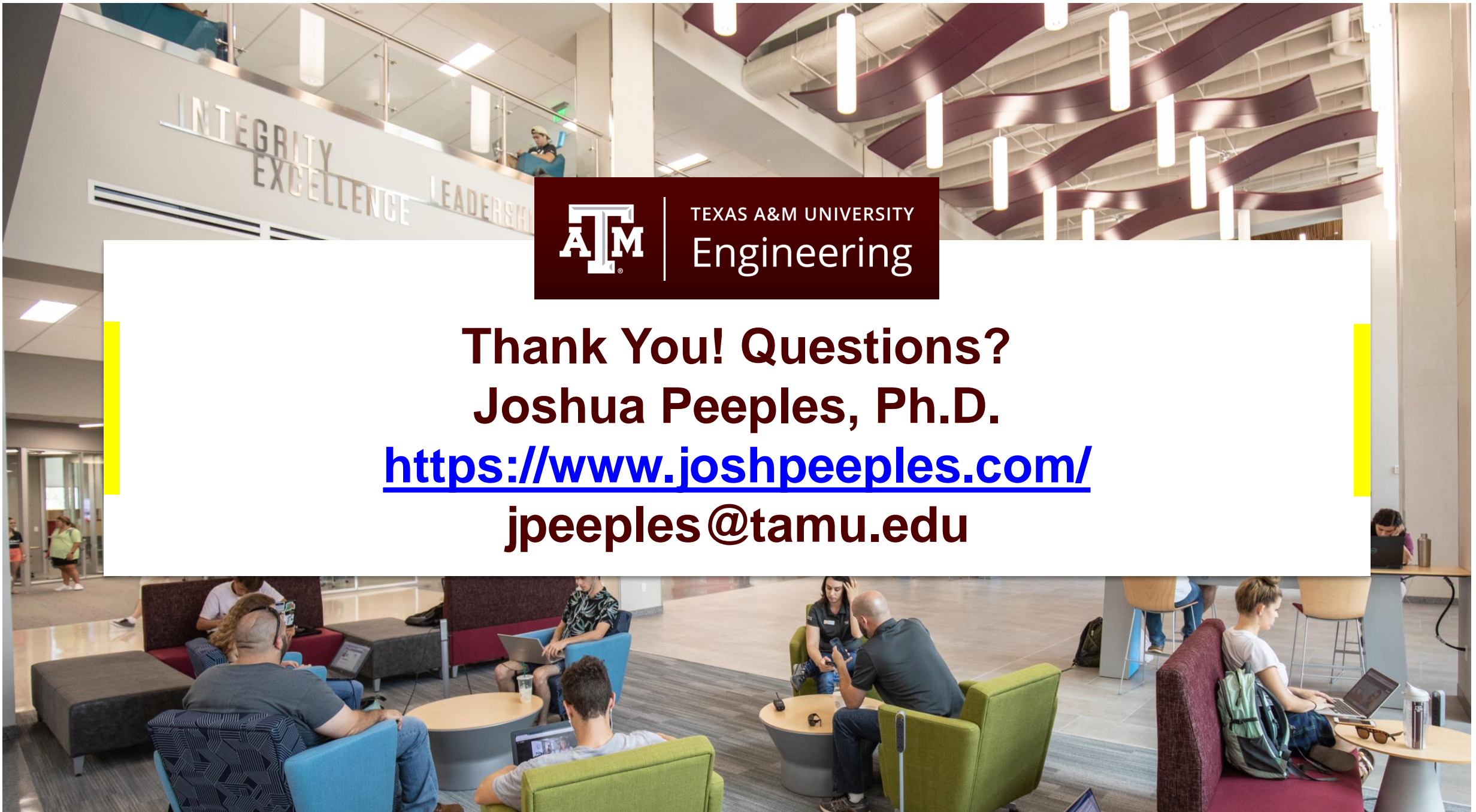
- Introduction to Classification II
 - Methodology
 - Metrics
 - Overfitting and underfitting
- Decision tree classification

INTEGRITY
EXCELLENCE LEADERSHIP



TEXAS A&M UNIVERSITY
Engineering

Thank You! Questions?
Joshua Peeples, Ph.D.
<https://www.joshpeeples.com/>
jpeeples@tamu.edu





TEXAS A&M UNIVERSITY
Engineering

Supplemental Slides

- [StatQuest: Bayes' Theorem](#)
- [StatQuest: Naive Bayes](#)
- [StatQuest: Gaussian Naive Bayes](#)