# Big Data and Data Mining

## Week 3/4: Classification



**Fenerbahce University**

# Instructors

Assist. Prof. Vecdi Emre Levent
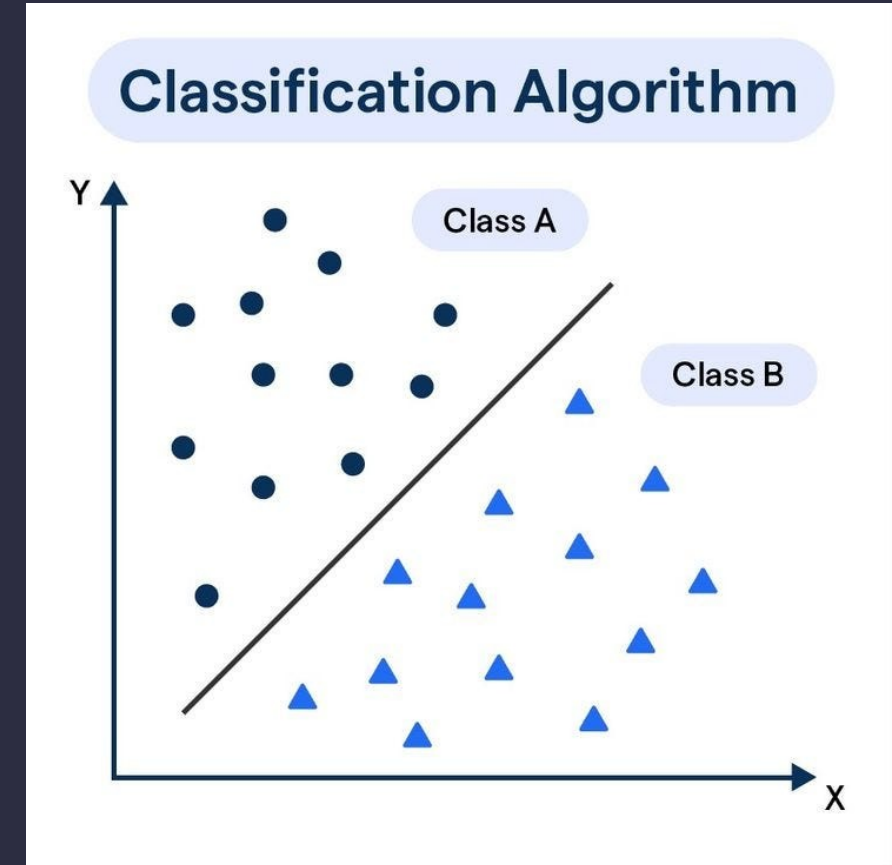
Office: 311

Email : emre.levent@fbu.edu.tr

# Classification

- What is classification?

- Issues regarding classification

- Bayesian Classification

- Classification by decision tree induction

- Classification by Neural Networks

- Classification by Support Vector Machines (SVM)

- Instance Based Methods

- Classification accuracy

- Summary

# Classification

- ## Classification:
  - predicts categorical class labels
  - classifies data (constructs a model) based on the training set and the values (class labels) in a classifying attribute and uses it in classifying new data

# Classification

- Classification:
- Typical Applications
  - credit approval
  - target marketing
  - medical diagnosis
  - treatment effectiveness analysis



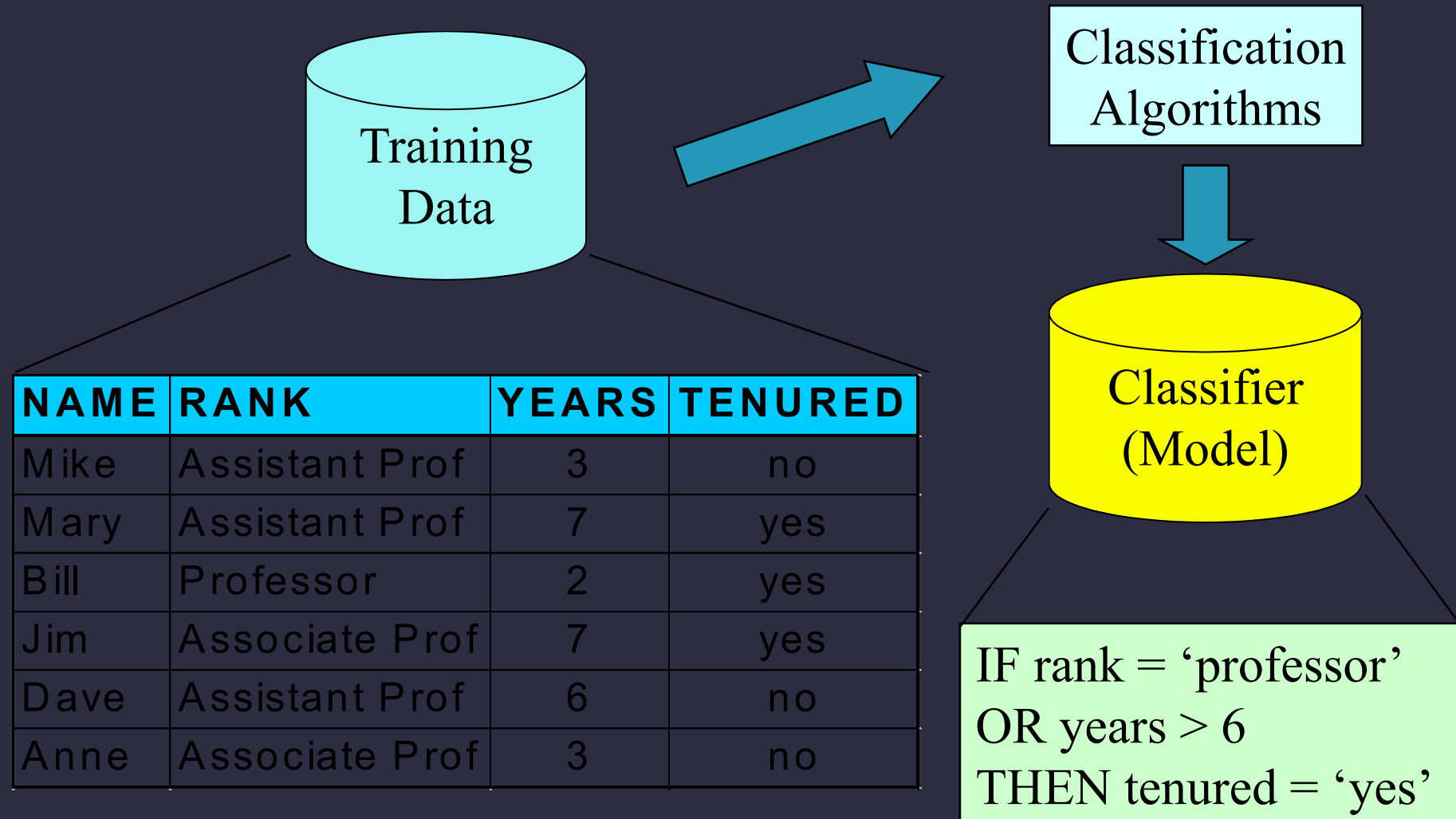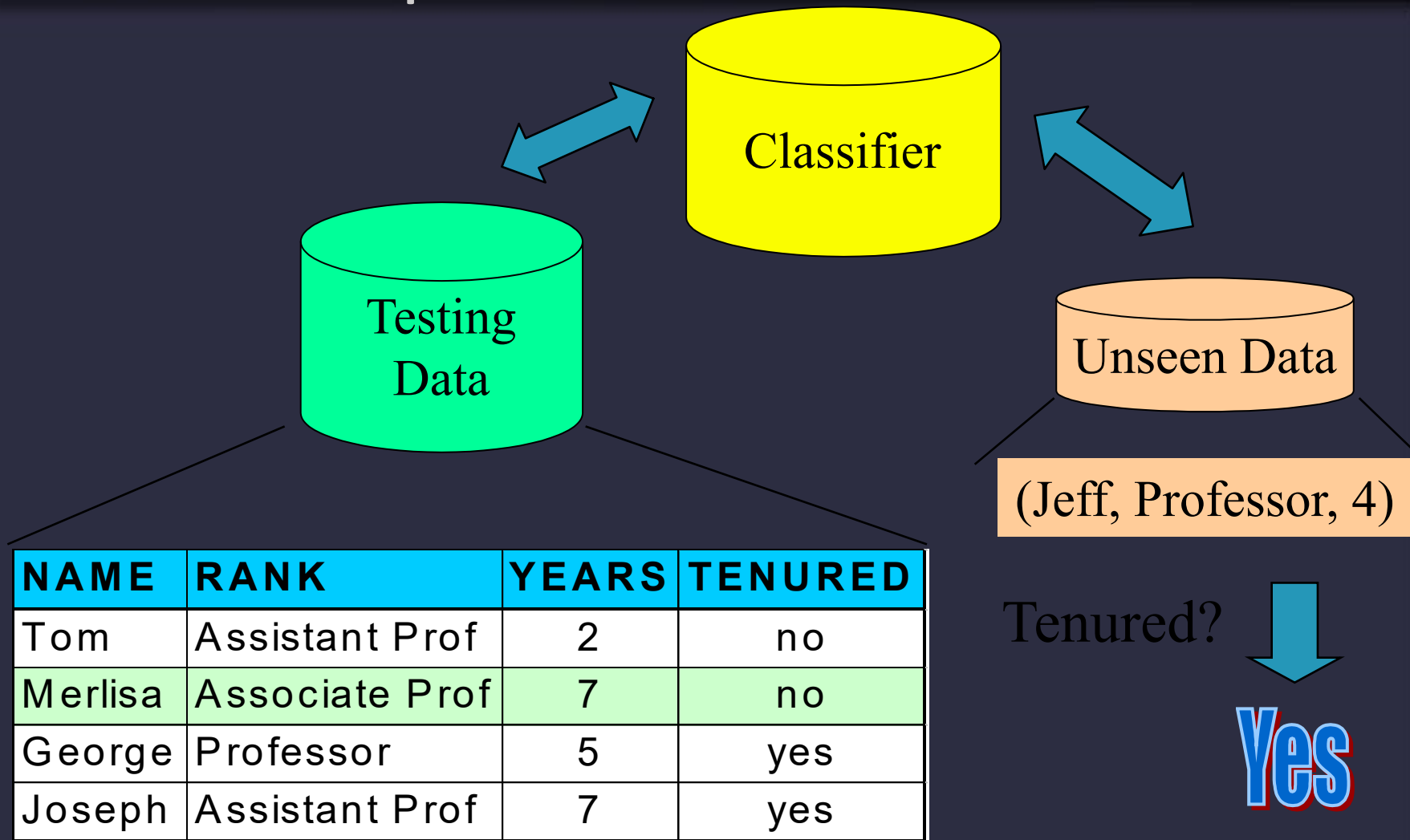COVID-19 Radiography Database — COVID-19, lung opacity, normal, viral pneumonia

# Classification—A Two-Step Process

- Model construction: describing a set of predetermined classes
  - Each sample is assumed to belong to a predefined class, as determined by the class label attribute
  - The set of sample used for model construction is training set
  - The model is represented as classification rules, decision trees, or mathematical formula or AI

Training Data

Classification Algorithms

| NAME | RANK | YEARS | TENURED |
|------|------|-------|---------|
| Mike | Assistant Prof | 3 | no |
| Mary | Assistant Prof | 7 | yes |
| Bill | Professor | 2 | yes |
| Jim | Associate Prof | 7 | yes |
| Dave | Assistant Prof | 6 | no |
| Anne | Associate Prof | 3 | no |

Classifier (Model)

IF rank = 'professor'
OR years > 6
THEN tenured = 'yes'

# Classification—A Two-Step Process
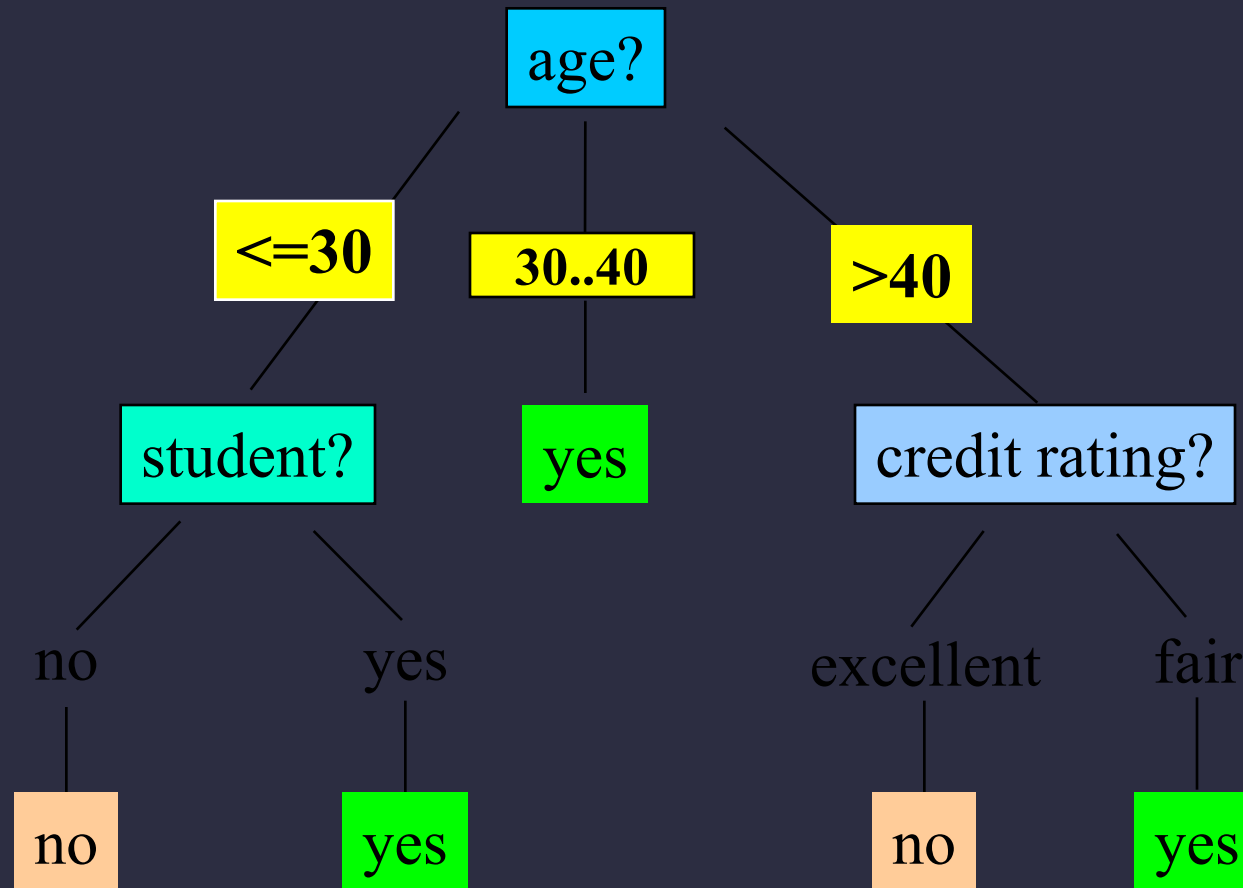
- Model usage: for classifying future or unknown objects
  - Estimate accuracy of the model
    - The known label of test sample is compared with the classified result from the model
    - Accuracy rate is the percentage of test set samples that are correctly classified by the model
    - Test set is independent of training set, otherwise over-fitting will occur
  - If the accuracy is acceptable, use the model to classify data tuples whose class labels are not known

**Classifier**

**Testing Data**

**Unseen Data**

(Jeff, Professor, 4)

Tenured?

Yes

| NAME | RANK | YEARS | TENURED |
|---|---|---|---|
| Tom | Assistant Prof | 2 | no |
| Merlisa | Associate Prof | 7 | no |
| George | Professor | 5 | yes |
| Joseph | Assistant Prof | 7 | yes |

# Dataset for computer buyers

| age | income | student | credit_rating | buys_computer |
|---|---|---|---|---|
| <=30 | high | no | fair | no |
| <=30 | high | no | excellent | no |
| 31…40 | high | no | fair | yes |
| >40 | medium | no | fair | yes |
| >40 | low | yes | fair | yes |
| >40 | low | yes | excellent | no |
| 31…40 | low | yes | excellent | yes |
| <=30 | medium | no | fair | no |
| <=30 | low | yes | fair | yes |
| >40 | medium | yes | fair | yes |
| <=30 | medium | yes | excellent | yes |
| 31…40 | medium | no | excellent | |
| 31…40 | high | yes | fair | |
| >40 | medium | no | excellent | |

# A Decision Tree for *"buys_computer"*

# Classification

- What is classification
- Issues regarding classification
- Bayesian Classification
- Classification by decision tree induction
- Classification by Neural Networks
- Classification by Support Vector Machines (SVM)
- Instance Based Methods
- Classification accuracy
- Summary

# Issues (1): Data Preparation

- Data cleaning
  - Preprocess data in order to reduce noise and handle missing values

# Issues (1): Data Preparation

- Relevance analysis (feature selection)
  - Remove the irrelevant or redundant attributes



(a) All features

(b) Informative features

# Issues (1): Data Preparation

- Data transformation
  - Generalize and/or normalize data



Generalisation        Overfitting

# Issues (1): Data Preparation

- Data transformation
  - Generalize and/or normalize data

- Predictive accuracy

# Issues (2): Evaluating Classification Methods

- Speed and scalability
  - time to construct the model
  - time to use the model



Model Execution Time (min) Before v.s. After Optimizations

# Issues (2): Evaluating Classification Methods

- Robustness
  - handling noise and missing values



**Handling Missing Data, Outliers and noisy data**

Outliers →

- Interpretability:
  - understanding and insight provided by the model

- Goodness of rules
  - decision tree size
  - compactness of classification rules

Grasshoppers

Katydids

Antenna Length

Abdomen Length

Dr. V. E. Levent    Big Data and Data Mining

# With a lot of data, we can build a histogram. Let us just build one for "Antenna Length" for now…



Katydids

Grasshoppers

We can leave the histograms as they are, or we can summarize them with two normal distributions.

• We want to classify an insect we have found. Its antennae are 3 units long. How can we classify it?

• We can just ask ourselves, give the distributions of antennae lengths we have seen, is it more *probable* that our insect is a **Grasshopper** or a **Katydid**.

• There is a formal way to discuss the most *probable* classification…

$p(c_j \mid d)$ = probability of class $c_j$, *given* that we have observed $d$

**3**

Antennae length is **3**

P(**Grasshopper** | **3** ) = 10 / (10 + 2) = 0.833

P(**Katydid** | **3** ) = 2 / (10 + 2) = 0.166



10

2

**3**

Antennae length is **3**

P(**Grasshopper** | **7** ) = 3 / (3 + 9) = 0.250

P(**Katydid** | **7** ) = 9 / (3 + 9) = 0.750



9

3

**7**

Antennae length is **7**

P(**Grasshopper** | **5** ) = 6 / (6 + 6)           = 0.500

P(**Katydid** | **5** )           = 6 / (6 + 6)           = 0.500



6  6

Antennae length is **5**

# Bayes Classifiers

That was a visual intuition for a simple case of the Bayes classifier, also called:

- Idiot Bayes
- Naïve Bayes
- Simple Bayes

*Find out the probability of the previously unseen instance belonging to each class, then simply pick the most probable class.*

# Bayes Classifiers

Bayesian classifiers use **Bayes theorem**, which says

$$p(c_j \mid d) = \frac{p(d \mid c_j)\, p(c_j)}{p(d)}$$

- $p(c_j \mid d)$ = probability of instance $d$ being in class $c_j$, This is what we are trying to compute

- $p(d \mid c_j)$ = probability of generating instance $d$ given class $c_j$,
  We can imagine that being in class $c_j$, causes you to have feature $d$ with some probability

- $p(c_j)$ = probability of occurrence of class $c_j$,
  This is just how frequent the class $c_j$, is in our database

- $p(d)$ = probability of instance $d$ occurring

Assume that we have two classes

$c_1$ = male, and $c_2$ = female.

We have a person whose sex we do not know, say "*drew*" or *d*.

Classifying *drew* as male or female is equivalent to asking is it more probable that *drew* is male or female, I.e which is greater $p(\text{male}\,|\,drew)$ or $p(\text{female}\,|\,drew)$

(Note: "Drew can be a male or female name")


Drew Barrymore


Drew Carey

What is the probability of being called "*drew*" given that you are a male?

What is the probability of being a male?

$$p(\text{male}\,|\,drew) = \frac{p(drew\,|\,\text{male})\,p(\text{male})}{p(drew)}$$

What is the probability of being named "*drew*"?
(actually irrelevant, since it is that same for all classes)

# Is Officer Drew a Male or Female?

Luckily, we have a small
database with names and
sex.

We can use it to apply
Bayes rule…

**Officer Drew**

$$p(c_j \mid d) = \frac{p(d \mid c_j)\, p(c_j)}{p(d)}$$

| Name | Sex |
|---------|--------|
| Drew | Male |
| Claudia | Female |
| Drew | Female |
| Drew | Female |
| Alberto | Male |
| Karin | Female |
| Nina | Female |
| Sergio | Male |

**Officer Drew**

| Name | Sex |
|------|-----|
| Drew | Male |
| Claudia | Female |
| Drew | Female |
| Drew | Female |
| Alberto | Male |
| Karin | Female |
| Nina | Female |
| Sergio | Male |

$$p(c_j \mid d) = \frac{p(d \mid c_j) \, p(c_j)}{p(d)}$$

$$p(\text{male} \mid drew) = \frac{p(drew \mid \text{male}) \, p(\text{male})}{p(drew)}$$

$$p(\text{male} \mid drew) = \frac{1/3 \ast 3/8}{3/8} = \frac{0.125}{3/8}$$

$$p(\text{female} \mid drew) = \frac{2/5 \ast 5/8}{3/8} = \frac{0.250}{3/8}$$

Officer Drew is more likely to be a Female.

So far we have only considered Bayes Classification when we have one attribute (the "*antennae length*", or the "*name*").
But we may have many features.
How do we use all the features?

$$p(c_j \mid d) = \frac{p(d \mid c_j)\, p(c_j)}{p(d)}$$

| Name | Over 170CM | Eye | Hair length | Sex |
|------|-----------|------|-------------|------|
| Drew | No | Blue | Short | Male |
| Claudia | Yes | Brown | Long | Female |
| Drew | No | Blue | Long | Female |
| Drew | No | Blue | Long | Female |
| Alberto | Yes | Brown | Short | Male |
| Karin | No | Blue | Long | Female |
| Nina | Yes | Brown | Short | Female |
| Sergio | Yes | Blue | Long | Male |

- To simplify the task, **naïve Bayesian classifiers** assume attributes have independent distributions, and thereby estimate

$$p(d|c_j) = p(d_1|c_j) * p(d_2|c_j) * ....* p(d_n|c_j)$$

The probability of class $c_j$ generating instance $d$, equals….

The probability of class $c_j$ generating the observed value for feature 1, multiplied by..

The probability of class $c_j$ generating the observed value for feature 2, multiplied by..

$$p(\text{male} \mid drew) = \frac{p(drew \mid \text{male}) \, p(\text{male})}{p(drew)}$$

- To simplify the task, **naïve Bayesian classifiers** assume attributes have independent distributions, and thereby estimate

$$p(d|c_j) = p(d_1|c_j) * p(d_2|c_j) * \ldots* p(d_n|c_j)$$

$$p(\text{officer drew}|c_j) = p(\text{over\_170}_{cm} = \text{yes}|c_j) * p(\text{eye} =blue|c_j) * \ldots.$$

Officer Drew is blue-eyed, over 170$_{cm}$ tall, and has long hair

$$p(\text{officer drew}| \text{Female}) = \; 2/5 \; * \; 3/5 \; * \; \ldots.$$

$$p(\text{officer drew}| \text{Male}) \; = \; 2/3 \; * \; 2/3 \; * \; \ldots.$$

The Naive Bayes classifiers is often represented as this type of graph…

Note the direction of the arrows, which state that each class causes certain features, with a certain probability

$Cj$

$p(d_1|c_j)$  $p(d_2|c_j)$  $\cdots$  $p(d_n|c_j)$

# Naïve Bayes is fast and space efficient

$$Cj$$

We can look up all the probabilities with a single scan of the database and store them in a (small) table…

$$p(d_1|c_j)$$  $$p(d_2|c_j)$$  $$\ldots$$  $$p(d_n|c_j)$$

| Sex | Over190$_{cm}$ | |
|------|------|------|
| Male | Yes | 0.15 |
| | No | 0.85 |
| Female | Yes | 0.01 |
| | No | 0.99 |

| Sex | Long Hair | |
|------|------|------|
| Male | Yes | 0.05 |
| | No | 0.95 |
| Female | Yes | 0.70 |
| | No | 0.30 |

| Sex | | |
|------|------|------|
| Male | | |
| | | |
| Female | | |
| | | |

Suppose we are trying to classify a persons sex based on several features, including eye color. (Of course, eye color is completely irrelevant to a persons gender)

$p(\text{Jessica} \mid c_j) = p(\text{eye} = \text{brown} \mid c_j) * p(\text{wears\_dress} = \text{yes} \mid c_j) * ….$

$p(\text{Jessica} \mid \text{Female}) = 9,000/10,000 \qquad * 9,975/10,000 \quad * ….$

$p(\text{Jessica} \mid \text{Male}) = 9,001/10,000 \qquad * 2/10,000 \qquad * ….$

Almost the same!

We can have an arbitrary number of classes, or feature values

$Cj$

$p(d_1|c_j)$    $p(d_2|c_j)$    . . .    $p(d_n|c_j)$

| Animal | Mass >10$_{kg}$ | |
|---|---|---|
| Cat | Yes | 0.15 |
| | No | 0.85 |
| Dog | Yes | 0.91 |
| | No | 0.09 |
| Pig | Yes | 0.99 |
| | No | 0.01 |

| Animal | Color | |
|---|---|---|
| Cat | Black | 0.33 |
| | White | 0.23 |
| | Brown | 0.44 |
| Dog | Black | 0.97 |
| | White | 0.03 |
| | Brown | 0.90 |
| Pig | Black | 0.04 |
| | White | 0.01 |

| Animal |
|---|
| Cat |
| Dog |
| Pig |

Katydids

Grasshoppers

Ants

Top plot (time series): Background noise, Bee begins to cross laser

Bottom plot: Single-Sided Amplitude Spectrum of Y(t), 60Hz interference, Peak at 197Hz, Harmonics; y-axis |Y(f)| ×10³; x-axis Frequency (Hz)

Frequency (Hz)

Frequency (Hz)

Frequency (Hz)

$$\frac{1}{\sqrt{2\pi\sigma^2}}e^{-\frac{(x-w)^2}{2\sigma^2}}$$

Wing Beat Frequency Hz

Wing Beat Frequency Hz

*Anopheles stephensi*: Female
mean =475, Std = 30

*Aedes aegyptii* : Female
mean =567, Std = 43

517

If I see an insect with a wingbeat frequency of 500, what is it?

517

What is the error rate?

Can we get more features?

12.2% of the area under the pink curve

8.02% of the area under the red curve

$$\frac{1}{\sqrt{2\pi}\,\sigma} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$

$$P(Anopheles|wingbeat = 500) = \frac{1}{\sqrt{2\pi}\,30} e^{-\frac{(500-475)^2}{2\times30^2}}$$

# Circadian Features



*Aedes aegypti* (yellow fever mosquito)

*dawn*

*dusk*

0

0    12    24

Midnight    Noon    Midnight

Suppose I observe an insect with a wingbeat frequency of 420 at 11:00am

What is it?

(Culex | [420Hz,11:00am])        = (6/ (6 + 6 + 0))  * (2/ (2 + 4 + 3))  = 0.111

(Anopheles | [420Hz,11:00am])  = (6/ (6 + 6 + 0))  * (4/ (2 + 4 + 3))  = 0.222

(Aedes | [420Hz,11:00am])        = (0/ (6 + 6 + 0))  * (3/ (2 + 4 + 3))  = 0.000

# Advantages/Disadvantages of Naïve Bayes

- Advantages:
  - Fast to train (single scan). Fast to classify
  - Not sensitive to irrelevant features
  - Handles real and discrete data
  - Handles streaming data well

- Disadvantages:
  - Assumes independence of features

# Naïve Bayes

Load Iris Dataset and goto Classify Tab

# Naïve Bayes

Select Naive Bayes Classifier

# Naïve Bayes

## Cross-Validation

- 10 percent will used for only validation

# Naïve Bayes

## Cross-Validation

- 10 percent will used for only validation

# Naïve Bayes

Start

# Naïve Bayes

Start

Percentage Split

# Nearest Neighbor Classifier

# Nearest Neighbor Classification

Given a training dataset $\mathcal{D} = \left\{ y^{(n)}, \mathbf{x}^{(n)} \right\}_{n=1}^{N}, \quad y \in \{1, \ldots, C\}, \quad \mathbf{x} \in \mathbb{R}^M$

and a test input $\mathbf{x}_{test}$, predict the class label

1) Find the closest point in the training data to $\mathbf{x}_{test}$
2) Return the class label of that closest point

Need distance function! What should $d(\mathbf{x}, \mathbf{z})$ be?

# 3-Nearest Neighbor (kNN) classifier

# 5-Nearest Neighbor (kNN) classifier



Whales

Seals

Sharks

# What is the best k?

How do we choose a learner that is accurate and also generalizes to unseen data?

- Larger k → predicted label is more stable
- Smaller k → predicted label is more affected by individual training points

But how to choose $k$?

# k-NN: Details

Inductive Bias:
1. Close points should have similar labels
2. All dimensions are created equally!



Example: two features for k-NN

**big problem:** feature scale could dramatically influence classification results

# KNN Mini Project

Visit to view project specifications

http://levent.tc/files/courses/big_data/mini_projects/knn/proje1_knn.pdf

Prepare a presentation

# KNN on Weka

## Select Classifier

Settings

# KNN on Weka

## Classification

# Decision Tree

## Problem Setting

- Set of possible instances $X$

- Set of possible labels $Y$

- Unknown target function $f : X \rightarrow Y$

- Set of function hypotheses $H = \{h \mid h : X \rightarrow Y\}$

**Input:** Training examples of unknown target function $f$
$$\{<\mathbf{x}_i, y_>\}^n_{i=1} = \{<\mathbf{x}_1, y_{1>}, \ldots, <\mathbf{x}_n, y_{n>}\}$$

**Output:** Best approximates $f$

# Sample Dataset

- Columns denote features $X_i$

- Rows denote labeled instances $\mathbf{x}_i, y_i$

- Class label denotes whether a tennis game was played

$\mathbf{x}_i, y_i$

| Predictors | | | | Response |
|---|---|---|---|---|
| **Outlook** | **Temperature** | **Humidity** | **Wind** | **Class** |
| Sunny | Hot | High | Weak | No |
| Sunny | Hot | High | Strong | No |
| Overcast | Hot | High | Weak | Yes |
| Rain | Mild | High | Weak | Yes |
| Rain | Cool | Normal | Weak | Yes |
| Rain | Cool | Normal | Strong | No |
| Overcast | Cool | Normal | Strong | Yes |
| Sunny | Mild | High | Weak | No |
| Sunny | Cool | Normal | Weak | Yes |
| Rain | Mild | Normal | Weak | Yes |
| Sunny | Mild | Normal | Strong | Yes |
| Overcast | Mild | High | Strong | Yes |
| Overcast | Hot | Normal | Weak | Yes |
| Rain | Mild | High | Strong | No |

# Decision Tree

- A possible decision tree for the data:



- Each internal node: test one attribute $X_i$
- Each branch from a node: selects one value for $X_i$
- Each leaf node: predict $Y$

- A possible decision tree for the data:



- What prediction would we make for

<outlook=sunny, temperature=hot, humidity=high, wind=weak> ?

- If features are continuous, internal nodes can test the value of a feature against a threshold

# Decision Tree Learning

**Problem Setting**:

- Set of possible instances $X$

  - each instance $x$ in $X$ is a feature vector
  - e.g., *<Humidity=low, Wind=weak, Outlook=rain, Temp=hot>*
- Unknown target function $f : X \rightarrow Y$
  - *Y* is discrete valued
- Set of function hypotheses $H$={ $h$ | $h : X \rightarrow Y$ }

  - each hypothesis $h$ is a decision tree
  - trees sorts $x$ to leaf, which assigns $y$

# Stages of (Batch) Machine Learning

**Given:** labeled training data $X, Y = \{h\mathbf{x}_i, y_i i\}_{i=1}^{n}$

- Assumes each $\mathbf{x}_i \leftarrow D(X)$ with $y_i = f_{target}(\mathbf{x}_i)$

**Train the model:**

$model \leftarrow classifier.\text{train}(X, Y)$

**Apply the model to new data:**

- Given: new unlabeled instance $\mathbf{x} \leftarrow D(X)$

$y_{\text{prediction}} \leftarrow model.\text{predict}(\mathbf{x})$

$X, Y$

learner

$\mathbf{x} \rightarrow$ $model$ $\rightarrow y_{\text{prediction}}$

# Decision Tree Induced Partition

# Decision Tree – Decision Boundary

- Decision trees divide the feature space into axis-parallel (hyper-)rectangles

- Each rectangular region is labeled with one label
  - or a probability distribution over labels



Decision boundary

# Expressiveness

- Decision trees can represent any boolean function of the input attributes



Truth table row → path to leaf

- In the worst case, the tree will require exponentially many nodes

# Expressiveness

Decision trees have a variable-sized hypothesis space

- As the #nodes (or depth) increases, the hypothesis space grows
  - Depth 1 ("decision stump"): can represent any boolean function of one feature
  - Depth 2: any boolean fn of two features; some involving three features (e.g., $(x_1 \land x_2) \lor (\lnot x_1 \land \lnot x_3)$ )
  - etc.

# Another Example: Restaurant Domain

Model a patron's decision of whether to wait for a table at a restaurant

| Example | Attributes | | | | | | | | | | Target |
|---------|-----|-----|-----|-----|------|-------|------|-----|--------|-------|-------|
| | Alt | Bar | Fri | Hun | Pat | Price | Rain | Res | Type | Est | Wait |
| $X_1$ | T | F | F | T | Some | $$$ | F | T | French | 0–10 | T |
| $X_2$ | T | F | F | T | Full | $ | F | F | Thai | 30–60 | F |
| $X_3$ | F | T | F | F | Some | $ | F | F | Burger | 0–10 | T |
| $X_4$ | T | F | T | T | Full | $ | F | F | Thai | 10–30 | T |
| $X_5$ | T | F | T | F | Full | $$$ | F | T | French | >60 | F |
| $X_6$ | F | T | F | T | Some | $$ | T | T | Italian | 0–10 | T |
| $X_7$ | F | T | F | F | None | $ | T | F | Burger | 0–10 | F |
| $X_8$ | F | F | F | T | Some | $$ | T | T | Thai | 0–10 | T |
| $X_9$ | F | T | T | F | Full | $ | T | F | Burger | >60 | F |
| $X_{10}$ | T | T | T | T | Full | $$$ | F | T | Italian | 10–30 | F |
| $X_{11}$ | F | F | F | F | None | $ | F | F | Thai | 0–10 | F |
| $X_{12}$ | T | T | T | T | Full | $ | F | F | Burger | 30–60 | T |

~7,000 possible cases

# A Decision Tree from Introspection



Is this the best decision tree?

# Decision Tree

- The smallest decision tree that correctly classifies all of the training examples is best

  - Finding the provably smallest decision tree is NP-hard
  - ...So instead of constructing the absolute smallest tree consistent with the training examples, construct one that is pretty small

*node* = root of decision tree

Main loop:

1. $A \leftarrow$ the "best" decision attribute for the next node.
2. Assign *A* as decision attribute for *node*.
3. For each value of *A*, create a new descendant of *node*.
4. Sort training examples to leaf nodes.
5. If training examples are perfectly classified, stop.
   Else, recurse over new leaf nodes.

How do we choose which attribute is best?

# Choosing the Best Attribute

**Key problem**: choosing which attribute to split a given set of examples

- Some possibilities are:
  - **Random:** Select any attribute at random
  - **Least-Values:** Choose the attribute with the smallest number of possible values
  - **Most-Values:** Choose the attribute with the largest number of possible values
  - **Max-Gain:** Choose the attribute that has the largest expected *information gain*
    - i.e., attribute that results in smallest expected size of subtrees rooted at its children

# Choosing an Attribute

**Idea**: a good attribute splits the examples into subsets that are (ideally) "all positive" or "all negative"



Which split is more informative: *Patrons?* or *Type?*

# ID3-induced Decision Tree

# Compare the Two Decision Trees

# Information Gain

Which test is more informative?



**Split over whether Balance exceeds 50K**

Less or equal 50K    Over 50K

**Split over whether applicant is employed**

Unemployed    Employed

# Information Gain

**Impurity/Entropy** (informal)

    – Measures the level of **impurity** in a group of examples

# Impurity

**Very impure group**

**Less impure**

**Minimum impurity**

# Select Decision Tree Based Algorithms

# Select Decision Tree Based Algorithms

# Neural Networks

- Analogy to Biological Systems (Indeed a great example of a good learning system)

- Massive Parallelism allowing for computational efficiency

- The first learning algorithm came in 1959 (Rosenblatt) who suggested that if a target output value is provided for a single neuron with fixed inputs, one can incrementally change weights to learn to produce these outputs using the perceptron learning rule

# Neural Network

"Combined logistic models"

Inputs

Age 34 .6

Gender 2 .1

Stage 4 .7

.5

.8

Σ

Output

0.6

"Probability of beingAlive"

Independent variables

Weights

Hidden Layer

Weights

Dependent variable

Prediction

Inputs

Age

34

Gender

2

Stage

4

.2

.3

.2

.5

.8

Σ

Output

0.6

"Probability of beingAlive"

Independent variables

Weights

Hidden Layer

Weights

Dependent variable

Prediction

Inputs

Age 34

Gender 1

Stage 4

.6
.2
.1
.3
.7
.2

.5
.8

Σ

Output

0.6

"Probability of beingAlive"

Independent variables

Weights

Hidden Layer

Weights

Dependent variable

Prediction

Age 34

Gender 2

Stage 4

.6
.2
.1
.3
.7
.2

Σ ∫ .4

Σ ∫ .2

.5
.8

Σ ∫ 0.6

"Probability of beingAlive"

*Independent variables*

Weights

Hidden Layer

Weights

*Dependent variable*

*Prediction*

# Neural Networks

– Example: Neural Network w/1 Hidden Layer

– Example: Neural Network w/2 Hidden Layers

– Example: Feed Forward Neural Network

# Neural Network

## Neural Network for Classification

Output

y

Hidden Layer

$z_1$   $z_2$   ...   $z_D$

Input

$x_1$   $x_2$   $x_3$   ...   $x_M$

(E) **Output (sigmoid)**
$$y = \frac{1}{1+\exp(-b)}$$

(D) **Output (linear)**
$$b = \sum_{j=0}^{D} \beta_j z_j$$

(C) **Hidden (sigmoid)**
$$z_j = \frac{1}{1+\exp(-a_j)}, \quad \forall j$$

(B) **Hidden (linear)**
$$a_j = \sum_{i=0}^{M} \alpha_{ji} x_i, \quad \forall j$$

(A) **Input**
Given $x_i, \quad \forall i$

23

- Question:

- Suppose you are training a one-hidden layer neural network with sigmoid activations for binary classification.



**True or False:** There is a unique set of parameters that maximize the likelihood of the dataset above.

# Neural Network Architectures

Even for a basic Neural Network, there are many design decisions to make:

1. # of hidden layers (depth)
2. # of units per hidden layer (width)
3. Type of activation function (nonlinearity)
4. Form of objective function

*Q: How many hidden units, D, should we use?*



Output

y

Hidden Layer

$z_1$   $z_2$   ...   $z_D$

D = M

Input

$x_1$   $x_2$   ...   $x_M$

*Q: How many hidden units, D, should we use?*



Output

y

Hidden Layer

$z_1$    $z_2$    $\cdots$    $z_D$

$D = M$

Input

$x_1$    $x_2$    $\cdots$    $x_M$

*Q: How many hidden units, D, should we use?*

Output

y

What method(s) is
this setting similar to?

Hidden Layer

$z_1$   $z_2$   ...   $z_D$

D < M

Input   $x_1$   $x_2$   $x_3$   ...   $x_M$

Q: How many hidden units, D, should we use?

Output

Hidden Layer

Input

$D > M$

What method(s) is
this setting similar to?

*Q: How many layers should we use?*

*Q: How many layers should we use?*

Output — y

Hidden Layer 3 — $c_1$ $c_2$ ... $c_F$

Hidden Layer 2 — $b_1$ $b_2$ ... $b_E$

Hidden Layer 1 — $a_1$ $a_2$ ... $a_D$

Input — $x_1$ $x_2$ $x_3$ ... $x_M$

## Q: How many layers should we use?

- **Theoretical answer:**
  - A neural network with 1 hidden layer is a **universal function approximator**
  - Cybenko (1989): For any continuous function g($x$), there exists a 1-hidden-layer neural net h$_\theta$($x$)
    s.t. | h$_\theta$($x$) – g($x$) | < ε for all $x$, assuming sigmoid activation functions

- **Empirical answer:**
  - Before 2015 : "Deep networks (e.g. 3 or more hidden layers) are too hard to train"
  - After 2015: "Deep networks are easier to train than shallow networks (e.g. 2 or fewer layers) for many problems"

    Big caveat: You need to know and use the right tricks.

# Different Levels of Abstraction

- We don't know the "right" levels of abstraction
- So let the model figure it out!



Feature representation

3rd layer "Objects"

2nd layer "Object parts"

1st layer "Edges"

Pixels

**Face Recognition:**

– Deep Network can build up increasingly higher levels of abstraction

– Lines, parts, regions



Feature representation

3rd layer "Objects"

2nd layer "Object parts"

1st layer "Edges"

Pixels

# Different Levels of Abstraction

# Activation Functions

Neural Network with sigmoid
activation functions



Output — y

Hidden Layer — $z_1$ $z_2$ ... $z_D$

Input — $x_1$ $x_2$ $x_3$ ... $x_M$

(F) **Loss**
$$J = \tfrac{1}{2}(y - y^*)^2$$

(E) **Output (sigmoid)**
$$y = \frac{1}{1+\exp(-b)}$$

(D) **Output (linear)**
$$b = \sum_{j=0}^{D} \beta_j z_j$$

(C) **Hidden (sigmoid)**
$$z_j = \frac{1}{1+\exp(-a_j)}, \quad \forall j$$

(B) **Hidden (linear)**
$$a_j = \sum_{i=0}^{M} \alpha_{ji} x_i, \quad \forall j$$

(A) **Input**
Given $x_i, \quad \forall i$

# Activation Functions

Neural Network with arbitrary nonlinear activation functions

Output

Hidden Layer

Input

**(F) Loss**
$$J = \tfrac{1}{2}(y - y^*)^2$$

**(E) Output (nonlinear)**
$$y = \sigma(b)$$

**(D) Output (linear)**
$$b = \sum_{j=0}^{D} \beta_j z_j$$

**(C) Hidden (nonlinear)**
$$z_j = \sigma(a_j), \ \forall j$$

**(B) Hidden (linear)**
$$a_j = \sum_{i=0}^{M} \alpha_{ji} x_i, \ \forall j$$

**(A) Input**
Given $x_i, \ \forall i$

## Sigmoid / Logistic Function

$$logistic(u) \equiv \frac{1}{1+e^{-u}}$$



Activation function (nonlinearity) is sigmoid function…

- A new change: modifying the nonlinearity
  - The logistic is not widely used in modern ANNs



Alternate 1:
tanh

Like logistic function but shifted to range [-1, +1]



Slide from William Cohen

# Understanding the difficulty of training deep feedforward neural networks

AI Stats 2010

# Decision Functions

## Neural Network for Classification

**Output** — y

**Hidden Layer** — $z_1$, $z_2$, ..., $z_D$

**Input** — $x_1$, $x_2$, $x_3$, ..., $x_M$

(E) **Output (sigmoid)**
$$y = \frac{1}{1+\exp(-b)}$$

(D) **Output (linear)**
$$b = \sum_{j=0}^{D} \beta_j z_j$$

(C) **Hidden (sigmoid)**
$$z_j = \frac{1}{1+\exp(-a_j)}, \quad \forall j$$

(B) **Hidden (linear)**
$$a_j = \sum_{i=0}^{M} \alpha_{ji} x_i, \quad \forall j$$

(A) **Input**
Given $x_i$, $\forall i$

# NN on Weka

# SVM – Support Vector Machines



Small Margin

Large Margin

Support Vectors

# Support vector machine(SVM).

- Classification is essentially finding the best boundary between classes.

- Support vector machine finds the best boundary points called support vectors and build classifier on top of them.

- Linear and Non-linear support vector machine.

The dots with shadow around them are support vectors. Clearly they are the best data points to represent the boundary. The curve is the separating boundary.

- In this case, class 1 and class 2 are separable.
- The representing points are selected such that the margin between two classes are maximized.
- Crossed points are support vectors.

$$x^T \beta + \beta_0 = 0$$

Example of Non-linear SVM

This classification problem clearly do not have a good optimal linear classifier.

Can we do better?
A non-linear boundary as shown will do fine.

# General SVM Cont.

- The idea is to map the feature space into a much bigger space so that the boundary is linear in the new space.

- Generally linear boundaries in the enlarged space achieve better training-class separation, and it translates to non-linear boundaries in the original space.

# SVM on Veka

# SVM on Veka

# SVM on Veka

# The Problem of Feature Selection

- Large number of features; sometimes greater than 100.

- The number of combinations can be well over a billion!



- Is there a way to search for an optimal set of features in reasonable time and with reasonable computation power?
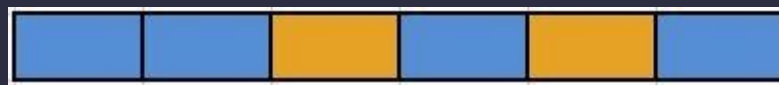
# Different ways to search for this needle

- Evaluate every possible combination to come up with the best combination – the brute force method!

- Step-up/step-down methods that add or remove a feature at a time and evaluate model performance.

- Use genetic algorithms (GA) for searching this huge solution space.
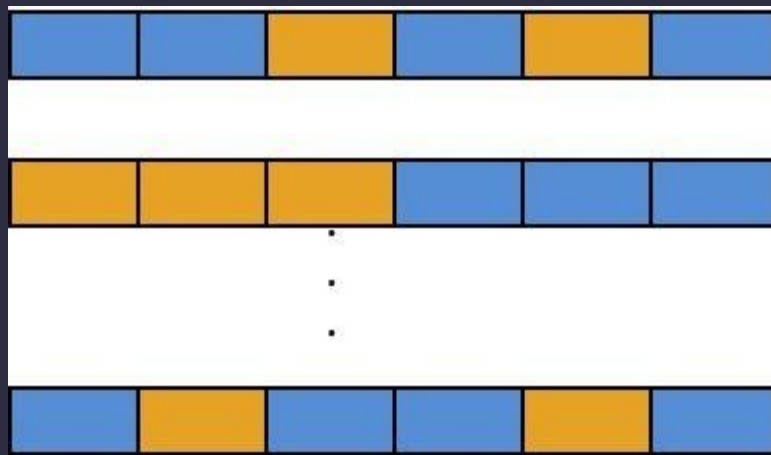
# Genetic Algorithms (GA)

- This is a high level simulation of a biologically inspired adaptive system – evolution.

- Using a simple set of rules, this system can have emergent behaviour that makes it useful for various applications.

- GA have been used in applications such as

  • predicting the structure of proteins

  • training neural networks

- Here, I will talk about the use of GA for searching through the feature space to select an optimal set of features.

- **Chromosome** – a potential solution to the problem. A common way to represent solutions is using binary numbers.

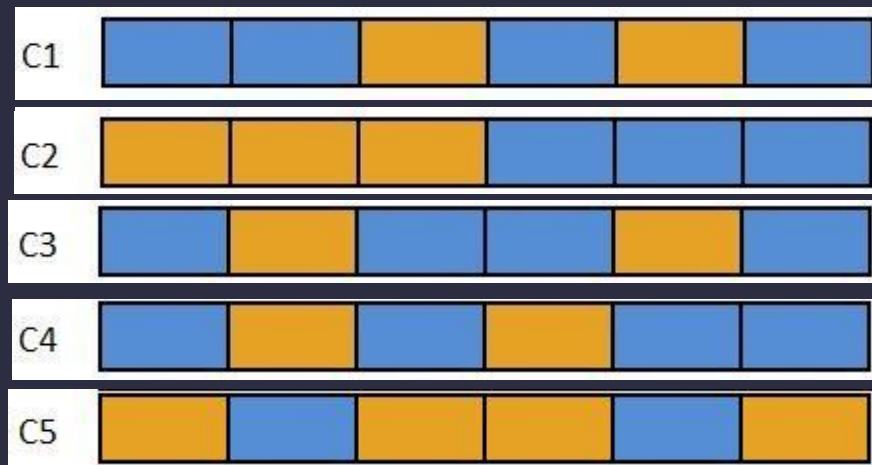

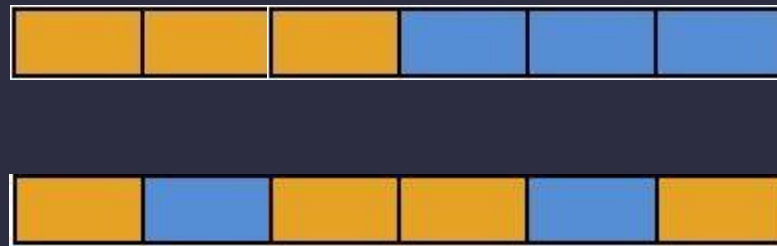- **Population** – a set of chromosomes belonging to a generation.

- **Fitness** – a metric to evaluate how well a particular solution solves the problem.

- **Generation** – each iteration of the algorithm.

- **Selection** – a process by which some chromosomes of a population are chosen for generating new solutions.
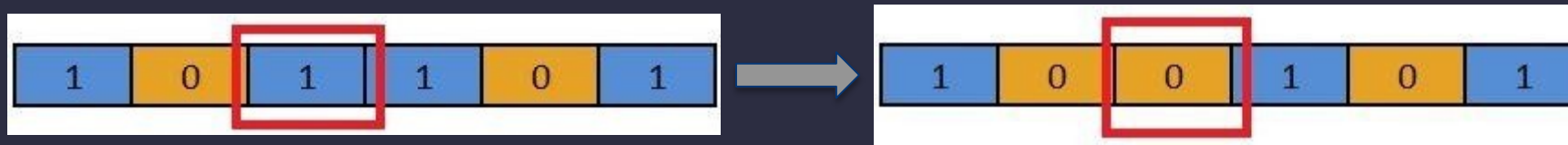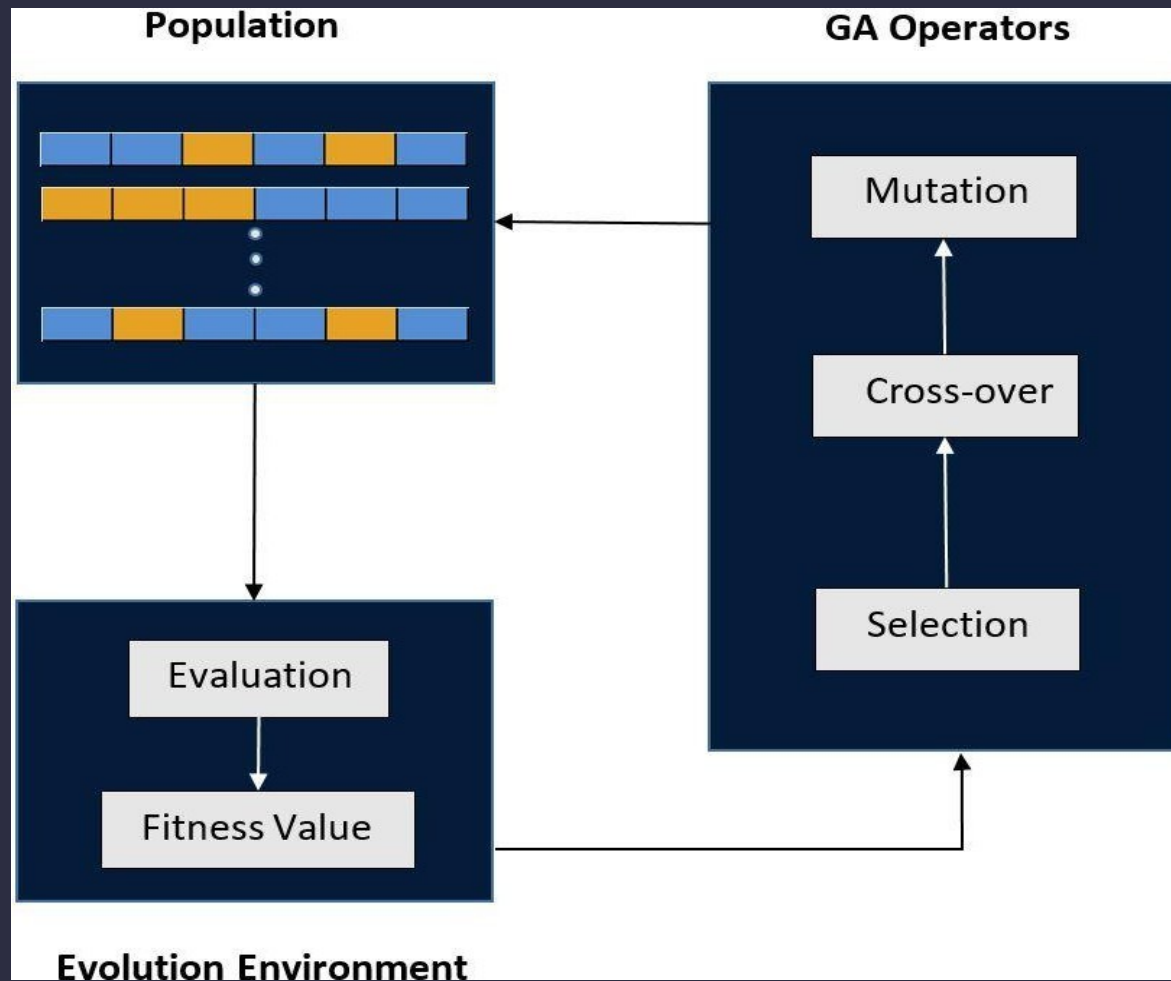
- **Cross-over** – is the process of exchange of information between selected chromosomes.



- **Mutation** – random changes in chromosomes.

# Performance