

# Statistical Methods in AI (CS7.403)

Lecture-4: Performance Measures, Benchmarking ; k-Nearest Neighbours

Ravi Kiran ([ravi.kiran@iiit.ac.in](mailto:ravi.kiran@iiit.ac.in))

<https://ravika.github.io>

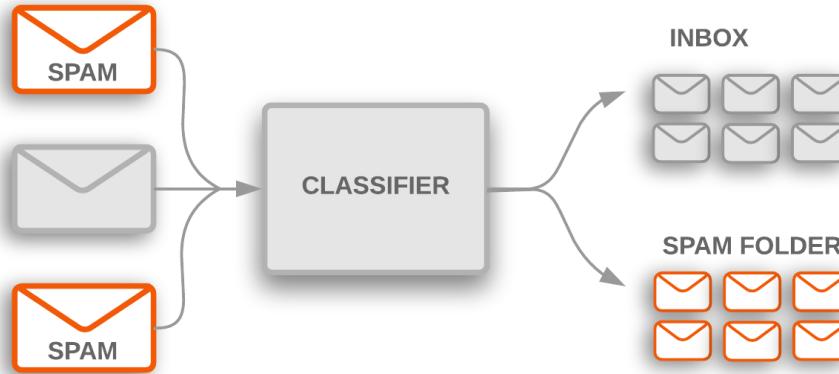


Center for Visual Information Technology (CVIT)  
IIIT Hyderabad

# Announcements

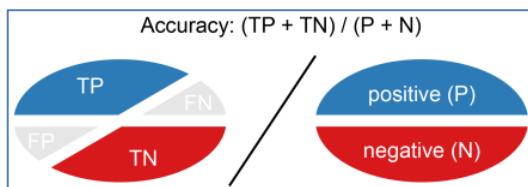
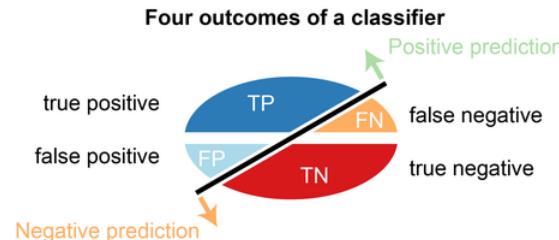
- For those who participated on Wednesday: Thanks !
- A1
  - Will be out today
  - Explanation session tomorrow in Tutorial
- Tutorial: Proof techniques

# Binary Classification

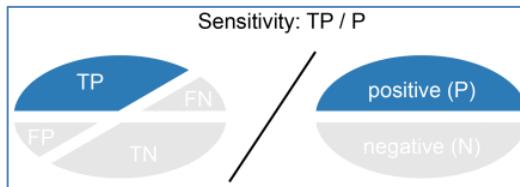


# Accuracy vs Precision vs Recall

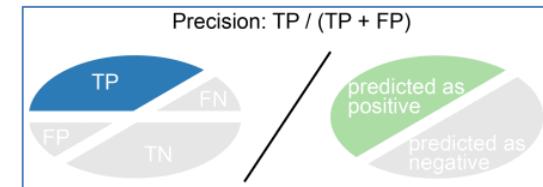
- Accuracy : Performance w.r.t both classes
- Recall : Performance w.r.t '+' class
- Precision : Reliability of predictions w.r.t '+' class



% of correct predictions



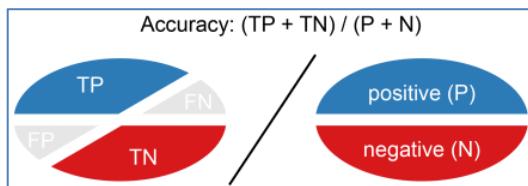
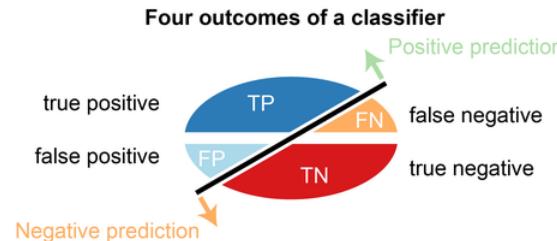
% of + class correctly predicted  
[aka Recall / TPR]



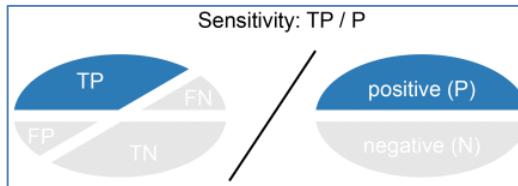
correct prediction of + class  
[aka Precision]

# Accuracy vs Precision vs Recall

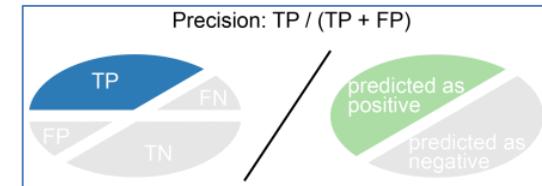
- Monitor **Precision** if a **false positive** carries higher cost.
- Monitor **Recall** if a **false negative** carries higher cost.



% of correct predictions



% of + class correctly predicted  
[aka Recall / TPR]



correct prediction of + class  
[aka Precision]

Precision, Recall, F1 score

## Classification

Binary

$\{0,1\}$

Multi-class

1-of-K

Multi-label

n-of-K

Structure

# Multi-class classification

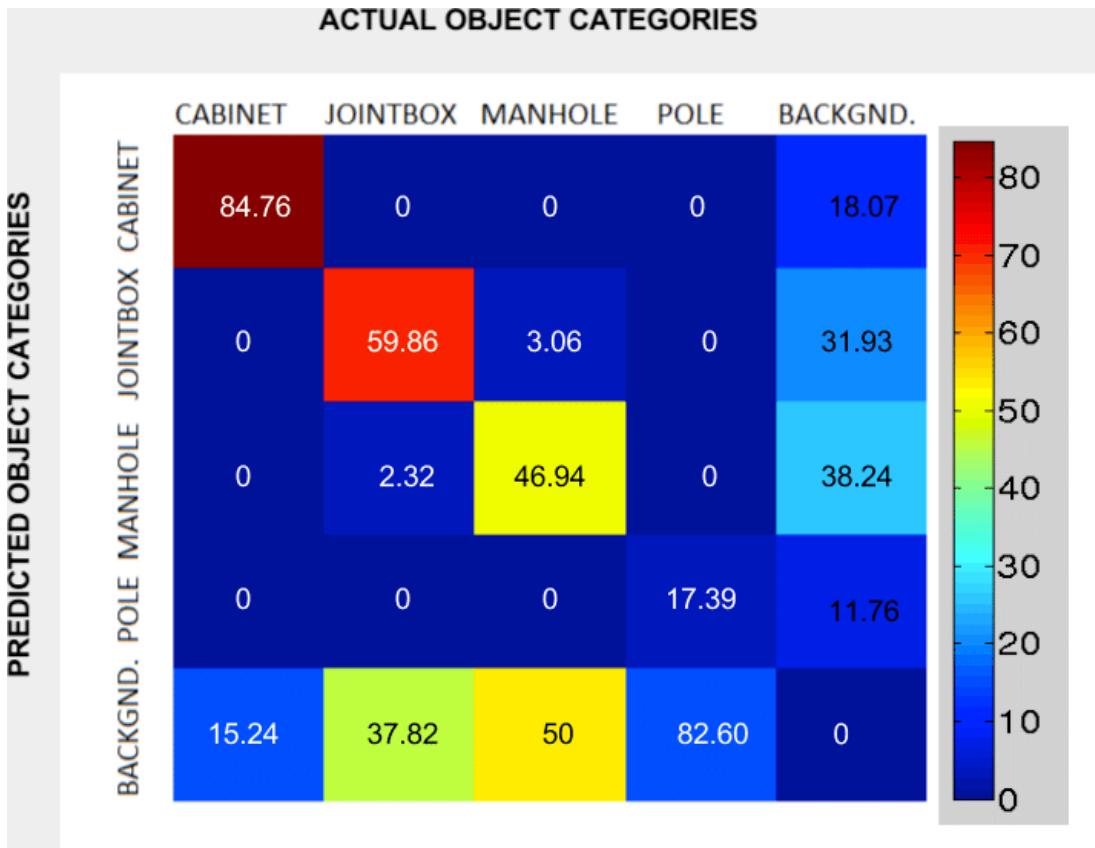
## Multi-class classification

# Multi-class Classification - Confusion matrix

| activity recognition from video |       |     |     |    |     |    |     |     |     |     |   |
|---------------------------------|-------|-----|-----|----|-----|----|-----|-----|-----|-----|---|
| actual class                    | bend  | 100 | 0   | 0  | 0   | 0  | 0   | 0   | 0   | 0   | 0 |
|                                 | jack  | 0   | 100 | 0  | 0   | 0  | 0   | 0   | 0   | 0   | 0 |
|                                 | jump  | 0   | 0   | 89 | 0   | 0  | 0   | 11  | 0   | 0   | 0 |
|                                 | pjump | 0   | 0   | 0  | 100 | 0  | 0   | 0   | 0   | 0   | 0 |
|                                 | run   | 0   | 0   | 0  | 0   | 89 | 0   | 11  | 0   | 0   | 0 |
|                                 | side  | 0   | 0   | 0  | 0   | 0  | 100 | 0   | 0   | 0   | 0 |
|                                 | skip  | 0   | 0   | 0  | 0   | 0  | 0   | 100 | 0   | 0   | 0 |
|                                 | walk  | 0   | 0   | 0  | 0   | 0  | 0   | 0   | 100 | 0   | 0 |
|                                 | wave1 | 0   | 0   | 0  | 0   | 0  | 0   | 0   | 67  | 33  |   |
|                                 | wave2 | 0   | 0   | 0  | 0   | 0  | 0   | 0   | 0   | 100 |   |
| predicted class                 |       |     |     |    |     |    |     |     |     |     |   |

- Reveals several biases of the classifier:
  - Most confusing pairs
  - Least / Most likely classes

# Multi-class Classification - Confusion matrix



- Reveals several biases of the classifier:
  - Most confusing pairs
  - Least / Most likely classes

# Multi-class Classification: Measures

- Mean <measure> +- standard deviation
- Median <measure> +- median absolute deviation

| Descriptor     | Spectral bands   |                  |     |
|----------------|------------------|------------------|-----|
|                | RGB              | PCA              | RGB |
| Gist           | $74.14 \pm 1.93$ | $77.76 \pm 2.62$ |     |
| MSIFT          | $88.92 \pm 1.39$ | $90.97 \pm 1.81$ |     |
| MBoW           | $88.60 \pm 1.70$ | $88.31 \pm 1.38$ |     |
| cSIFT          | $88.17 \pm 1.17$ | $88.76 \pm 1.74$ |     |
| rgSIFT         | $88.24 \pm 1.89$ | $87.71 \pm 1.33$ |     |
| BoWV [8]       | 71.86            | N/A              |     |
| SPMK [12]      | 74.00            | N/A              |     |
| SPCK++ [8]     | 76.05            | N/A              |     |
| Dense SIFT [2] | $81.67 \pm 1.23$ | N/A              |     |

# (Sanity Check) Baselines

- 0 cost-to-build classifiers
- Binary: Assume equal # of samples / class
  - What is random guessing accuracy ?

# (Sanity Check) Baselines

- 0 cost-to-build classifiers
- Binary: Assume unequal # of samples / class [n=#samples, x= % of + class ( $\neq 50\%$ ) , - class is majority in training set]
  - Classifier: Predict majority class for all samples.
  - Accuracy ?

# (Sanity Check) Baselines

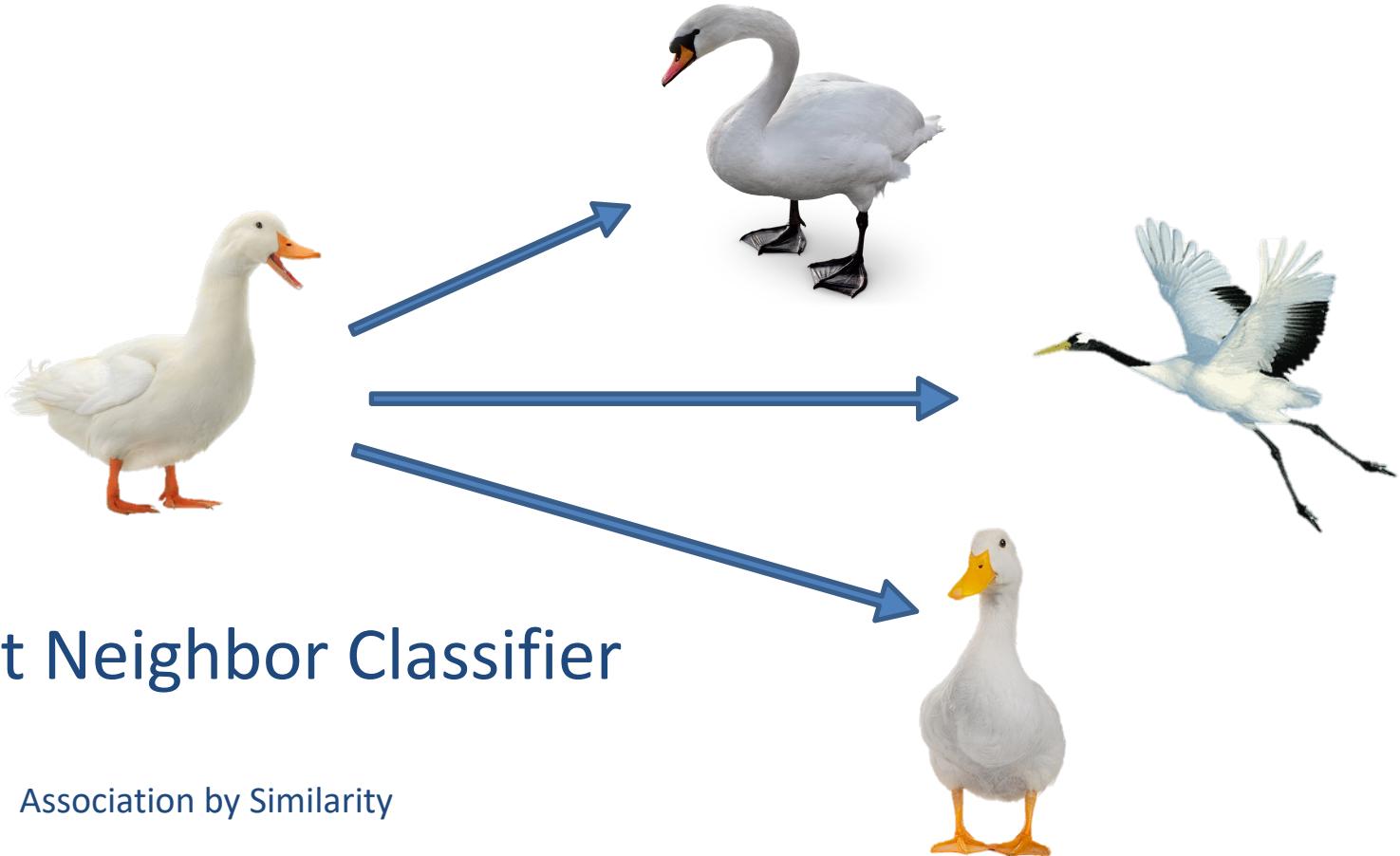
- Multi-class: Assume equal # of samples / class [k classes, N total samples]
  - What is random guessing accuracy ?

# Summary

- Many metrics:
  - Accuracy, TP, FP, Precision, Recall
  - Class imbalance and decision-cost imbalance must be taken into account
- Confusion Matrix: Important to analyze and refine solution.

# References and Reading

- Code
  - [https://scikit-learn.org/stable/modules/model\\_evaluation.html#classification-metrics](https://scikit-learn.org/stable/modules/model_evaluation.html#classification-metrics)



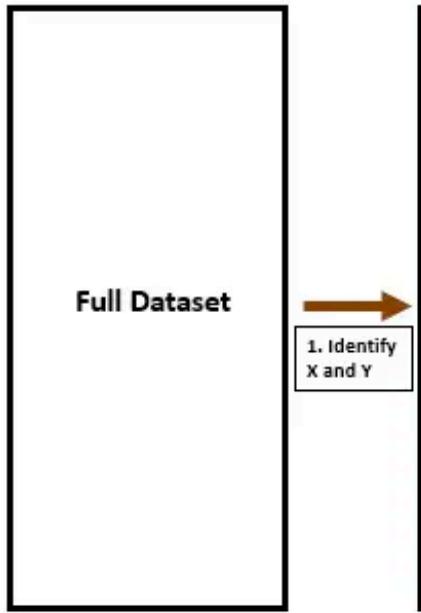
## Nearest Neighbor Classifier

Association by Similarity

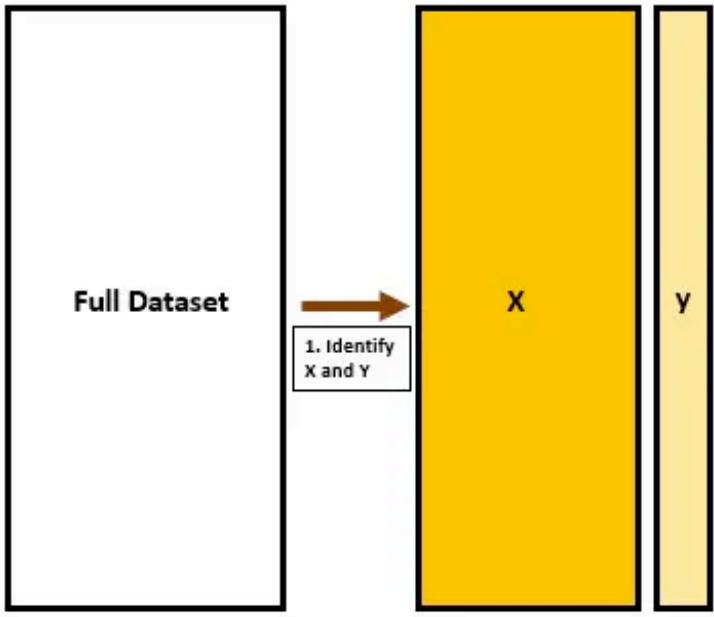
# Goal of ML Classification

- Given:
  - A set of training  $n$  training samples:  $(\mathbf{x}_i, y_i)$
- Correctly predict
  - Labels of **unseen** test samples:  $\mathbf{x}_t$
- Goal:
  - Maximize the **accuracy** on unseen test samples
- How do we find the **accuracy** on unseen samples?

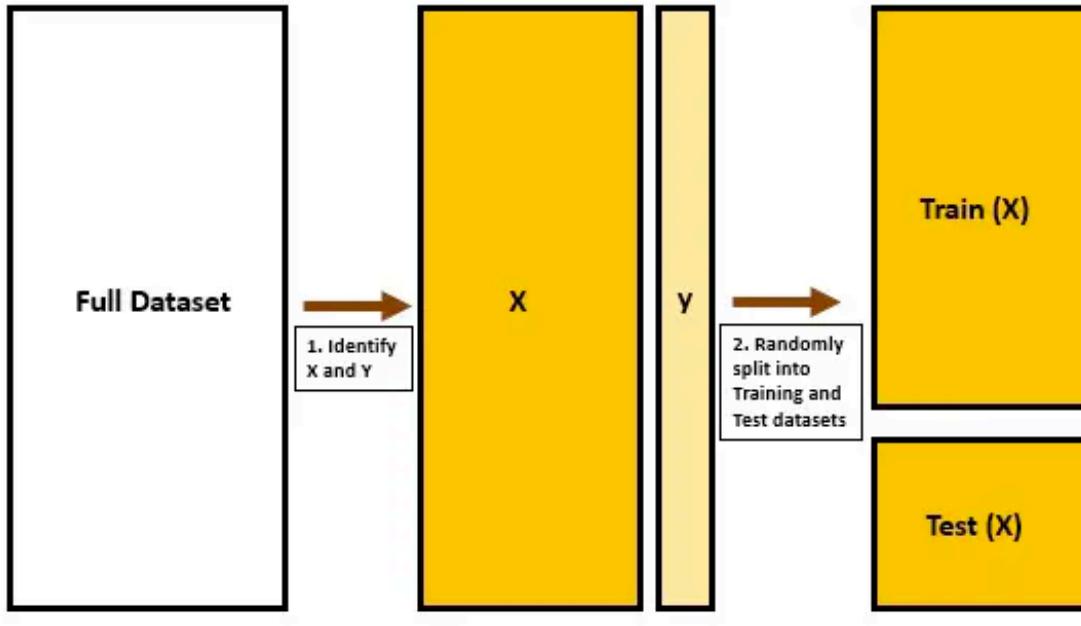
# ML Train-Test Setup



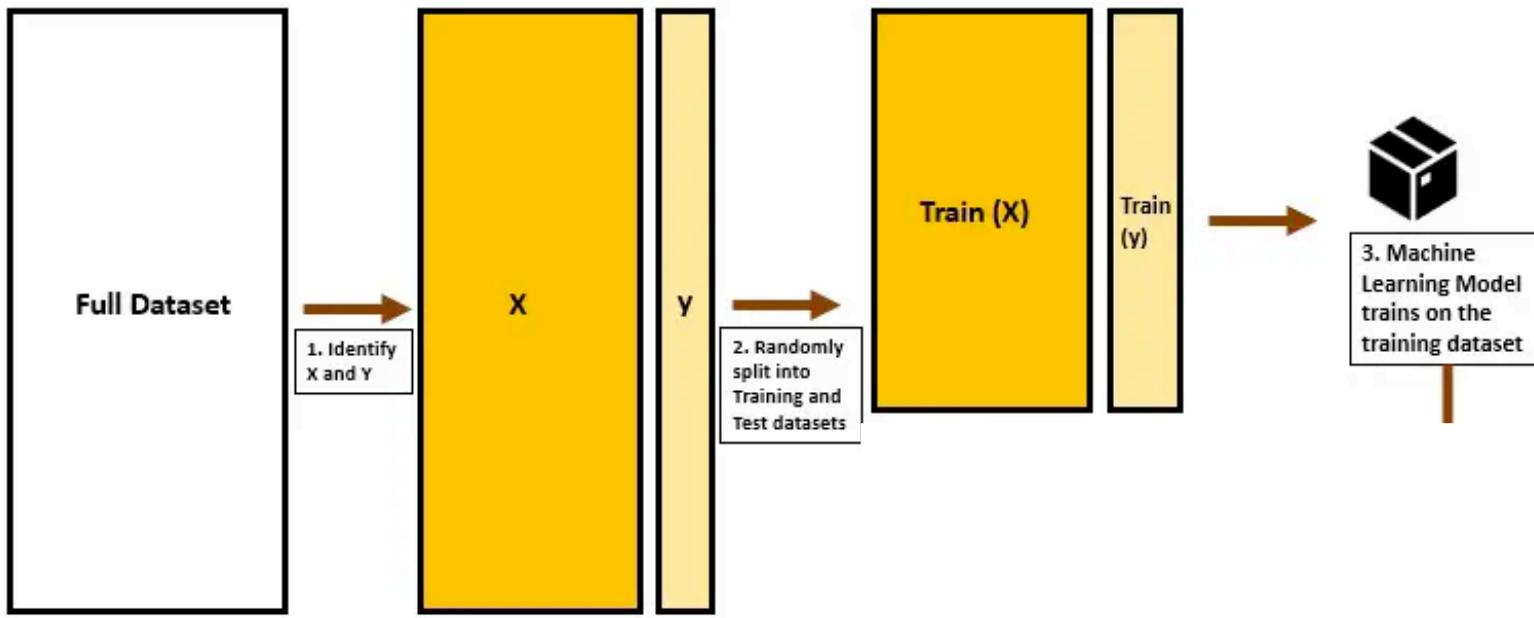
# ML Train-Test Setup



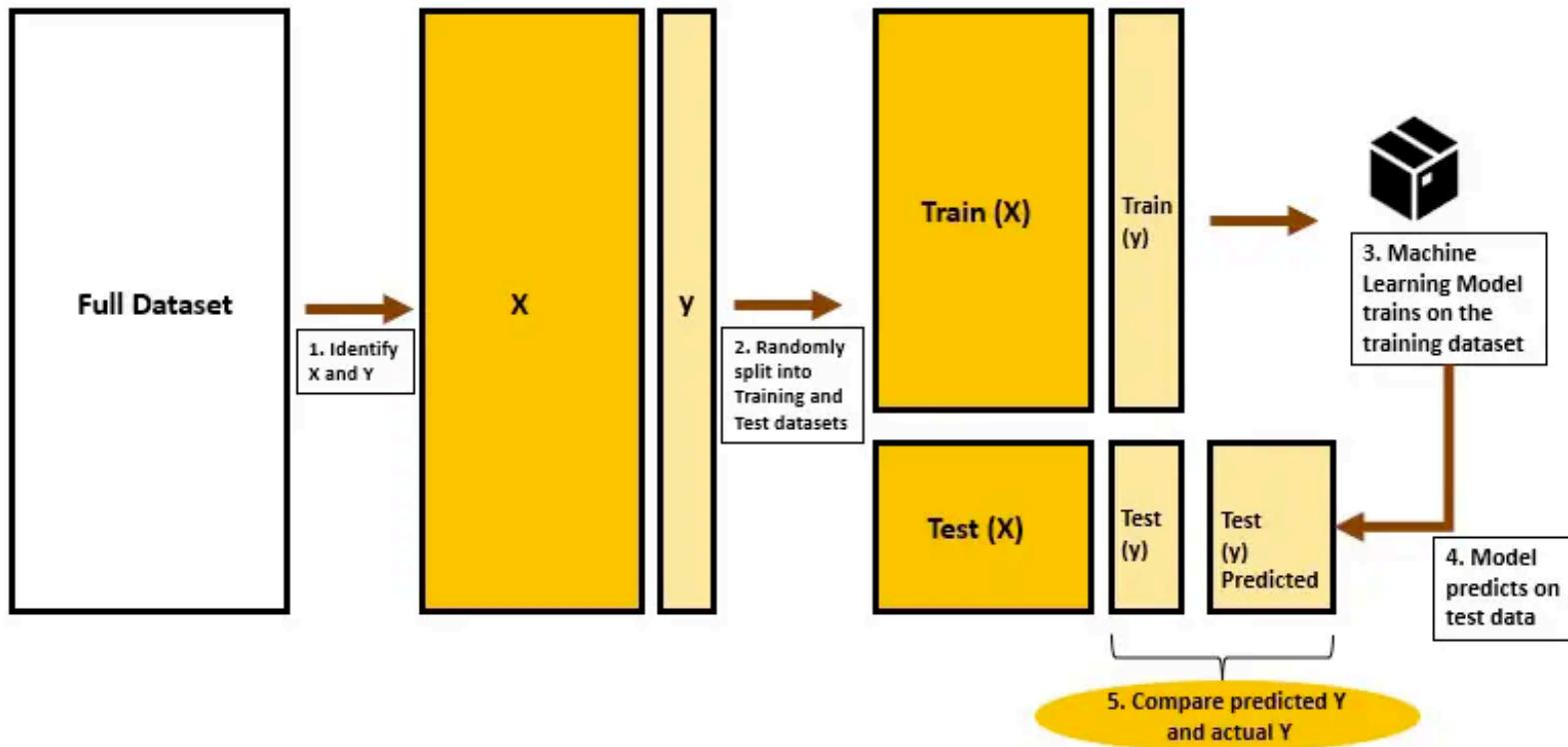
# ML Train-Test Setup



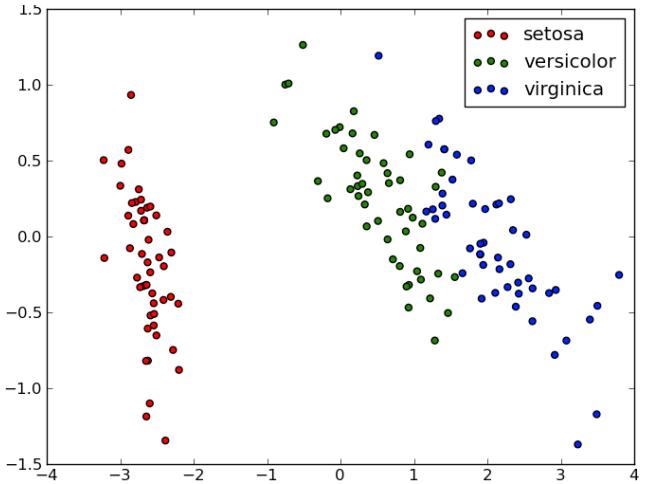
# ML Train-Test Setup



# ML Train-Test Setup

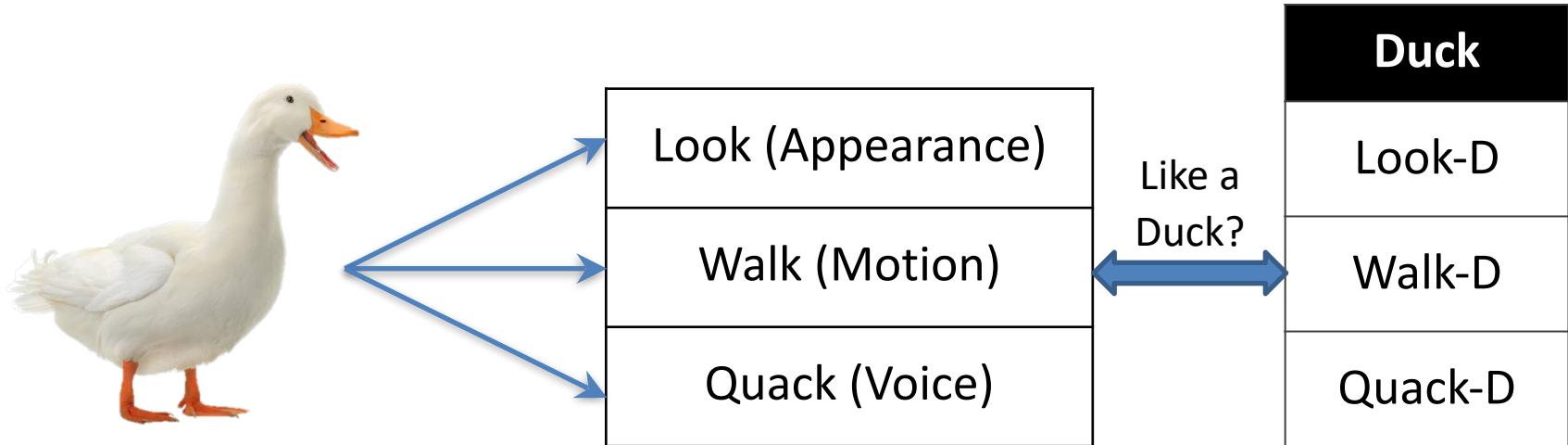


# Data Samples → Vectors



# How do we compare?

If it looks like a duck, walks like a duck and quacks like a duck,



- Find distance to feature vectors of known classes

# Nearest Neighbor Classifier

- Assign label of that sample which is nearest to the test sample

$X_{\text{test}}$  :

[30.9, 15.1, 1.32]

The test sample  
is a **Duck**

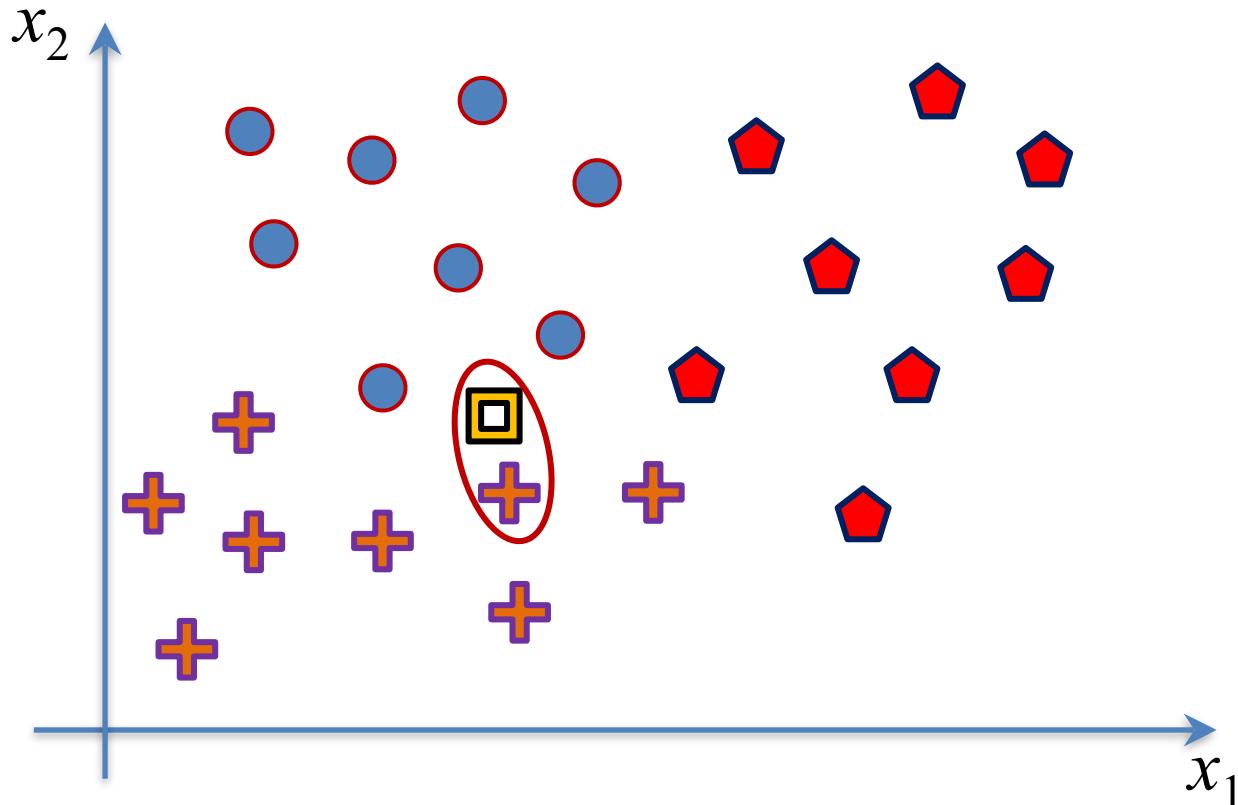
| Distance |
|----------|
| 1.30     |
| 5.78     |
| 13.54    |
| 4.71     |
| 15.21    |
| 3.04     |
| 13.58    |

Training Samples

| Feature Vector           | Label |
|--------------------------|-------|
| $X_1 [32.1, 14.6, 1.42]$ | Duck  |
| $X_2 [25.3, 16.3, 2.11]$ | Swan  |
| $X_3 [42.2, 7.7, 0.38]$  | Crane |
| $X_4 [26.7, 17.1, 2.04]$ | Swan  |
| $X_5 [44.1, 7.6, 0.32]$  | Crane |
| $X_6 [31.4, 12.1, 1.29]$ | Duck  |
| $X_7 [41.9, 7.2, 0.35]$  | Crane |

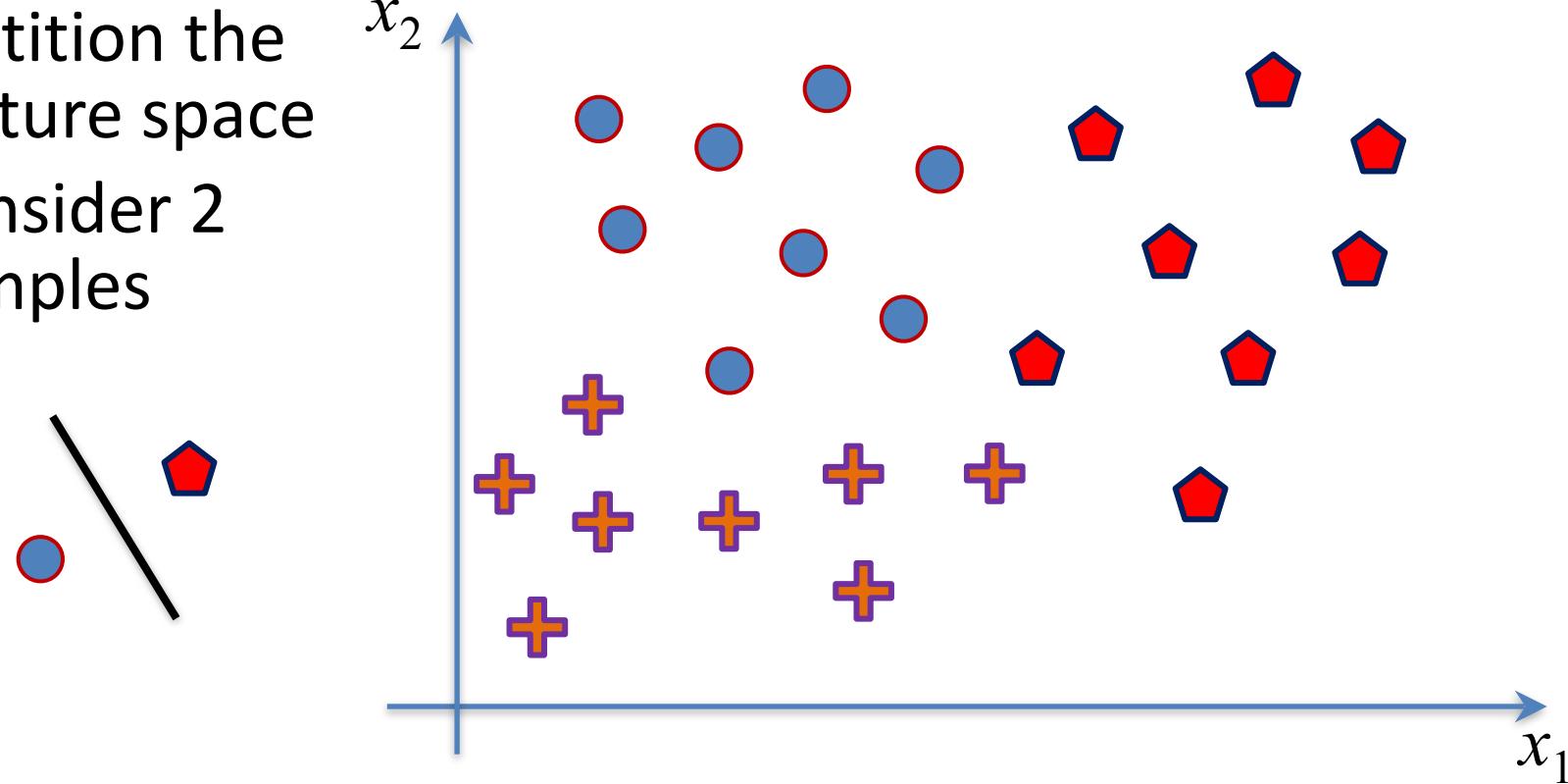
# Visualization in Feature Space

- The nearest train sample is a +.
- We assign the test sample to the + class.



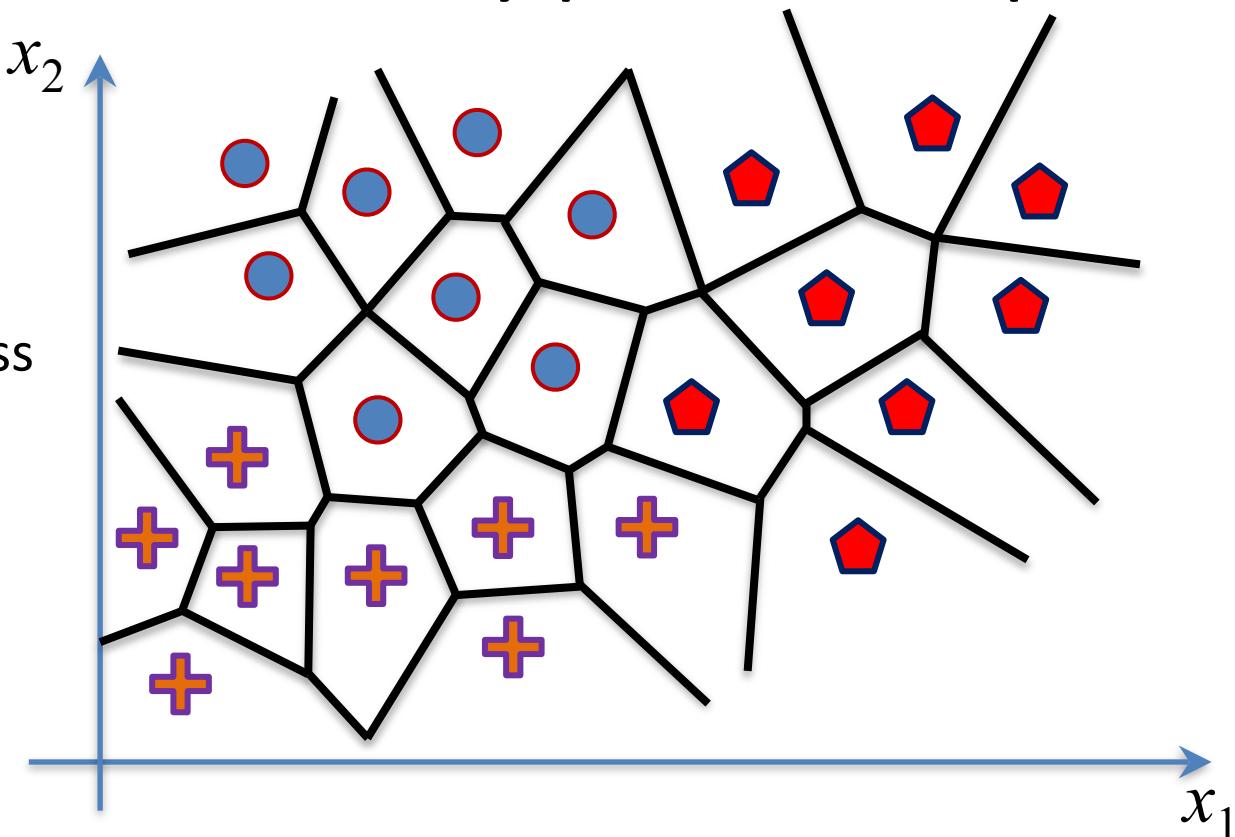
# Class Boundaries

- Partition the feature space
- Consider 2 samples



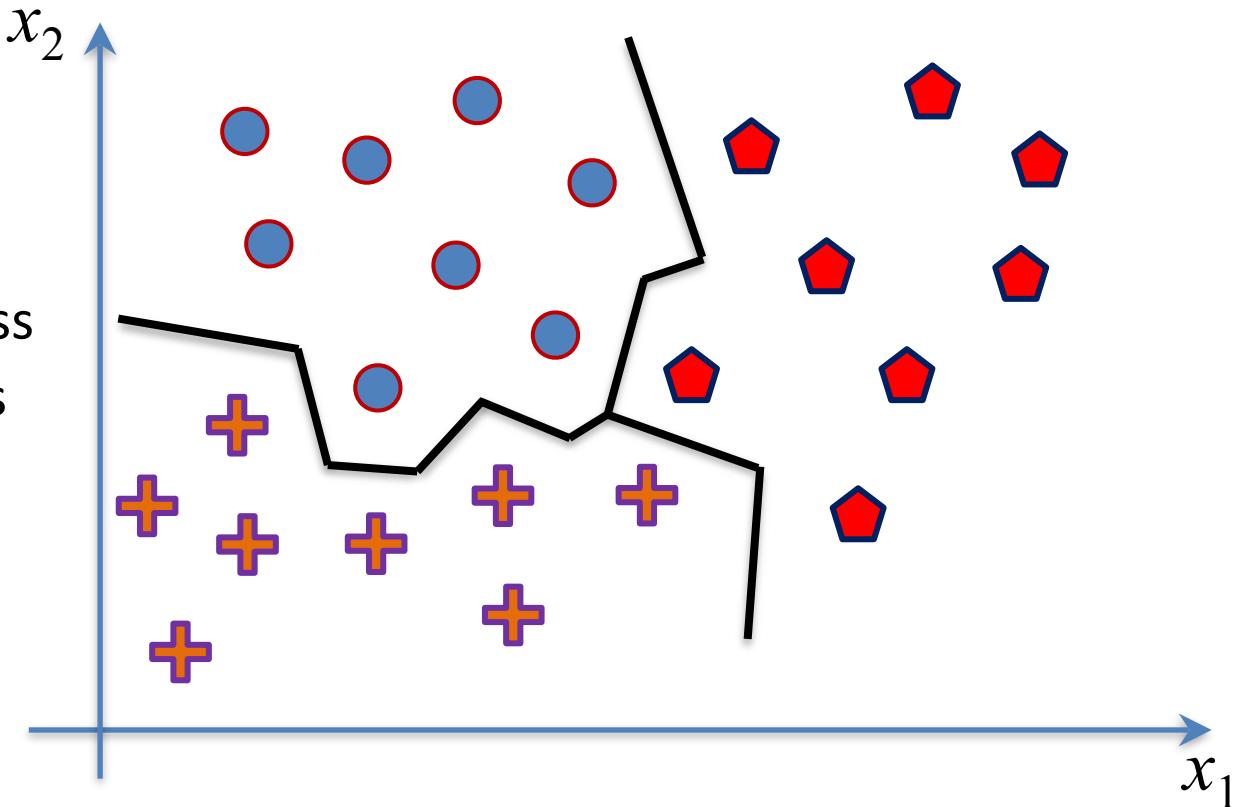
# Boundaries between every pair of samples

- Voronoi Tessellation
- We can ignore boundaries between samples of same class



# Boundaries between every pair of samples

- Voronoi Tessellation
- We can ignore boundaries between samples of same class
- Decision Boundary is piece-wise Linear



# The 1-NN Classification Algorithm

## Problem

- Given:
  - A set of training  $n$  training samples:  $(\mathbf{x}_i, y_i)$
  - A test sample:  $\mathbf{x}_t$
- Find:
  - $\text{label}(\mathbf{x}_t)$

## Algorithm

```
for  $\mathbf{x}_i$  in  $\{\mathbf{X}_{\text{train}}\}$ :  
    dist = Dist( $\mathbf{x}_i$ ,  $\mathbf{x}_t$ )  
    if (dist < minDist):  
        minDist = dist  
        nearest =  $\mathbf{x}_i$   
return ( $y_i$ )
```

# 1-NN

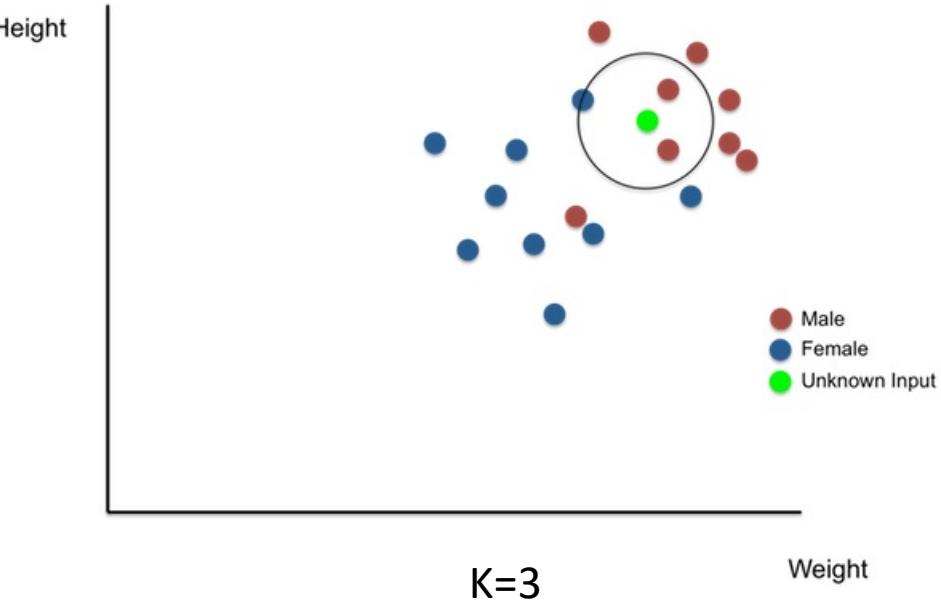
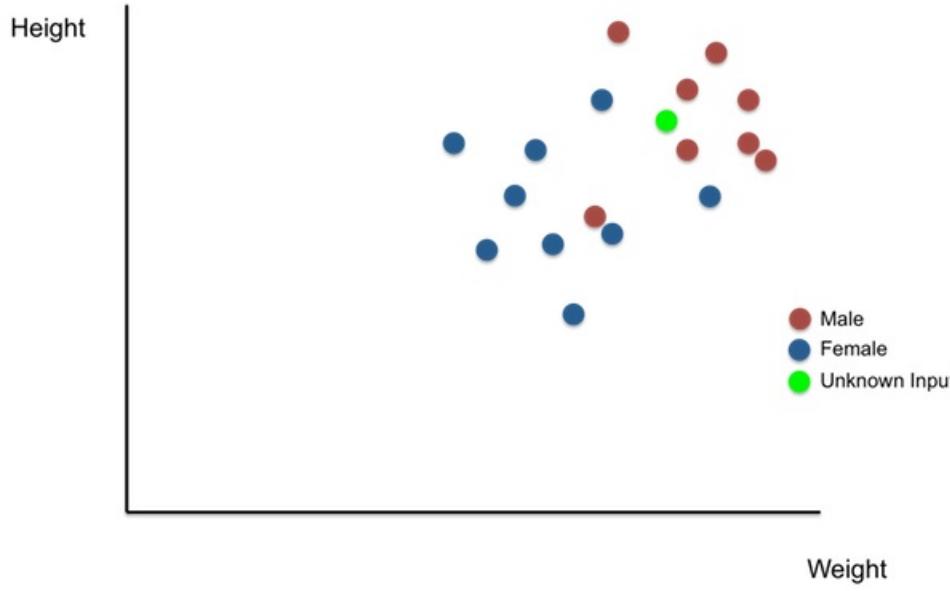
## Algorithm:

1. Find example  $(\mathbf{x}^*, t^*)$  (from the stored training set) closest to the test instance  $\mathbf{x}$ . That is:

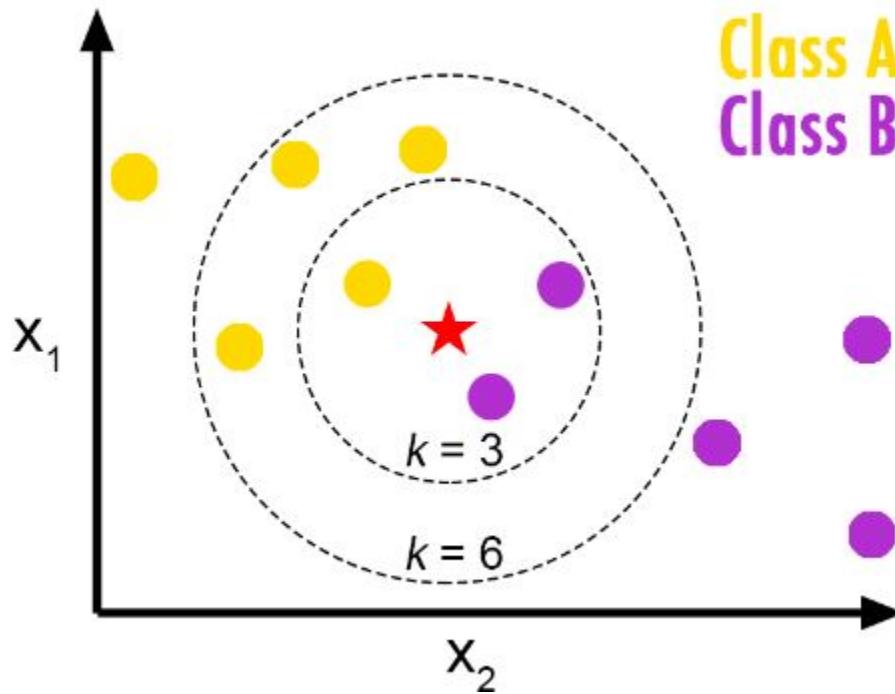
$$\mathbf{x}^* = \operatorname{argmin}_{\mathbf{x}^{(i)} \in \text{train. set}} \text{distance}(\mathbf{x}^{(i)}, \mathbf{x})$$

2. Output  $y = t^*$

# k-nearest neighbor classifier



# k-nearest neighbor classifier



K is usually an odd number

# k-NN

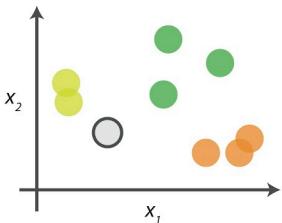
## Algorithm (kNN):

1. Find k examples  $\{\mathbf{x}^{(i)}, t^{(i)}\}$  closest to the test instance  $\mathbf{x}$
2. Classification output is majority class

$$y = \arg \max_{t^{(z)}} \sum_{r=1}^k \delta(t^{(z)}, t^{(r)})$$

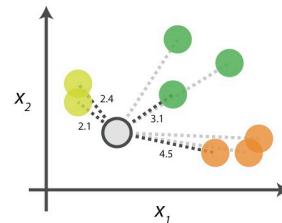
# k-NN algorithm in pictures

## 0. Look at the data



Say you want to classify the grey point into a class. Here, there are three potential classes - lime green, green and orange.

## 1. Calculate distances



Start by calculating the distances between the grey point and all other points.

## 2. Find neighbours

| Point | Distance | Rank   |
|-------|----------|--------|
| ●     | 2.1      | 1st NN |
| ●     | 2.4      | 2nd NN |
| ●     | 3.1      | 3rd NN |
| ●     | 4.5      | 4th NN |

Next, find the nearest neighbours by ranking points by increasing distance. The nearest neighbours (NNs) of the grey point are the ones closest in dataspace.

## 3. Vote on labels

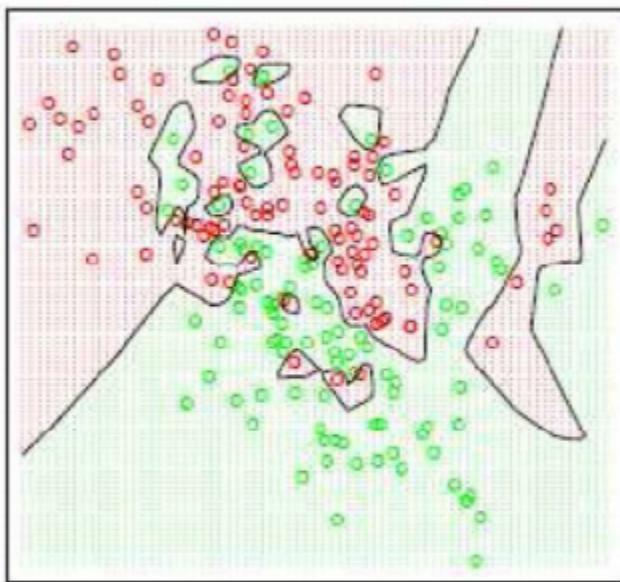
| Class | # of votes |
|-------|------------|
| ●     | 2          |
| ●     | 1          |
| ●     | 1          |

Class ● wins the vote!  
Point ○ is therefore predicted to be of class ●.

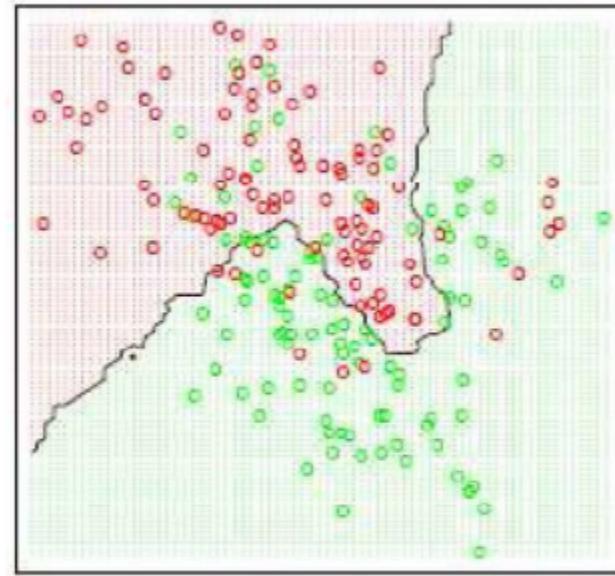
Vote on the predicted class labels based on the classes of the k nearest neighbours. Here, the labels were predicted based on the k=3 nearest neighbours.

# Effect of K

K=1



K=15



Figures from Hastie, Tibshirani and Friedman (Elements of Statistical Learning)

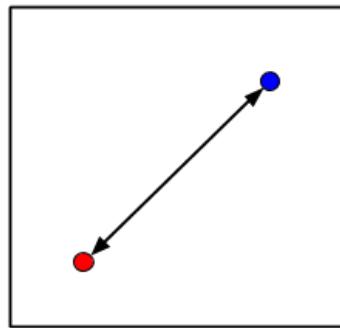
Larger k produces smoother boundary effect and can reduce the impact of class label noise.

# Complexity of k-NN

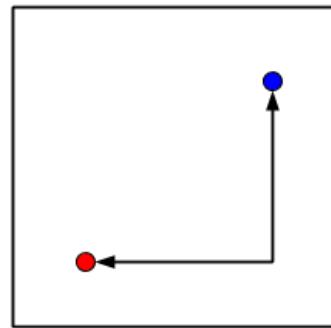
- Training
  - Time:
  - Space:
- Testing
  - Time:
  - Space:

# Distance measures

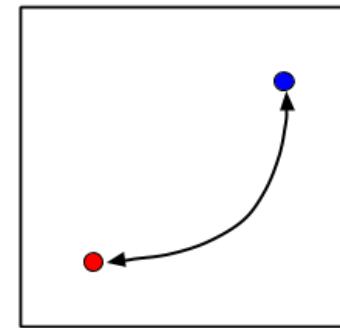
Euclidean



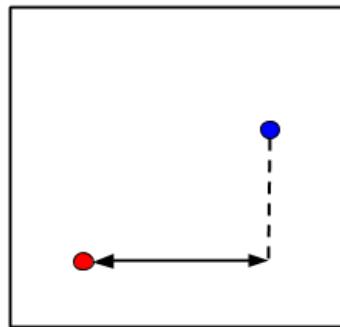
Manhattan



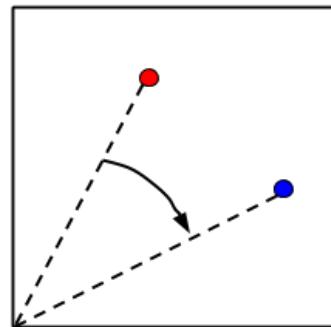
Minkowski



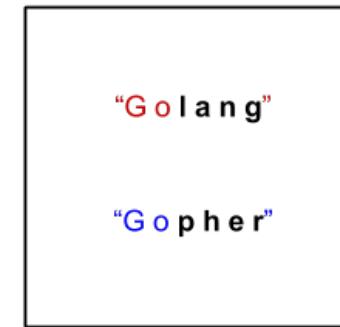
Chebychev



Cosine Similarity



Hamming

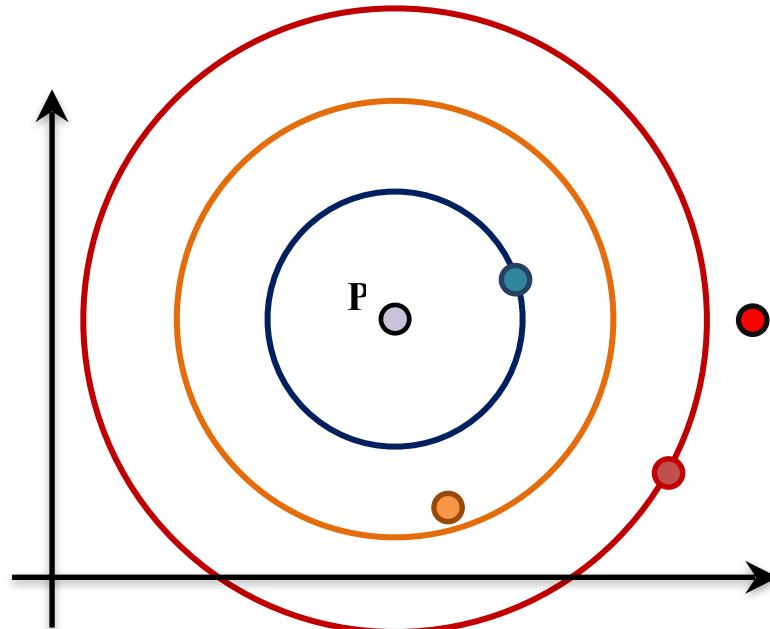
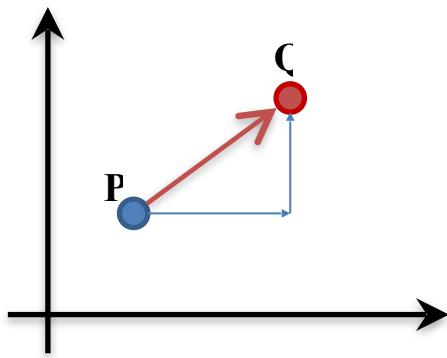


# Euclidean Distance [L2]

Iso-surfaces

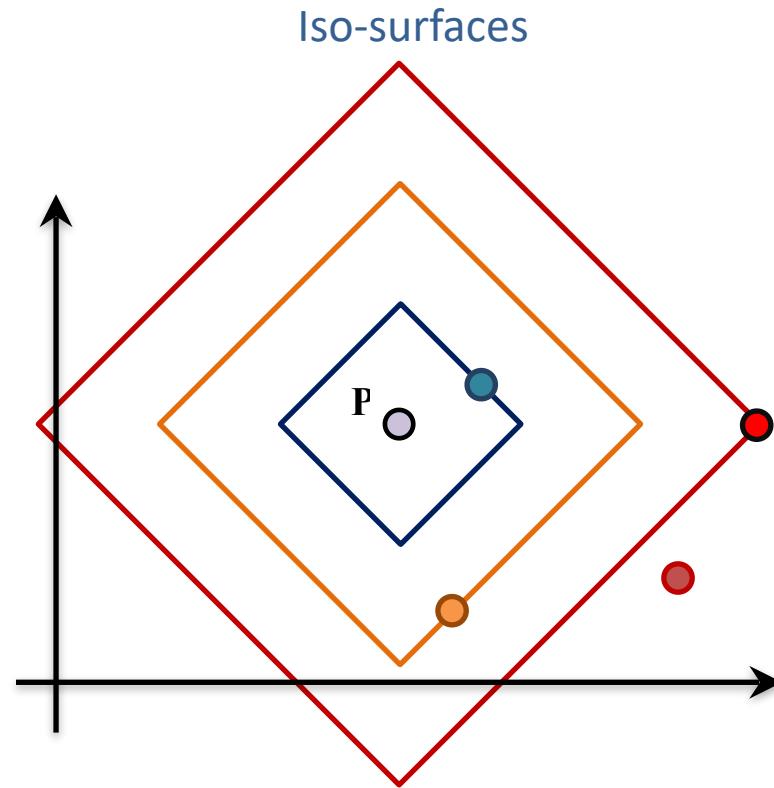
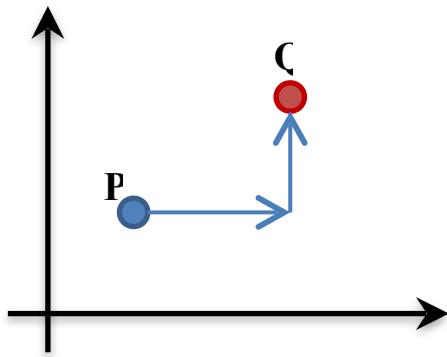
$$d(P, Q)^2 = \sum_{i=1}^d (p_i - q_i)^2$$

$$d(P, Q)^2 = (P - Q)^T (P - Q)$$



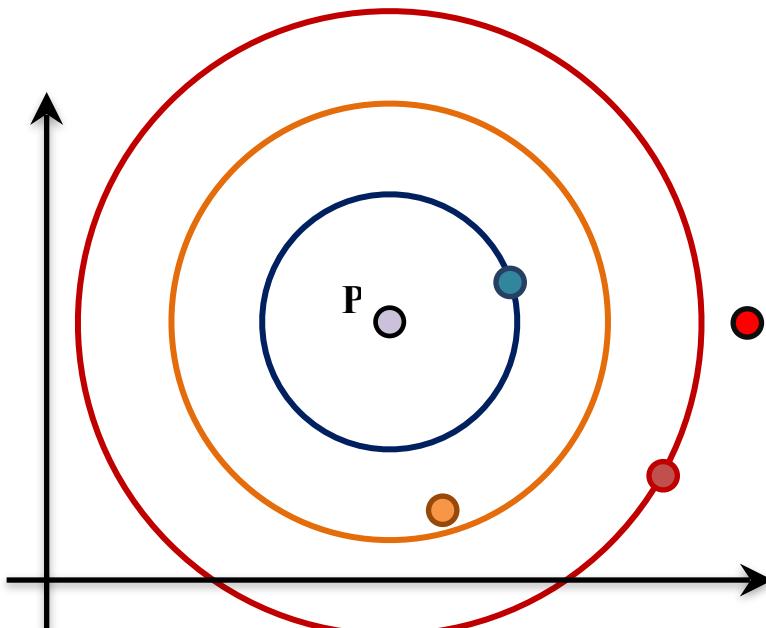
# Manhattan Distance [L1]

$$d(P, Q) = \sum_{i=1}^d |(p_i - q_i)|$$



# Distance measure can affect k-NN classification

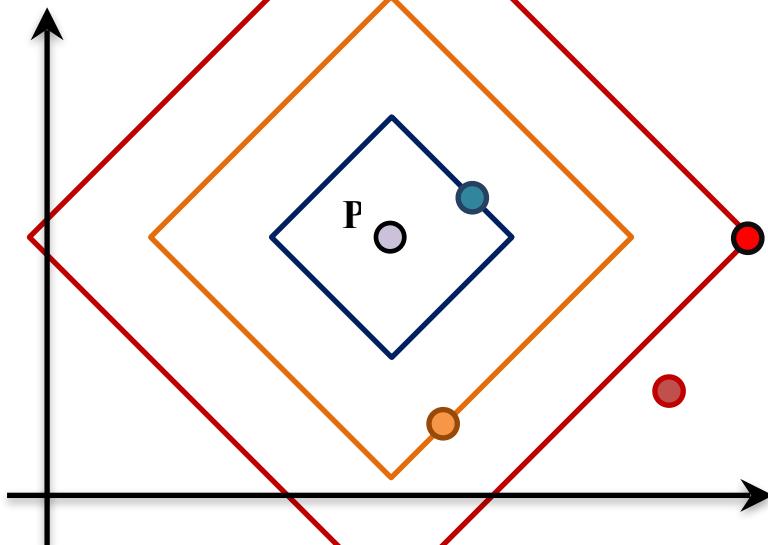
Iso-surfaces



Euclidean

$$d(P, Q)^2 = \sum_{i=1}^d (p_i - q_i)^2$$

Iso-surfaces



Manhattan

$$d(P, Q) = \sum_{i=1}^d |(p_i - q_i)|$$

# Minkowski Dist.

$$d(P, Q)^r = \sum_{i=1}^d |p_i - q_i|^r$$

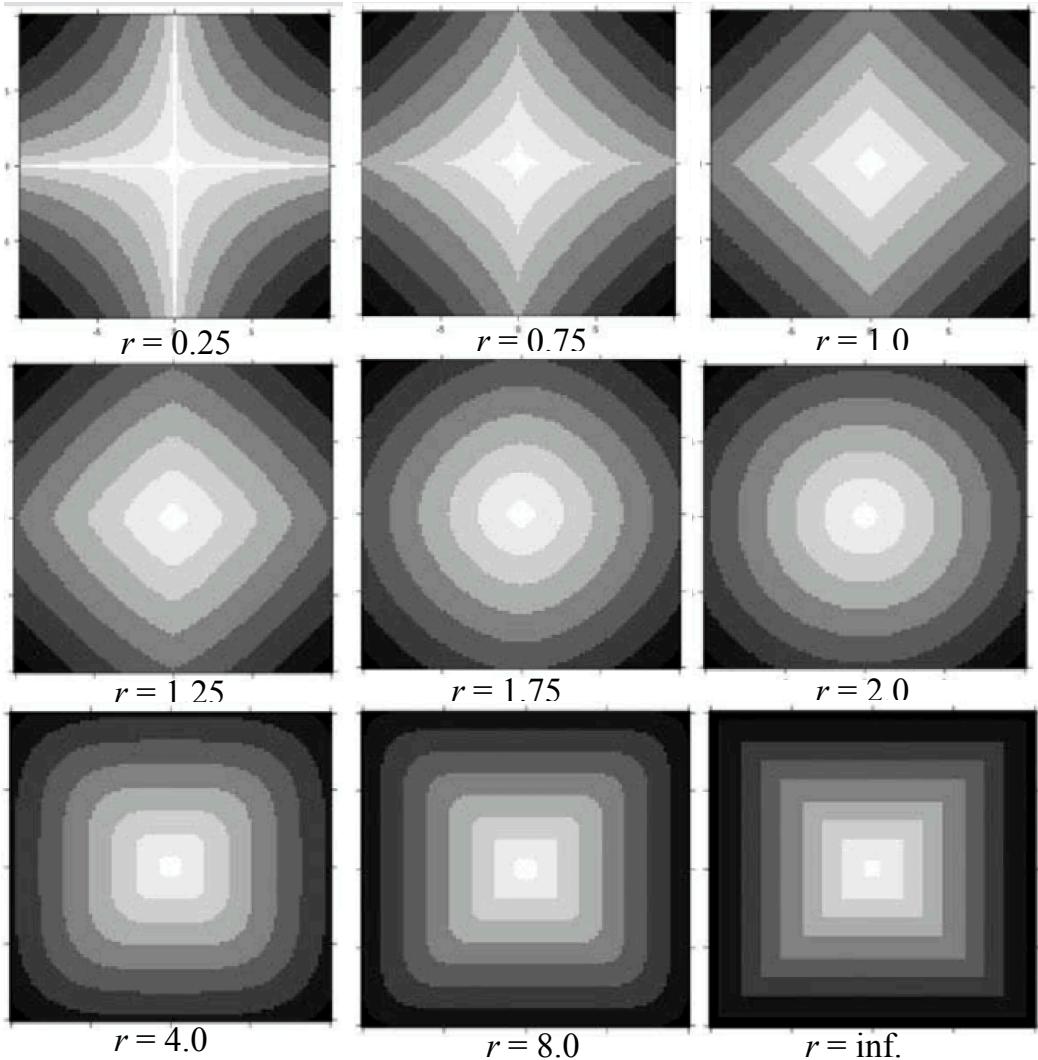
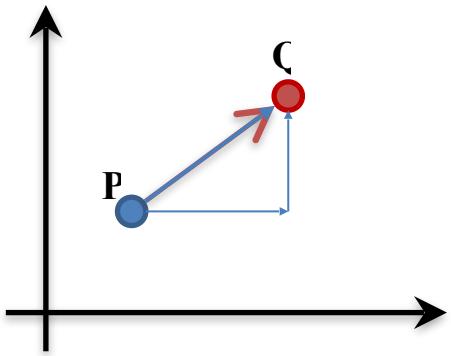
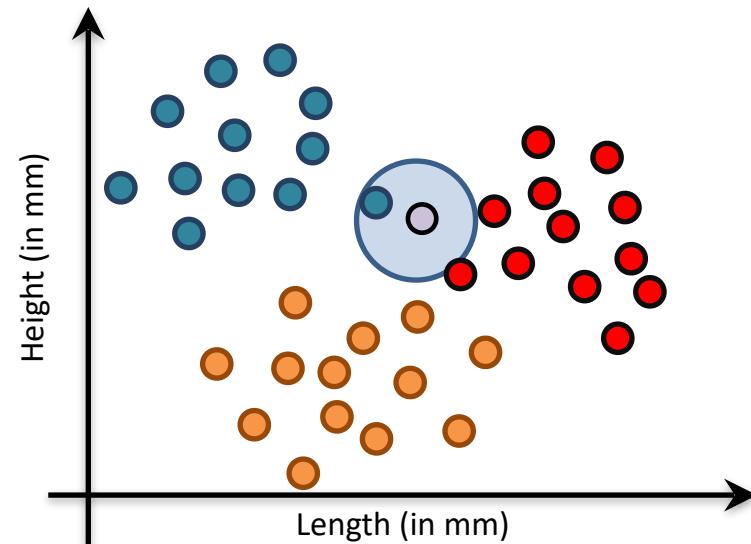
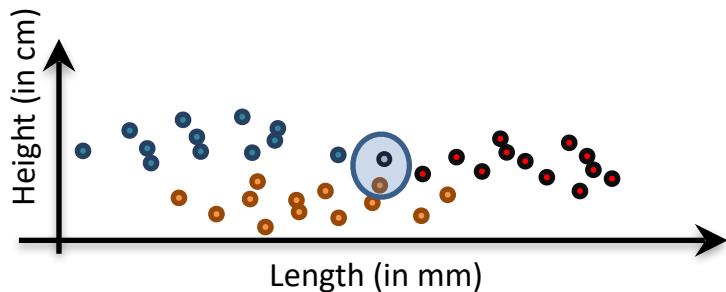


Fig. by Lu et al., "The Minkowski Approach for Choosing the Distance Metric in Geographically Weighted Regression"

# Note: Feature Normalization

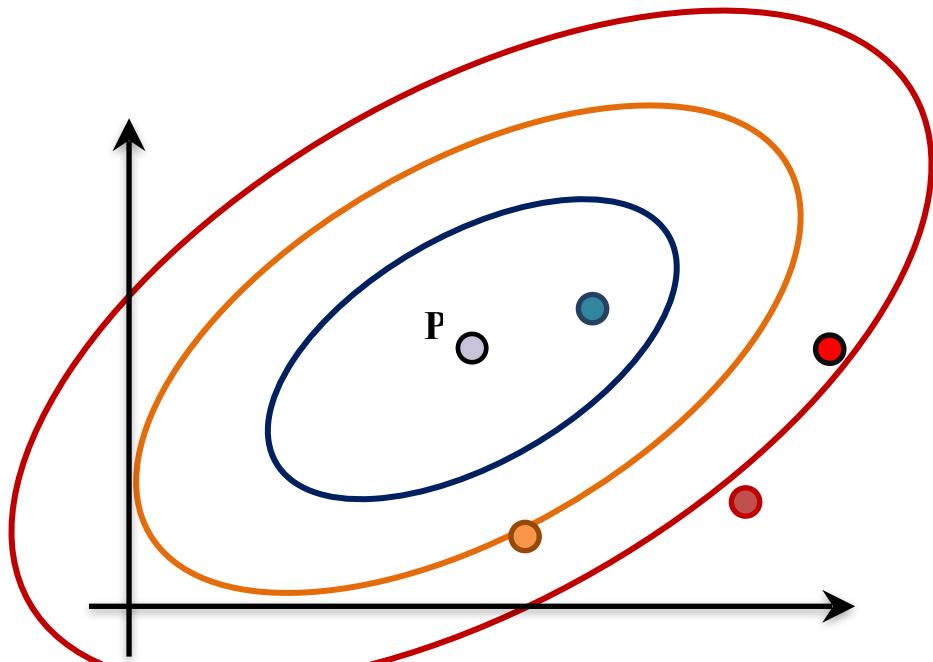
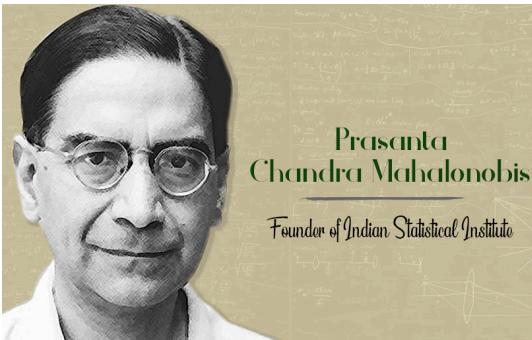
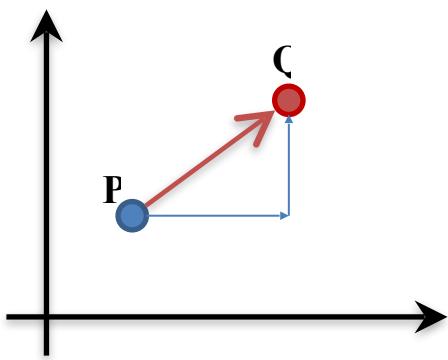
- If different features have different variances
  - Some features will dominate distance computation
  - Normalization can reduce this feature bias



# Mahalanobis Distance

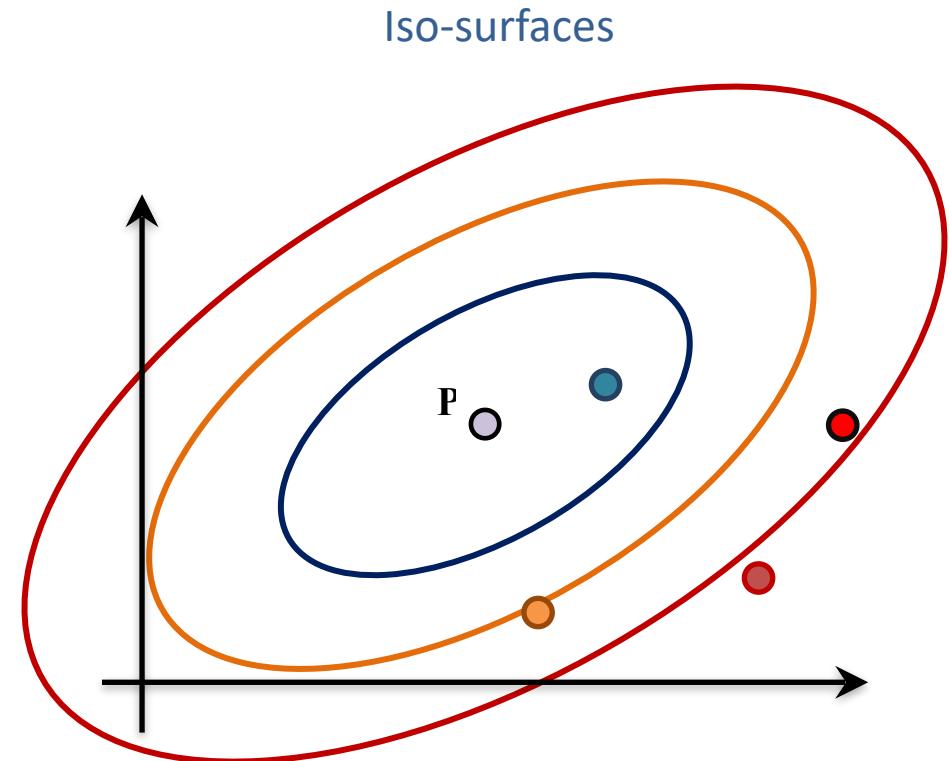
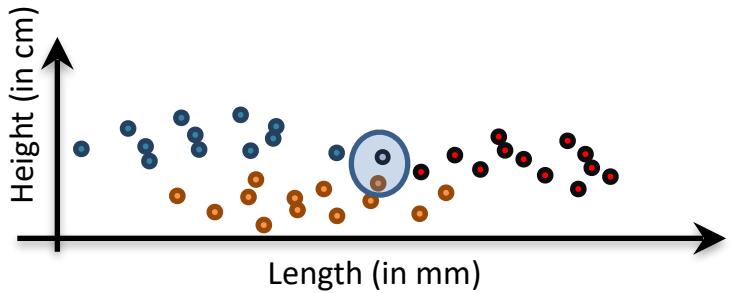
Iso-surfaces

$$d(P, Q)^2 = (P - Q)^T \mathbf{S}^{-1} (P - Q)$$



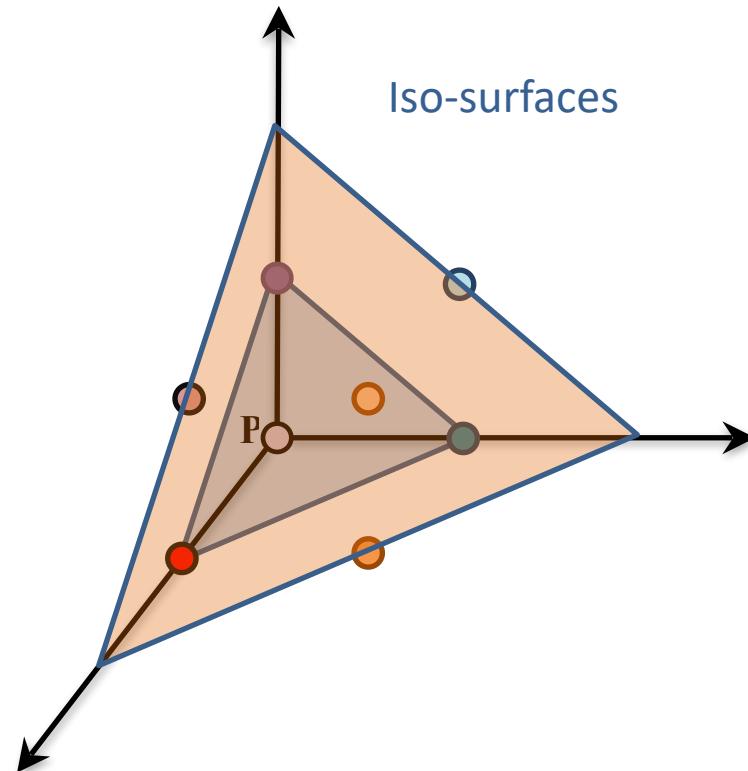
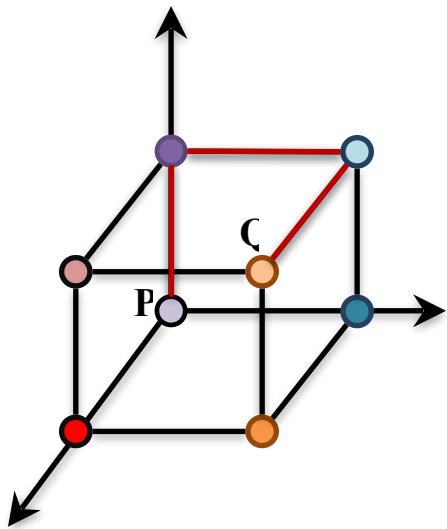
# Mahalanobis Distance

$$d(P, Q)^2 = (P - Q)^T \mathbf{S}^{-1} (P - Q)$$



# Hamming Distance

$$d(P, Q) = \sum_{i=1}^d I(p_i \neq q_i)$$



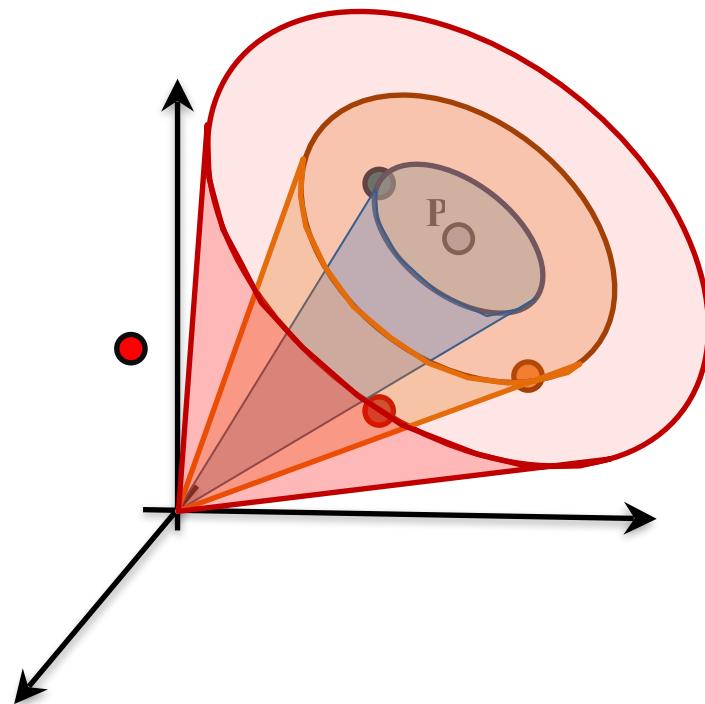
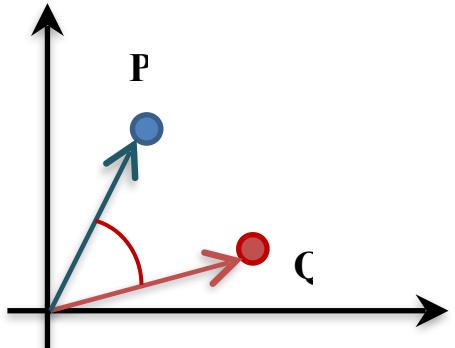
|   |   |   |   |   |   |   |
|---|---|---|---|---|---|---|
| 0 | 1 | 1 | 0 | 1 | 1 | 0 |
| 0 | 0 | 1 | 0 | 1 | 0 | 0 |

# Cosine Distance

Iso-surfaces

$$s(P, Q) = \cos\theta = \frac{\mathbf{P} \cdot \mathbf{Q}}{\|\mathbf{P}\| \|\mathbf{Q}\|}$$

$$d(P, Q) = 1 - s(P, Q)$$



# Summary of Distances

- Distance Metrics decide Neighborhoods
- Need not be a "metric"
  - Jaccard Distance
  - Edit Distance
- Feature vectors need not be of same length
- Selection of metric depends on the nature of feature vector