



数据挖掘导论

Introduction to Data Mining

Basic Classification



数据智能实验室
DATA INTELLIGENCE LABORATORY



浙江大学
Zhejiang University

Agenda

□ Classification: Basic Concepts

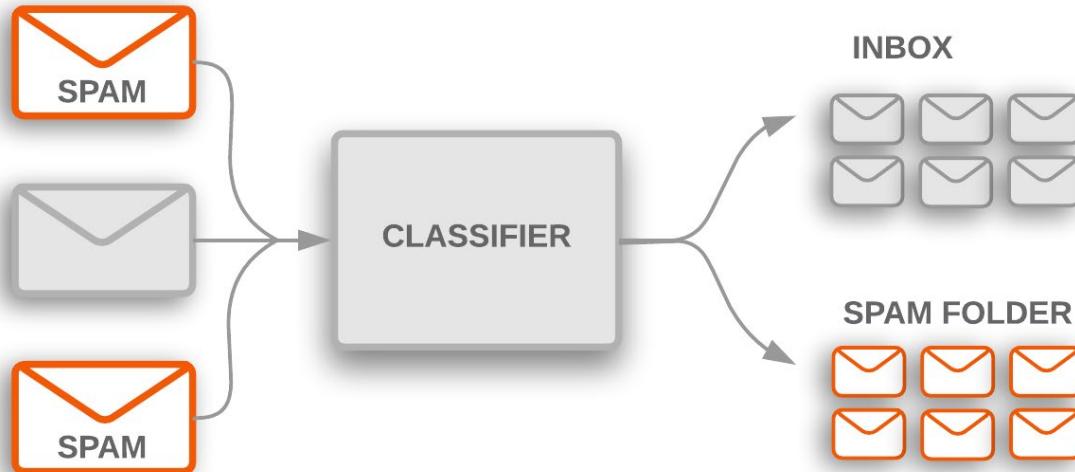
□ kNN Classifier

□ Decision Tree

□ Bayes Classification

□ Support Vector Machine

Classification on Text Data



系统通知

admin : 用户 zhangdongxiang@zju.edu.cn 维护原因为进一步提升邮件系统的安全性，我部门于近日陆续更新

本次将继承往届会议的优点，继续为领域内的专家学者提供高水准的国际交流平台

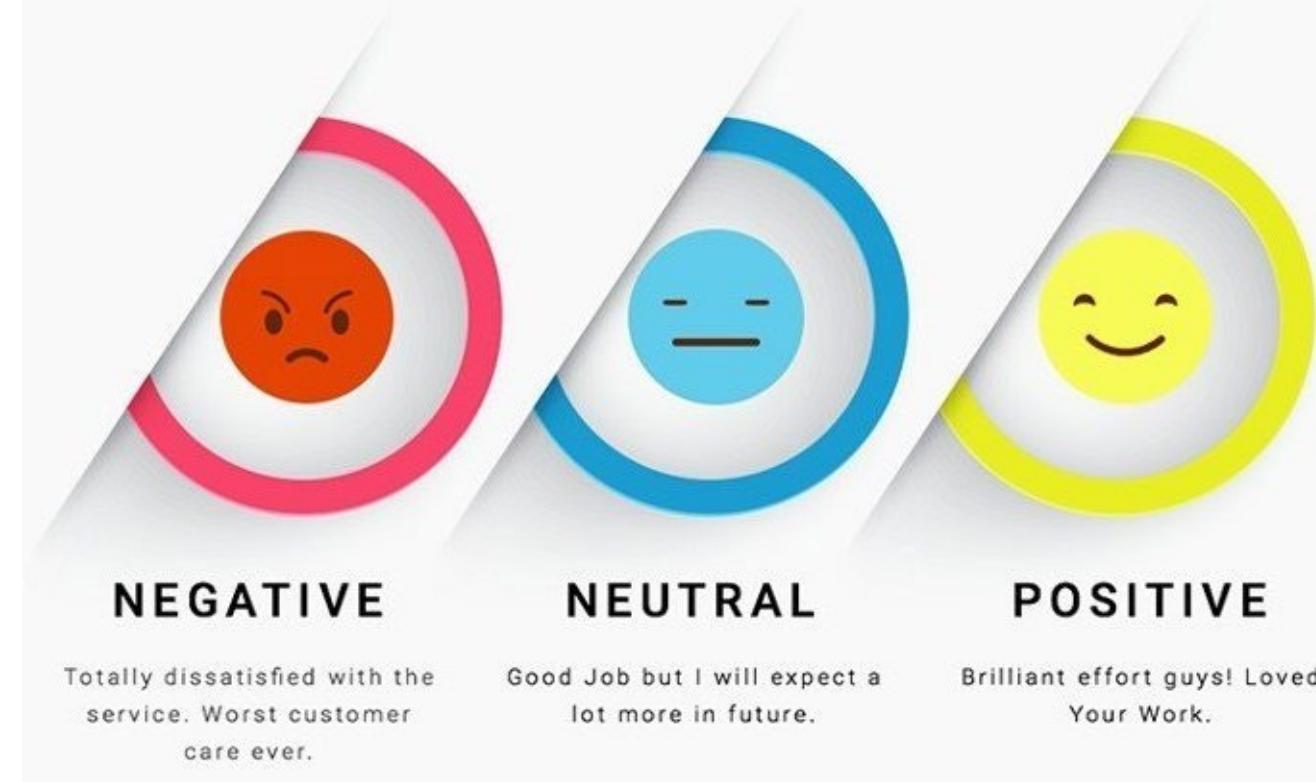
EI-- ICNISC-9th : 计算机网络-信息系统会议/时间：2023年10月27日-29日会议/地/点：中国武汉会议/语/言

早(12)

精美配图助力冲击高水平期刊

松迪科研绘图/科学可视化 : 如果邮件内容无法正常显示请点击[这里](#) 退订 投诉 培训安排: 5月10日-5月12日 合肥

SENTIMENT ANALYSIS



Classification on Image Data

The figure displays three sequential screenshots from a mobile application for flower identification:

- 手动缩放 (Manual Zoom):** The first screen shows a pink rose flower with a dashed white circle highlighting the area being analyzed. A red banner at the bottom provides instructions: "双指移动并缩放一朵花至虚线框内" (Move and zoom one flower into the dashed box) and features a hand icon with arrows indicating the zoom action.
- 匹配结果 (Matching Result):** The second screen shows the识别结果 (Identification Result) for the flower. It displays a large image of the flower and a card below it stating "100% MATCH" with "玫瑰 Rosa rugosa" and "花语:热恋".
- 花朵知识 (Flower Knowledge):** The third screen shows the花卉信息 (Flower Information) for the flower. It displays a large image of the flower and text: "玫瑰 Rosa rugosa", "花语:热恋", and a poem: "我看天看地，看路边凶猛生长的野草，看长河用力捞下的落日，却全都不及你。"

手动缩放
手动缩放花朵 · 更精准

匹配结果
百分制匹配结果 · 最专业

花朵知识
有趣的花朵介绍 · 涨知识

Classification on Video Data

PART FOUR 赛题升级 真知比拼

赛题描述

多模态短视频分类是视频理解领域的基础技术之一，在安全审核、推荐运营、内容搜索等领域有着非常广泛的应用。

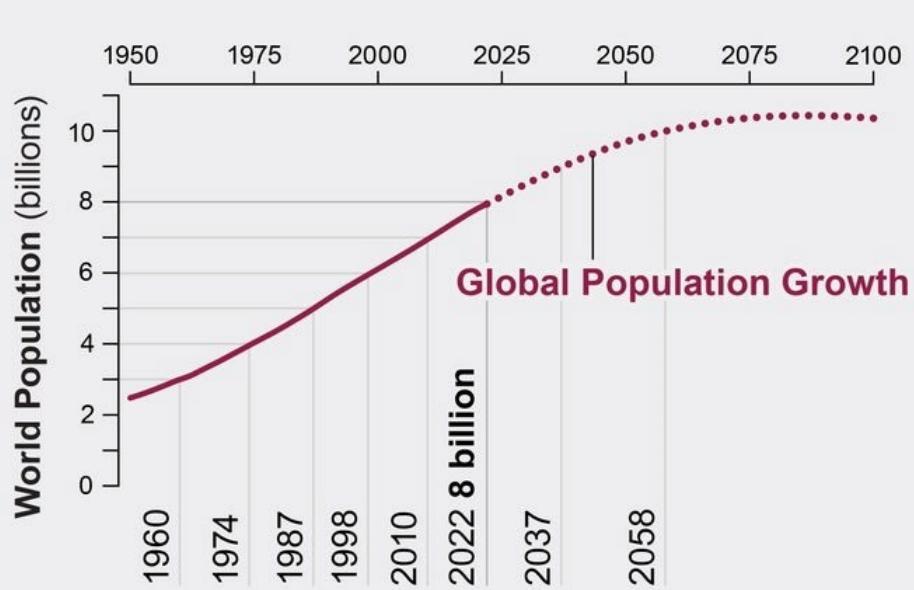
赛题要求

参赛队伍基于授权后的微信视频号短视频数据以及对应的分类标签标注，采用合理的机器学习技术对制定的测试短视频进行分类预测。

PART FIVE 王者对决 勇夺荣誉

-  总奖金池**56万**
-  学生团队奖项
-  周周星奖励
-  荣誉证书

Regression



Global Population Prediction

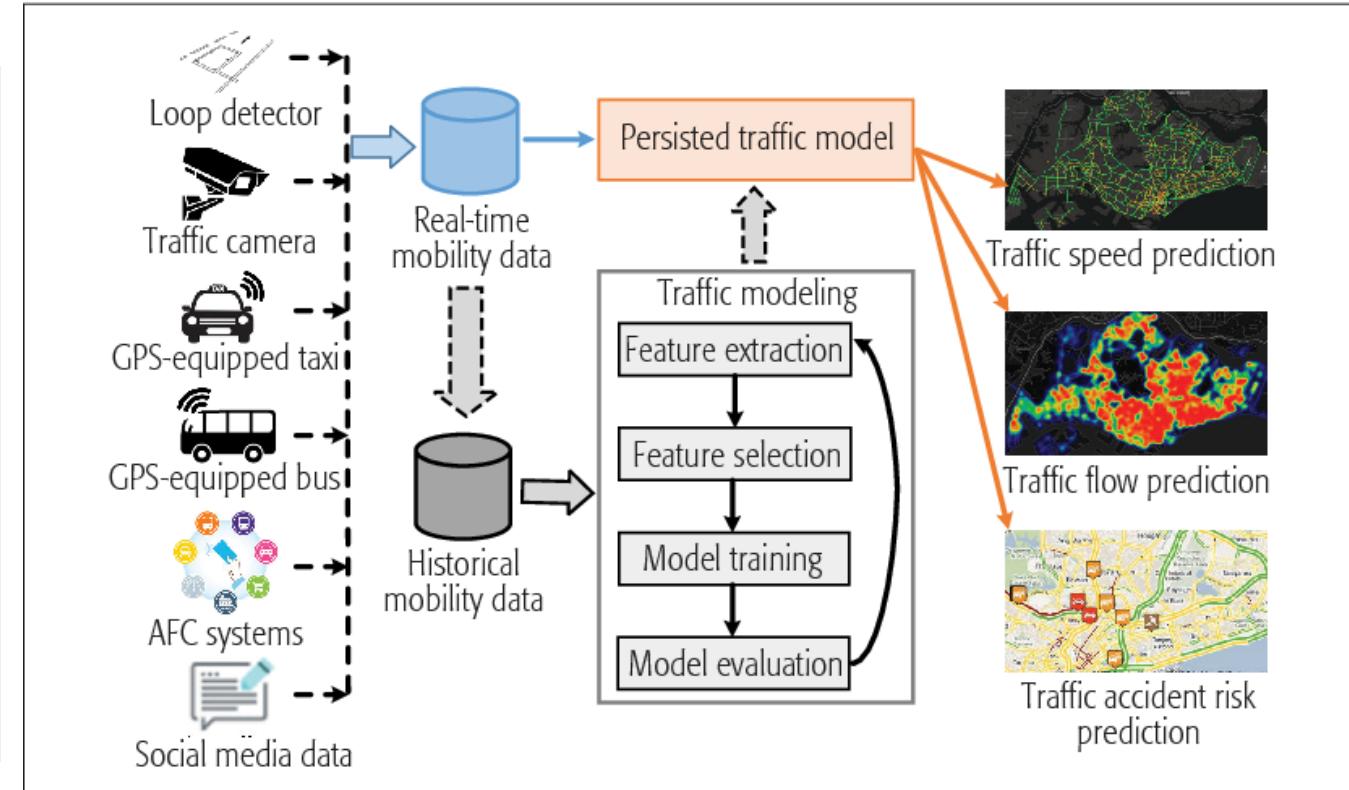


FIGURE 1. The basic components of urban traffic prediction.

Traffic Prediction

Prediction: Classification vs. Regression

□ Classification

- ✓ predicts categorical class labels

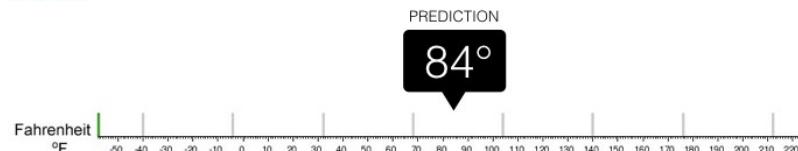
□ Regression

- ✓ models continuous-valued functions



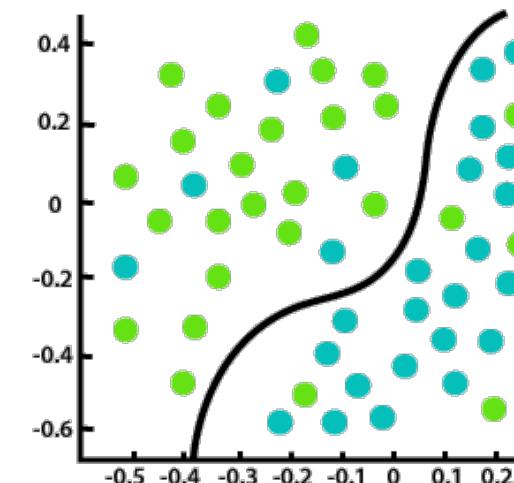
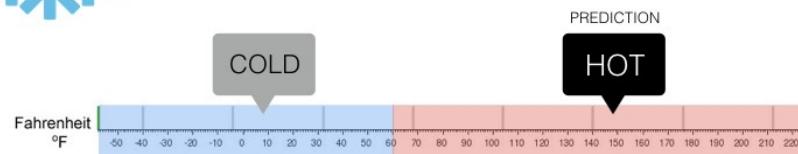
Regression

What is the temperature going to be tomorrow?

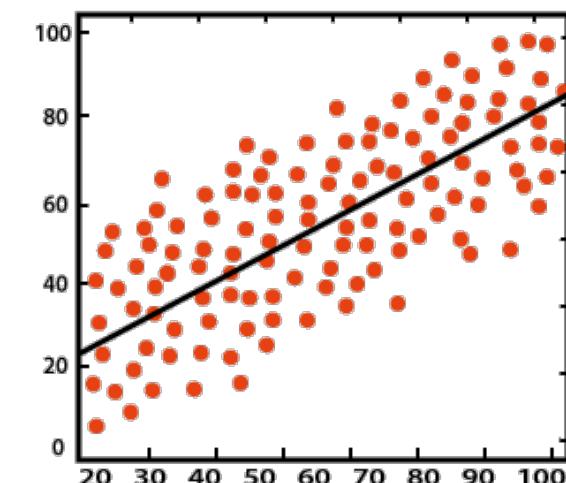


Classification

Will it be Cold or Hot tomorrow?

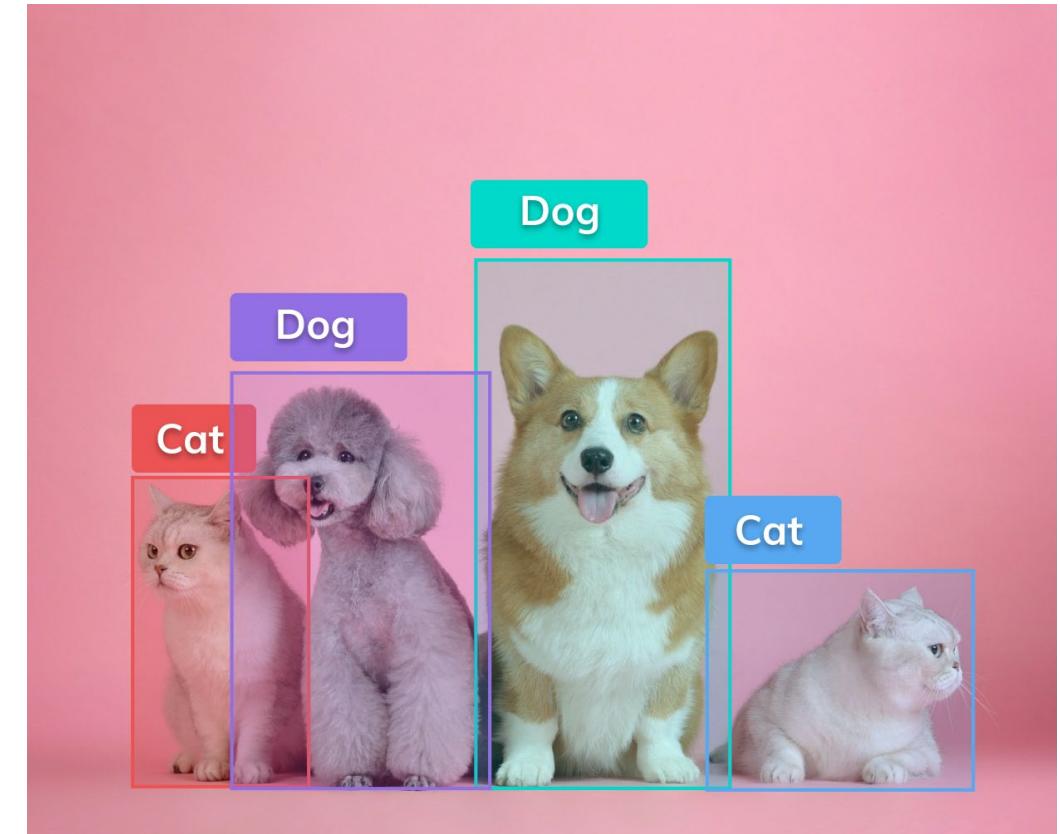
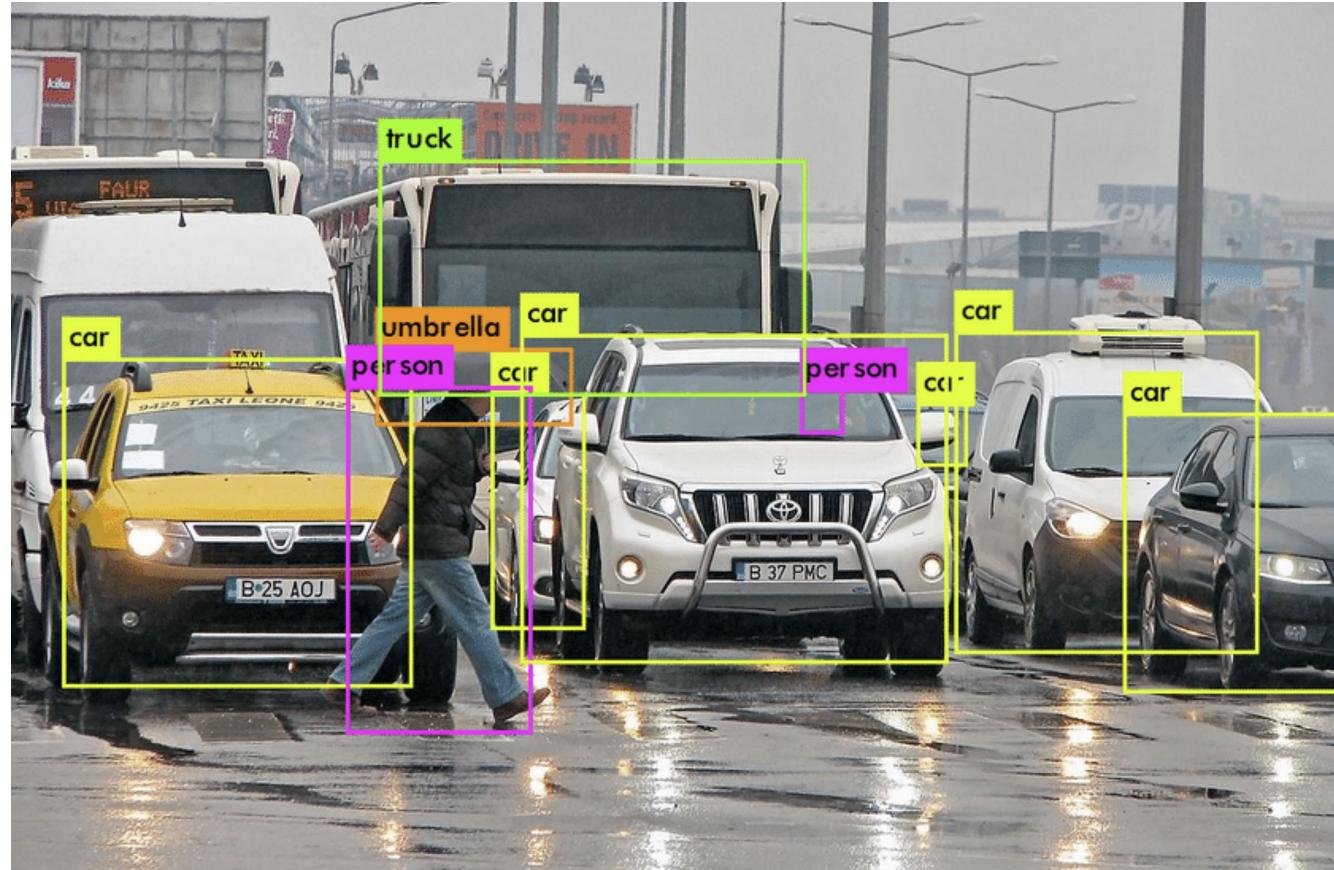


Classification



Regression

Object Detection



Supervised vs. Unsupervised Learning

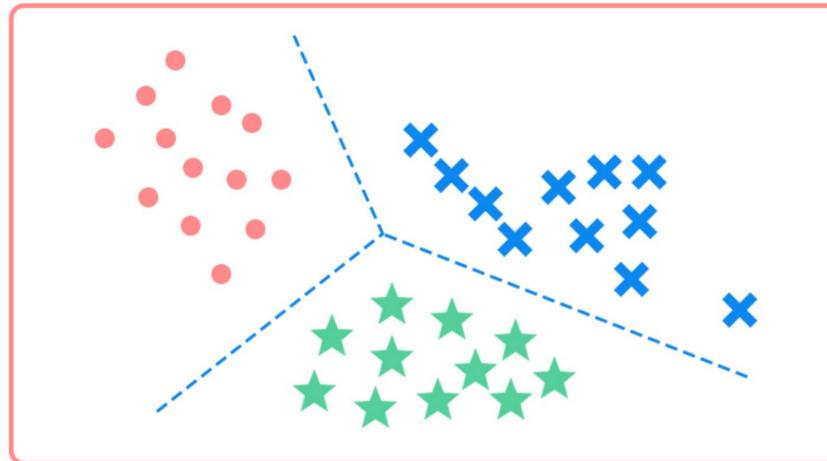
□ Supervised learning

- ✓ Supervision: The training data (observations, measurements, etc.) are accompanied by labels indicating the class of the observations
- ✓ New data is classified based on the training set

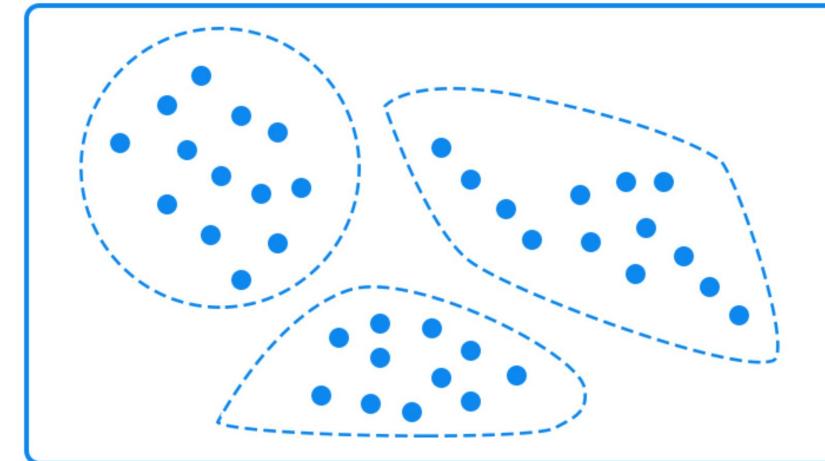
□ Unsupervised learning

- ✓ The class labels of training data is not available

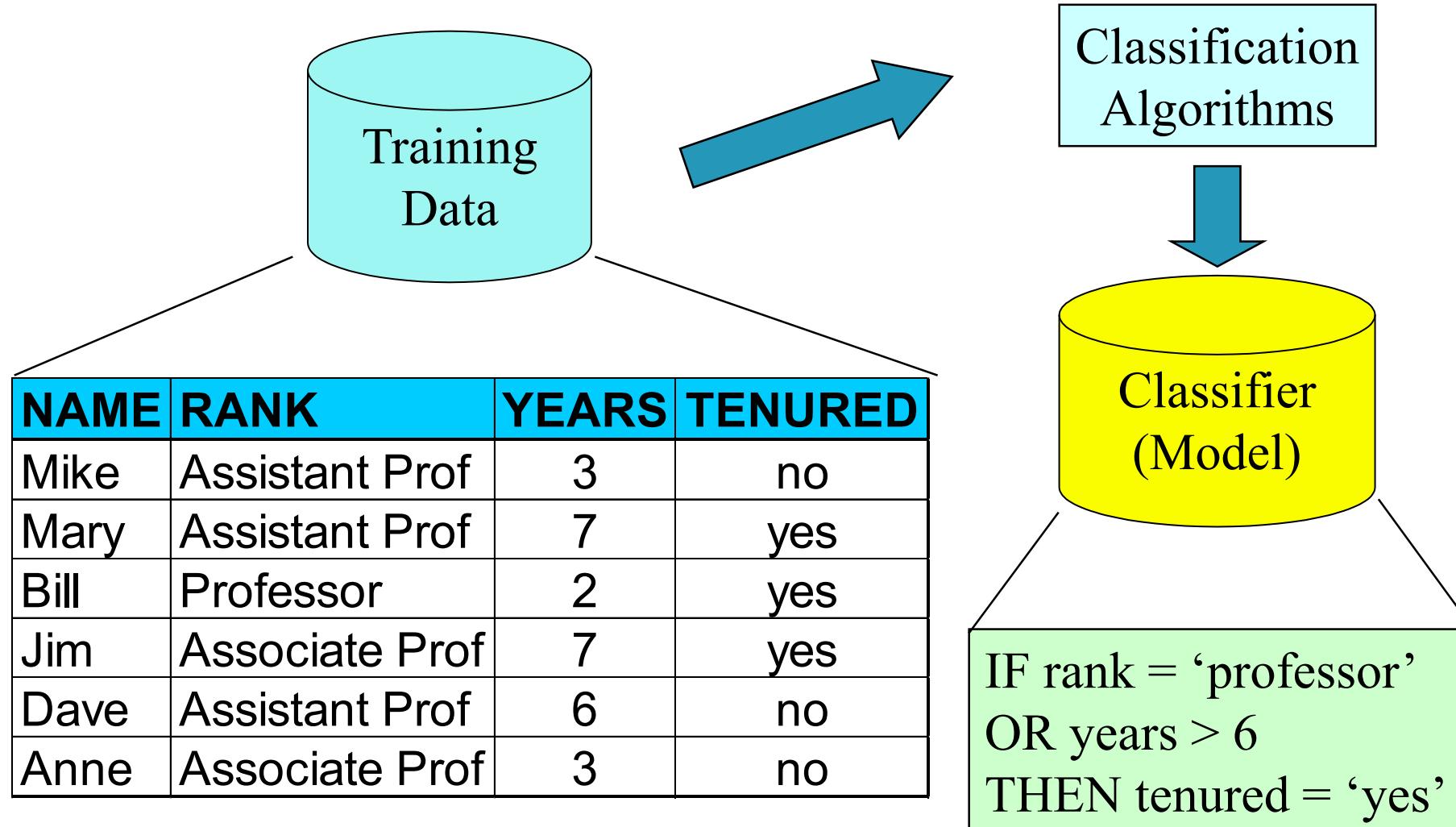
Classification



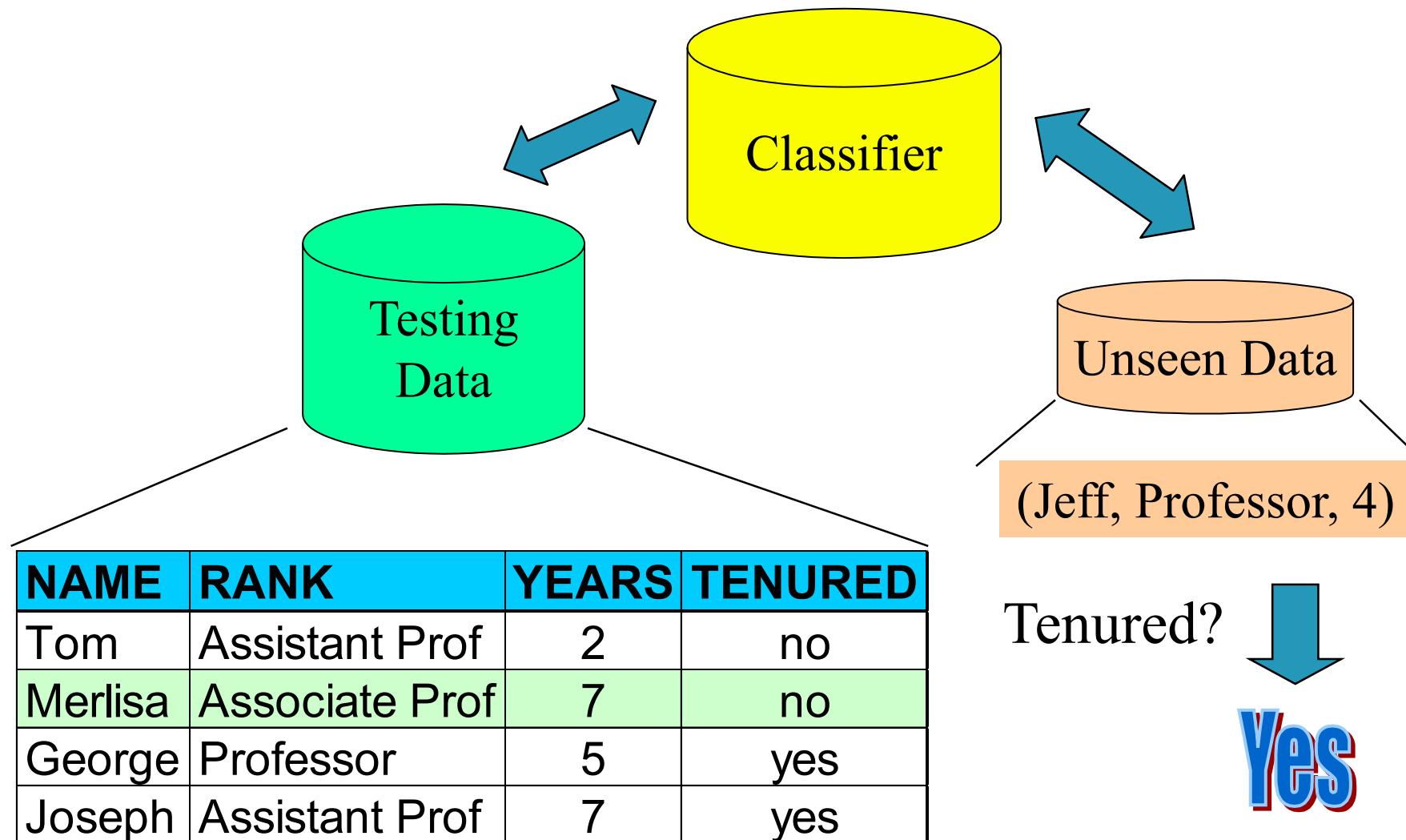
Clustering



Process (1): Model Construction



Process (2): Using the Model in Prediction



Data Splitting

□ Training Set

- ✓ The training set is the set of data we use to train a model

□ Validation Set

- ✓ The validation set is a set of data that we did not use when training our model that we use to assess how well the model performs on new data
- ✓ The class labels of data in the validation set are known in advance

□ Test Set

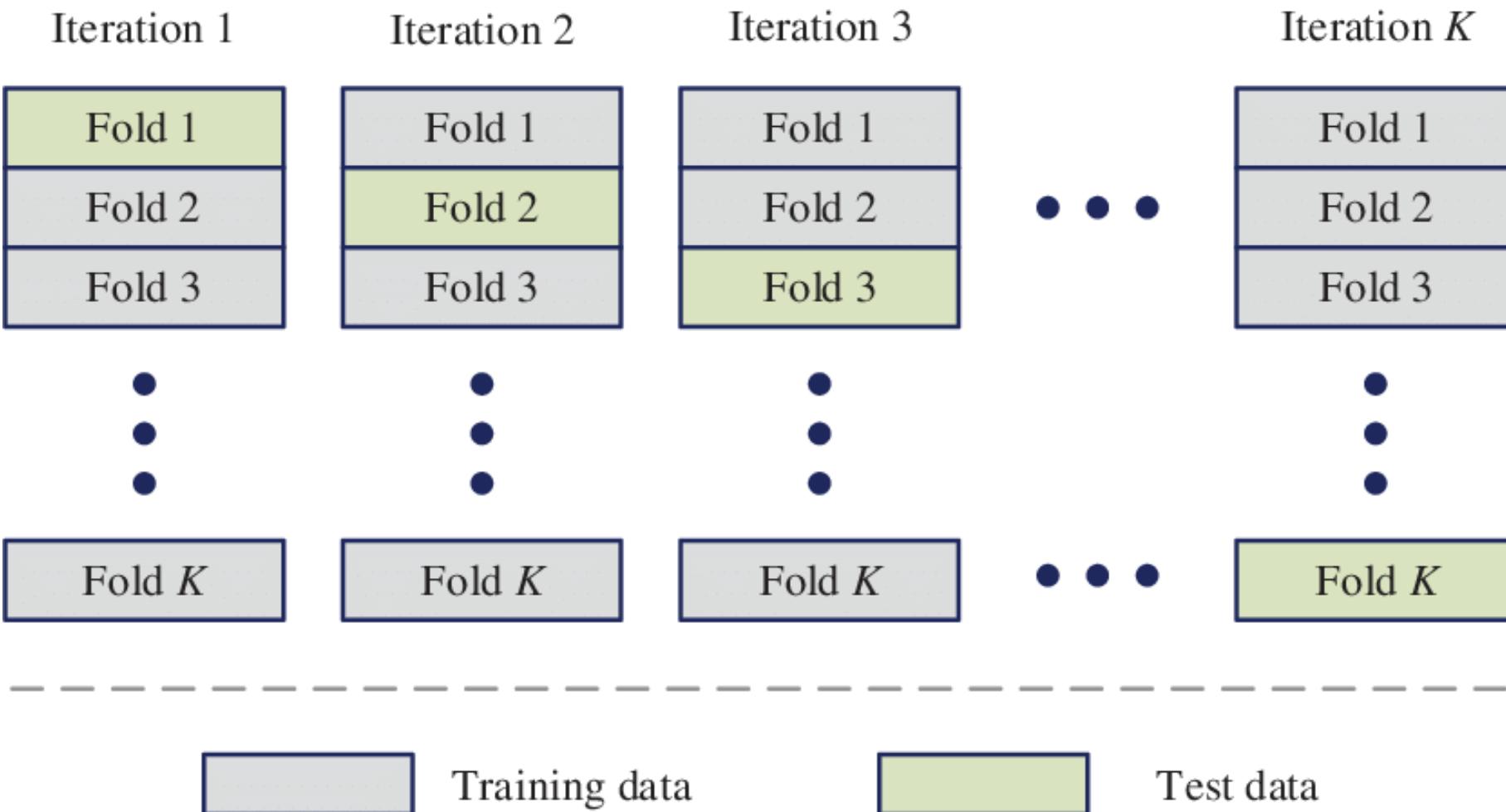
- ✓ The test set is a set of data we did not use to train our model or use in the validation set to inform our choice of parameters/input features
- ✓ In many Kaggle tasks, you are unaware of the labels for the test samples

Data Explorer

219.53 MB

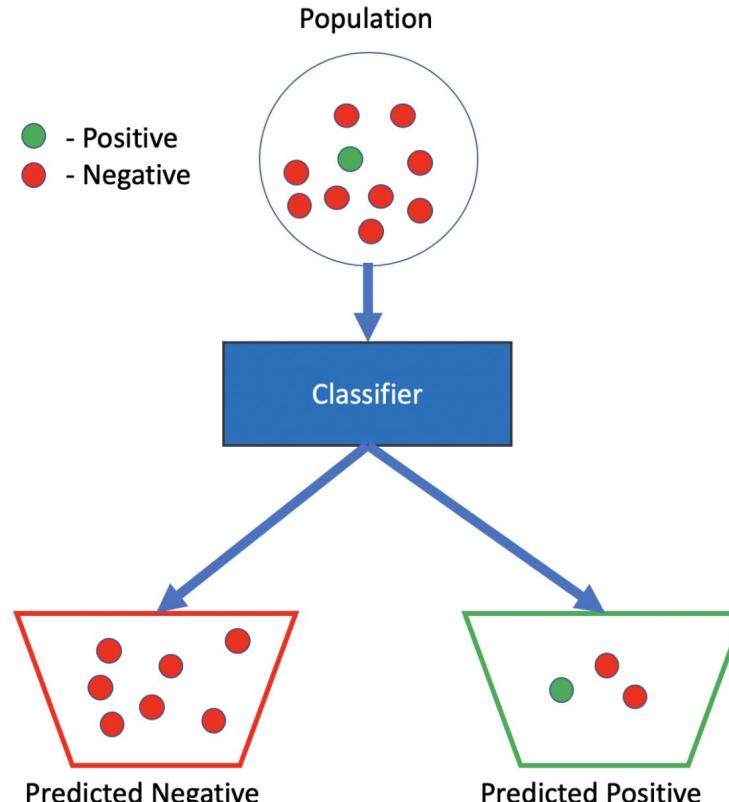
- ▼ test
 - ▶ images
 - ▶ train
- ▼ val
 - ▶ images
 - val_annotations.txt

Cross Validation



Performance Metrics of Binary Classification

一个医院新开发了一套癌症AI诊断系统，想评估其性能好坏。我们把病人得了癌症定义为Positive，没得癌症定义为Negative。那么，到底该用什么指标进行评估呢？



True Positive (TP): 把正样本成功预测为正。

True Negative (TN): 把负样本成功预测为负。

False Positive (FP): 把负样本错误地预测为正。

False Negative (FN): 把正样本错误的预测为负。

Real			
		Positive	Negative
Predicted	Positive	1	2
	Negative	0	7

Performance Metrics of Binary Classification

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN}$$

$$Precision = \frac{TP}{TP + FP}$$

$$Recall = \frac{TP}{TP + FN}$$

$$F1-score = \frac{2 \times Precision \times Recall}{Precision + Recall}$$

在诊断为癌症的一堆人中，到底有多少人真得了癌症？

在一堆真的癌症病人中，有多少人被成功检测出癌症？

在一堆癌症病人和正常人中，有多少人被系统给出了正确诊断结果（患癌或没患癌）？

		Real	
		Positive	Negative
Predicted	Positive	1	2
	Negative	0	7

$$\text{precision} = \frac{tp}{tp + fp} = \frac{1}{3} = 33\%$$

$$\text{recall} = \frac{tp}{tp + fn} = \frac{1}{1} = 100\%$$

Performance Metrics of Multiclass Classification

		Predicted		
		Cat	Dog	Pig
Actual	Cat	40	20	10
	Dog	35	85	40
	Pig	0	10	20

	TP	FN	FP
Cat	40	30	35
Dog	85	75	30
Pig	20	10	50

→ 20只Cat和10只Pig被错误预测成Dog

↓ 10只Pig被错误预测成Dog

$$\text{Accuracy: } (40+85+20)/260=0.558$$

$$P_{cat} = 8/15, P_{dog} = 17/23, P_{pig} = 2/7 \quad (\text{P代表Precision})$$

$$R_{cat} = 4/7, R_{dog} = 17/32, R_{pig} = 2/3 \quad (\text{R代表Recall})$$

$$\text{Macro-Precision} = \frac{P_{cat} + P_{dog} + P_{pig}}{3} = 0.5194$$

$$\text{Macro-Recall} = \frac{R_{cat} + R_{dog} + R_{pig}}{3} = 0.5898$$

Performance Metrics of Regression

$$MAE = \frac{1}{n} \sum_{i=1}^n |\hat{y}_i - y_i|$$

$$MSE = \frac{1}{n} \sum_{i=1}^n (\hat{y}_i - y_i)^2$$

$$RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^n (\hat{y}_i - y_i)^2}$$

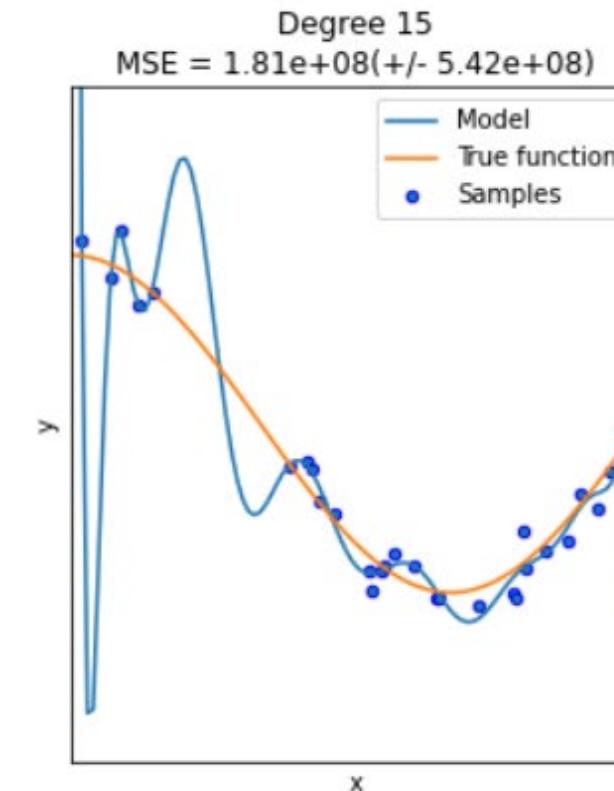
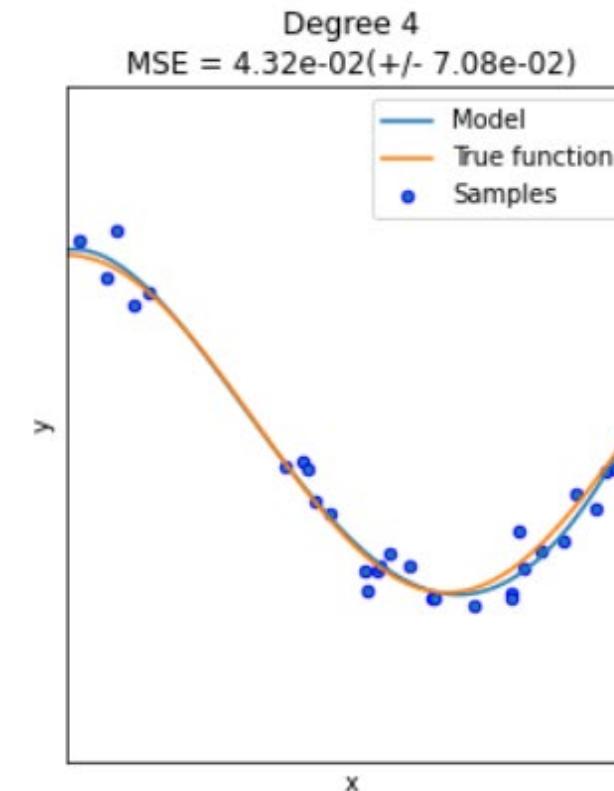
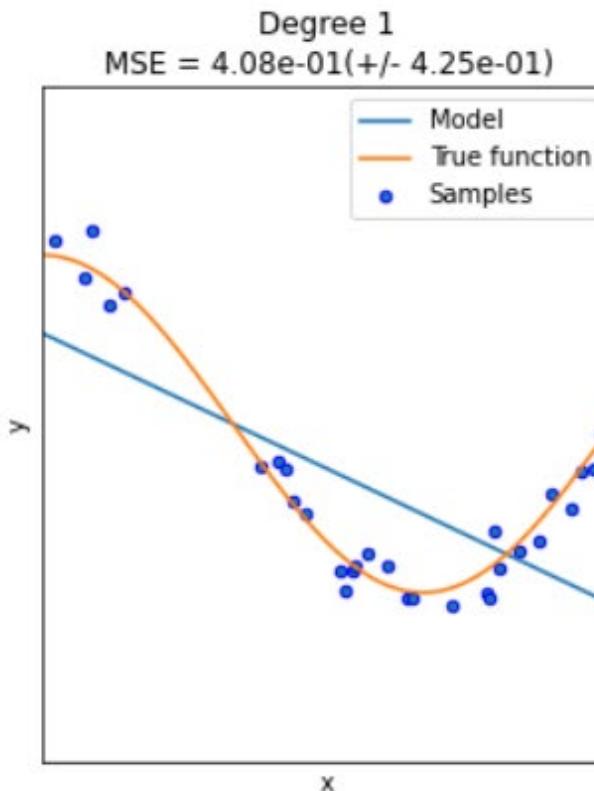
Overfitting and Underfitting

□ Overfitting

- ✓ The model performs well on the training data but does not perform well on the evaluation data

□ Underfitting

- ✓ The model performs poorly on the training



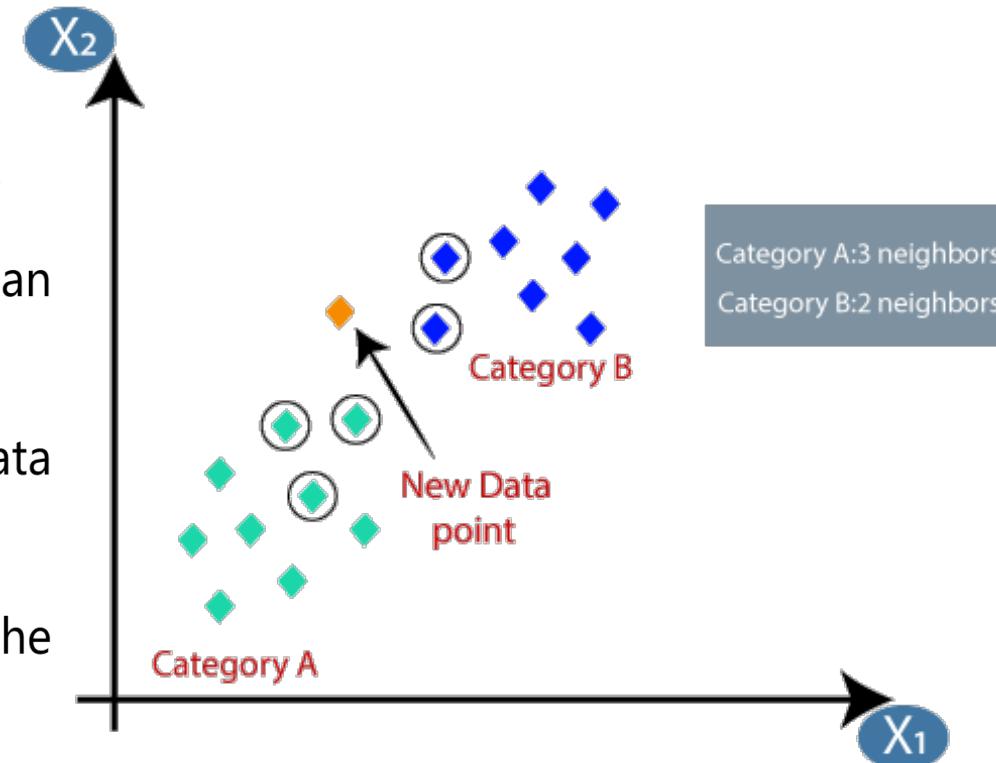
Agenda

- Classification: Basic Concepts
- kNN Classifier
- Decision Tree
- Bayes Classification
- Support Vector Machine

kNN Classifier

□ Non-parametric algorithm, which means it does not make any assumption on underlying data

- Step-1: Select the number K of the neighbors
- Step-2: Calculate the Euclidean distance of **K number of neighbors**
- Step-3: Take the K nearest neighbors as per the calculated Euclidean distance.
- Step-4: Among these k neighbors, count the number of the data points in each category.
- Step-5: Assign the new data points to that category for which the number of the neighbor is maximum.



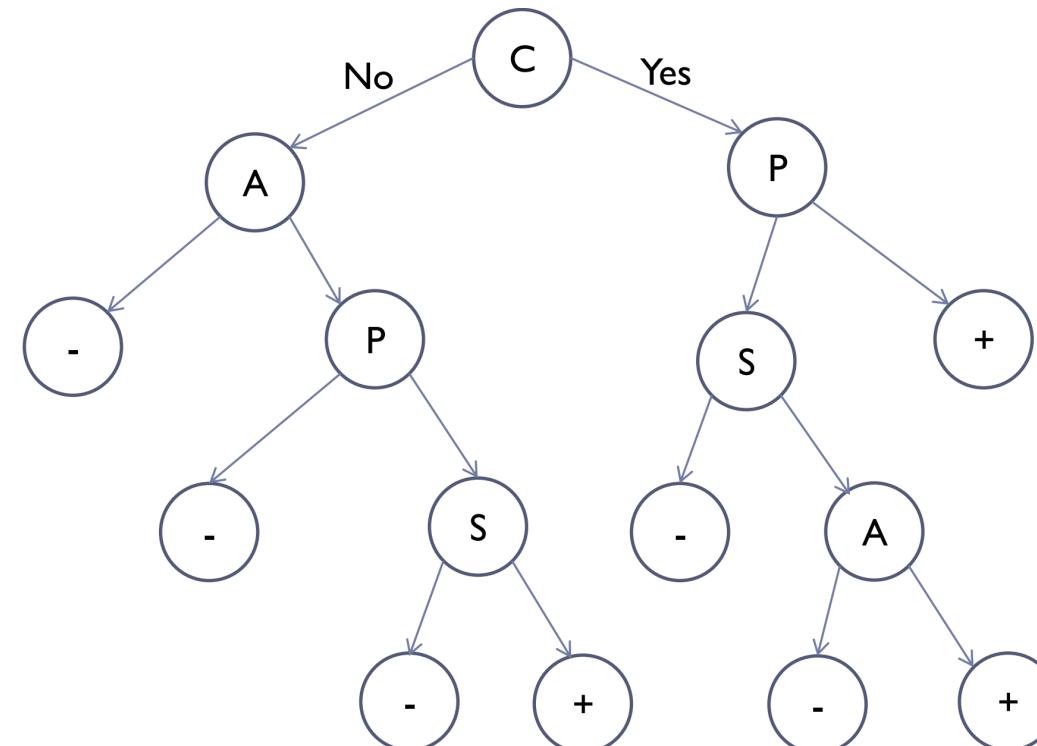
Agenda

- Classification: Basic Concepts
- kNN Classifier
- Decision Tree
- Bayes Classification
- Support Vector Machine

Decision Tree Structure

- Each leaf node represents a class label
- Each internal node corresponds to denotes a test on an attribute
 - ✓ Edges to children for each of the possible values of that attribute

- ▶ Attributes:
 - ▶ A: age>40
 - ▶ C: chest pain
 - ▶ S: smoking
 - ▶ P: physical test



- ▶ Label:
 - ▶ Heart disease (+), No heart disease (-)

Decision Tree Learning

- The most common strategy for DT learning is a greedy top-down approach

Basic decision tree building algorithm:

- Pick some feature/attribute (how to pick the “best”?)
- Partition the data based on the value of this attribute
- Recurse over each new partition (when to stop?)

Attribute Selection

□ Max-Gain: Choose the attribute that has the largest expected information gain

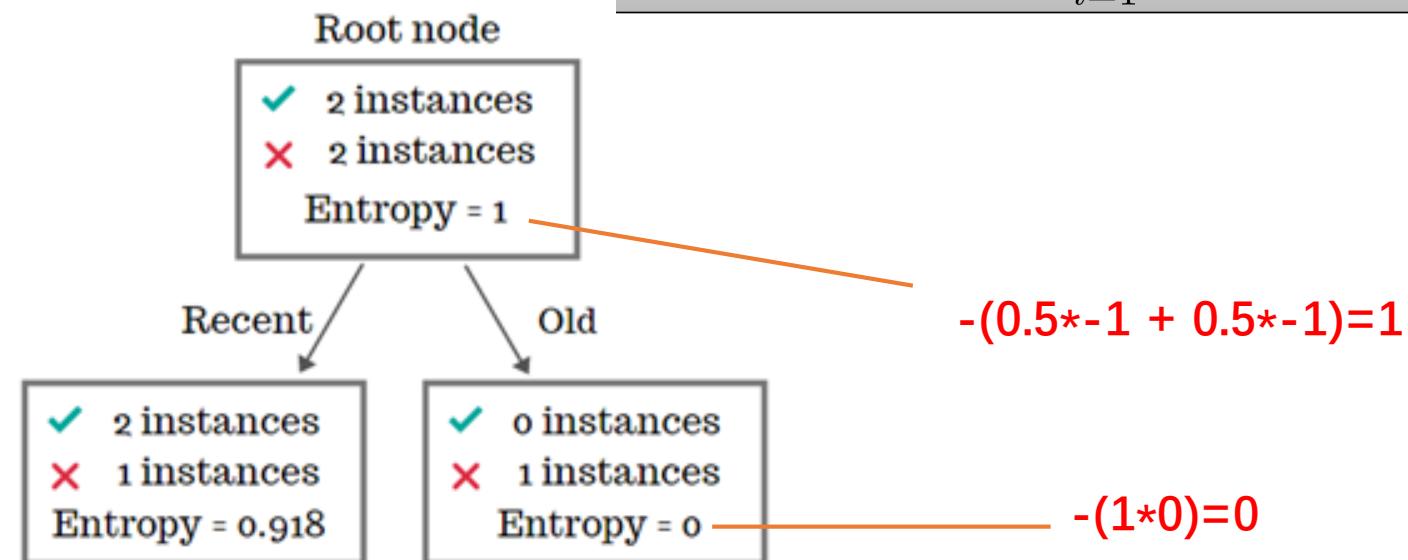
Information Gain (IG)

$$\text{IG} = \text{Entropy}(\text{Parent}) - \text{weighted_avg} * \text{Entropy}(\text{Children})$$

$$\text{Entropy}(\text{Parent}) = 1$$

Information gain for Age

Age	Mileage	Road Tested	Buy
Recent	Low	Yes	Buy ✓
Recent	High	Yes	Buy ✓
Old	Low	No	Don't buy ✗
Recent	High	No	Don't buy ✗



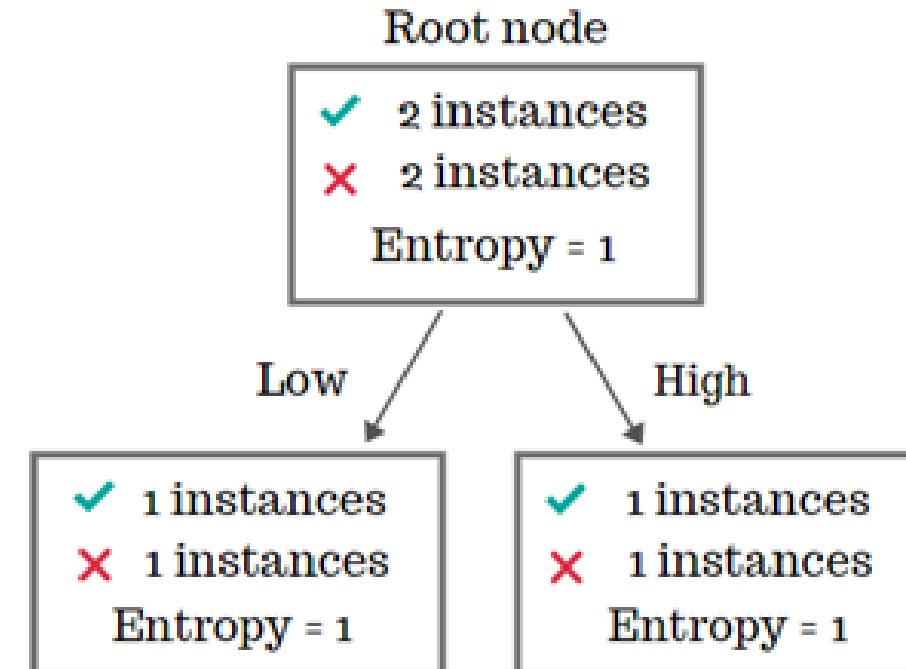
$$\text{Children Entropy: } 0.75 * 0.918 + 0.25 * 0 = 0.688$$

$$\text{Information Gain: } 1 - 0.688 = 0.312$$

Attribute Selection

- Max-Gain: Choose the attribute that has the largest expected information gain
- Information gain for Milage

Age	Mileage	Road Tested	Buy
Recent	Low	Yes	Buy ✓
Recent	High	Yes	Buy ✓
Old	Low	No	Don't buy ✗
Recent	High	No	Don't buy ✗



$$\text{Children Entropy} = \frac{1}{2} (1) + \frac{1}{2} (1) = 1$$

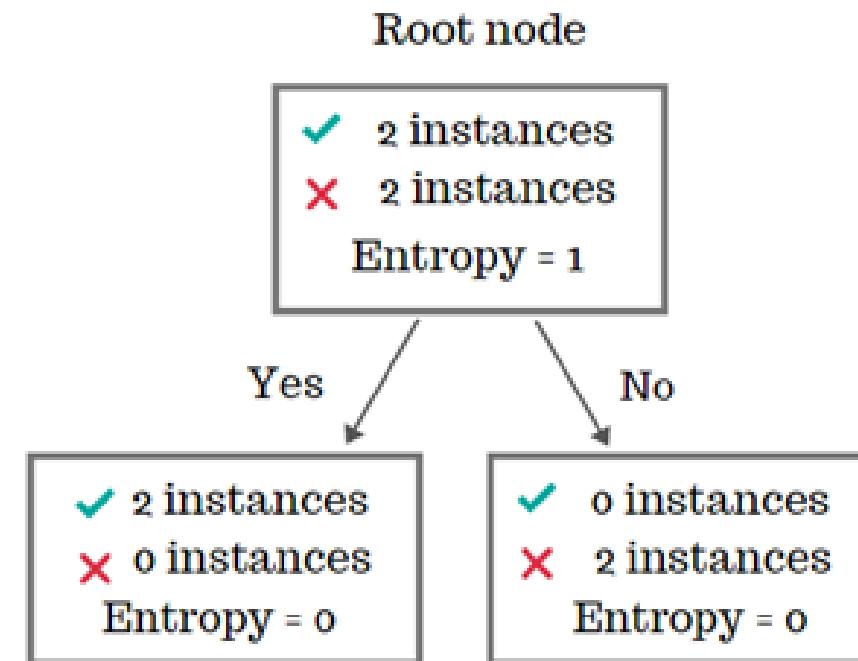
$$\text{Information gain} = 1 - 1 = 0$$

Attribute Selection

- Max-Gain: Choose the attribute that has the largest expected information gain

Information gain for Road Tested

Age	Mileage	Road Tested	Buy
Recent	Low	Yes	Buy ✓
Recent	High	Yes	Buy ✓
Old	Low	No	Don't buy ✗
Recent	High	No	Don't buy ✗



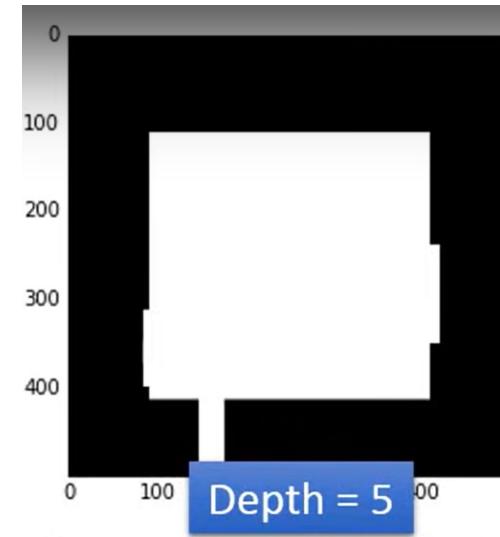
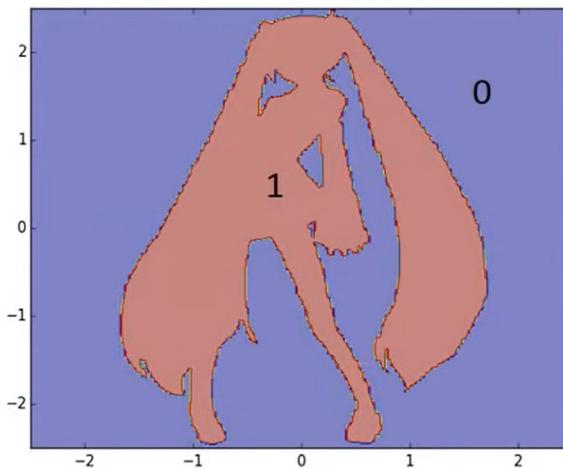
$$\text{Children Entropy} = \frac{1}{2}(0) + \frac{1}{2}(0) = 0$$

$$\text{Information gain} = 1 - 0 = 1$$

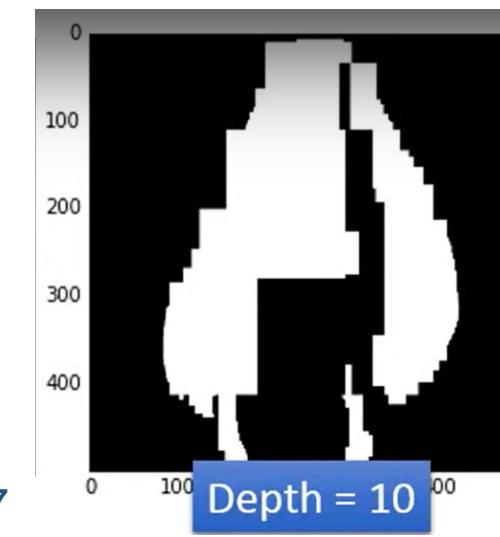
Decision Tree Learning

□ What functions can be represented by decision trees?

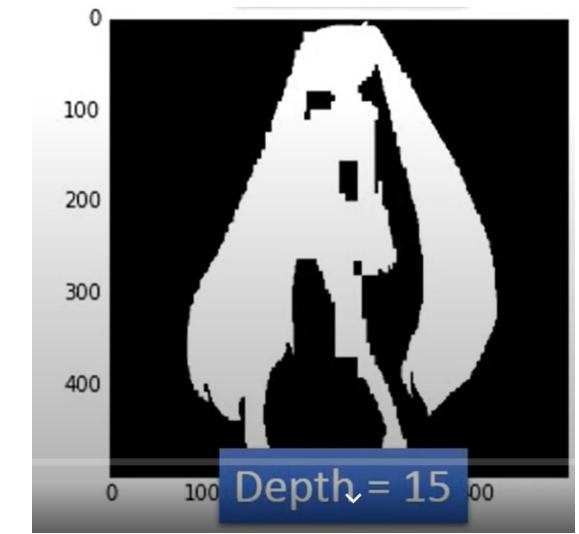
- ✓ Every function can be represented by a sufficiently complicated decision tree



27



Depth = 10



Depth = 15

When to Stop

□ It stops at smallest acceptable tree

- ✓ Occam's razor: prefer the simplest hypothesis that fits the data
- ✓ If all attributes have small information gain, then don't recurse

Agenda

□ Classification: Basic Concepts

□ kNN Classifier

□ Decision Tree

□ Bayes Classification

□ Support Vector Machine

Naïve Bayesian Classification

	Animals	Size of Animal	Body Color	Can we Pet them
0	Dog	Medium	Black	Yes
1	Dog	Big	White	No
2	Rat	Small	White	Yes
3	Cow	Big	White	Yes
4	Cow	Small	Brown	No
5	Cow	Big	Black	Yes
6	Rat	Big	Brown	No
7	Dog	Small	Brown	Yes
8	Dog	Medium	Brown	Yes
9	Cow	Medium	White	No
10	Dog	Small	Black	Yes
11	Rat	Medium	Black	No
12	Rat	Small	Brown	No
13	Cow	Big	White	Yes

Test Data: (Cow, Medium, Black)

$P(\text{Yes}|\text{Cow, Medium, Black})$

$P(\text{No}|\text{Cow, Medium, Black})$

Which probability is higher?



Naïve Bayesian Classification

□ Why is it called Naïve?

- ✓ It assumes features are independent of each other, which is not true in real life

□ Why is it popular?

- ✓ It is a probabilistic approach and the predictions can be made instantly
- ✓ It can be used for both binary and multi-class classification problems

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

↓

$$P(Y = k|X_1, X_2, \dots, X_n) = \frac{P(X_1|Y = k) * P(X_2|Y = k) * \dots * P(X_n|Y = k) * P(Y = k)}{P(X_1) * P(X_2) * \dots * P(X_n)}$$

Naïve Bayesian Classification

	Animals	Size of Animal	Body Color	Can we Pet them
0	Dog	Medium	Black	Yes
1	Dog	Big	White	No
2	Rat	Small	White	Yes
3	Cow	Big	White	Yes
4	Cow	Small	Brown	No
5	Cow	Big	Black	Yes
6	Rat	Big	Brown	No
7	Dog	Small	Brown	Yes
8	Dog	Medium	Brown	Yes
9	Cow	Medium	White	No
10	Dog	Small	Black	Yes
11	Rat	Medium	Black	No
12	Rat	Small	Brown	No
13	Cow	Big	White	Yes

Test Data: (Cow, Medium, Black)

$$P1 = P(\text{Cow} | \text{Yes}) = 3/8$$

$$P2 = P(\text{Medium} | \text{Yes}) = 2/8$$

$$P3 = P(\text{Black} | \text{Yes}) = 3/8$$

$$P(\text{Yes}) = 8/14$$

$$P(\text{Yes} | \text{Cow, Medium, Black})$$

$$= P1 * P2 * P3 * P(\text{Yes}) / P(\text{Cow, Medium, Black})$$

$$= 3/8 * 2/8 * 3/8 * 8/14 / P(\text{Cow, Medium, Black})$$

$$= 0.02 / P(\text{Cow, Medium, Black})$$

$$P(\text{No} | \text{Cow, Medium, Black})$$

$$= 2/6 * 2/6 * 1/6 * 6/14 / P(\text{Cow, Medium, Black})$$

$$= 0.0079 / P(\text{Cow, Medium, Black})$$

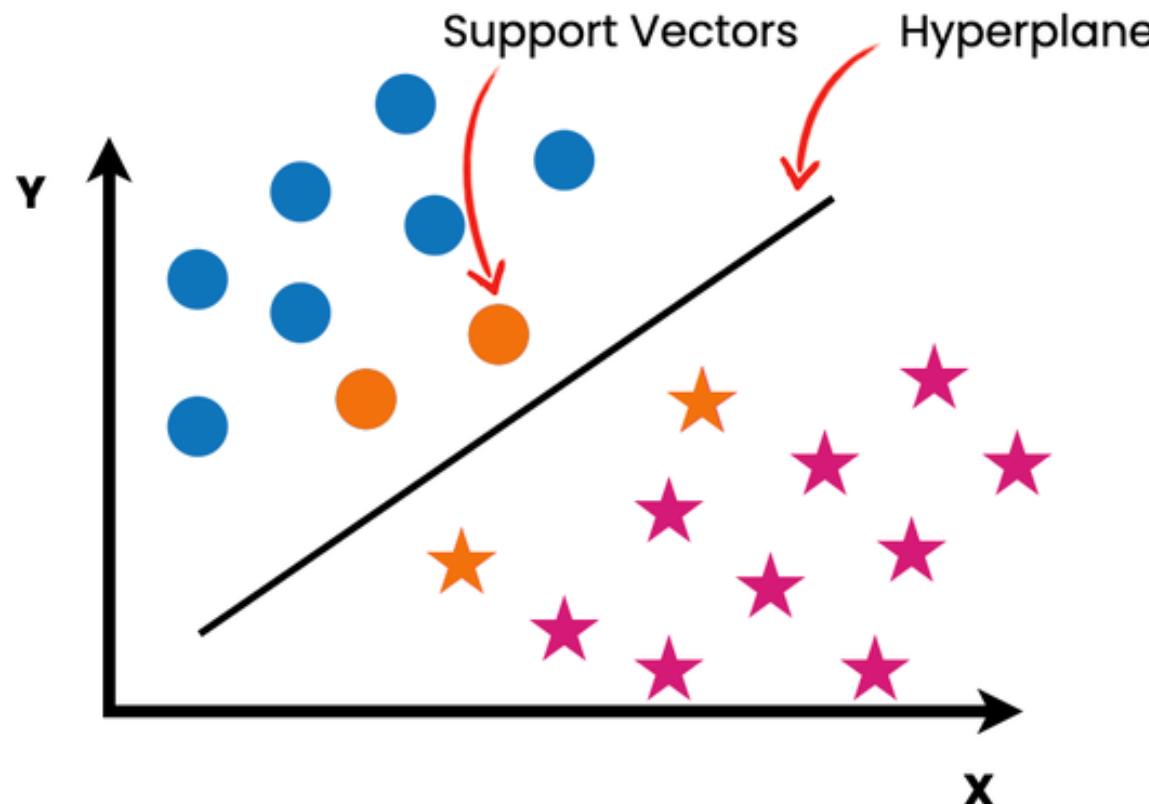


Agenda

- Classification: Basic Concepts
- kNN Classifier
- Decision Tree
- Bayes Classification
- Support Vector Machine

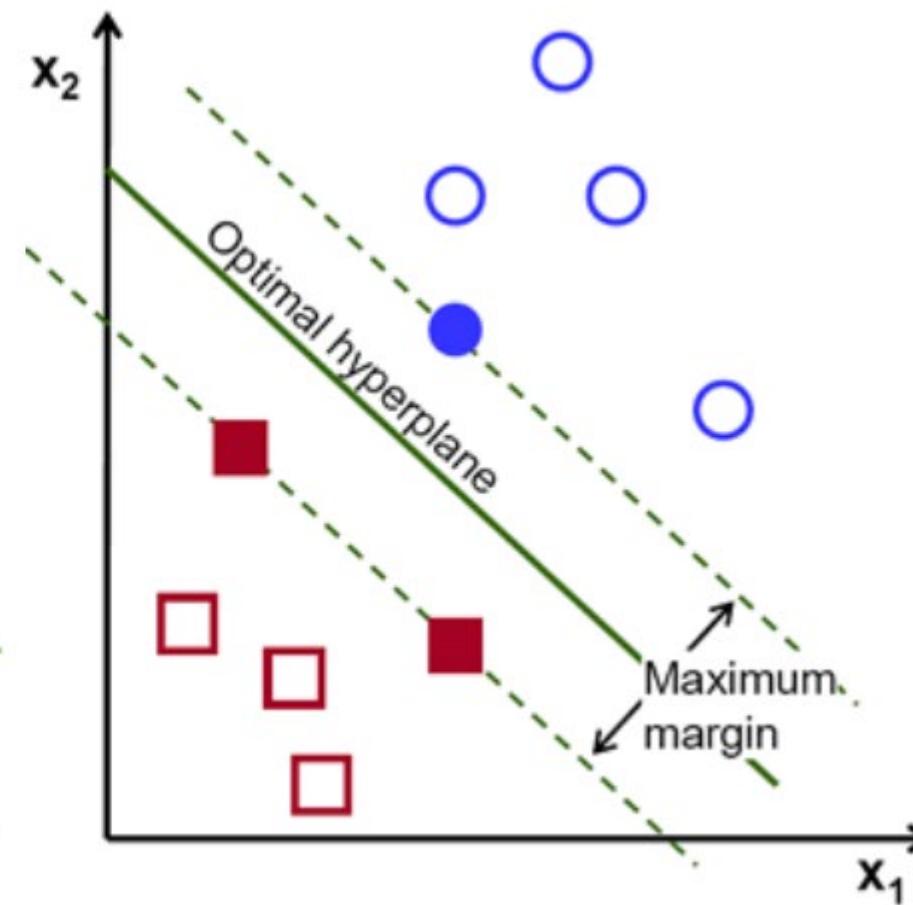
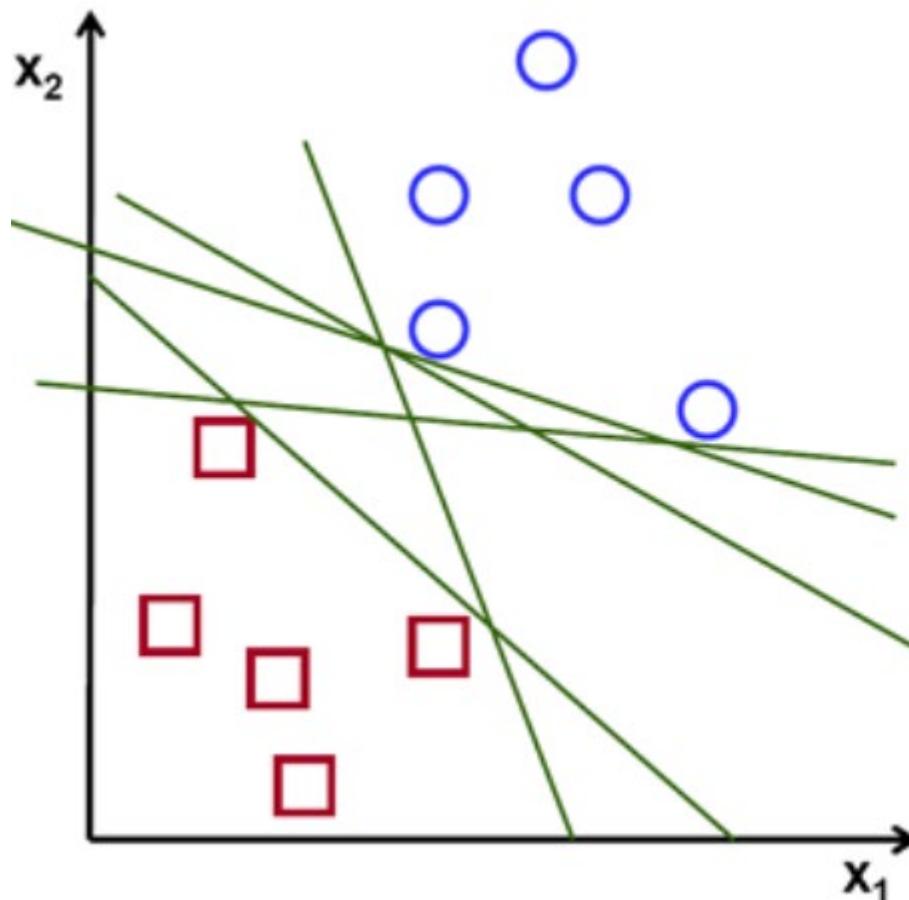
Support Vector Machine

□ Basic idea



Support Vector Machine

□ Which line is the best?



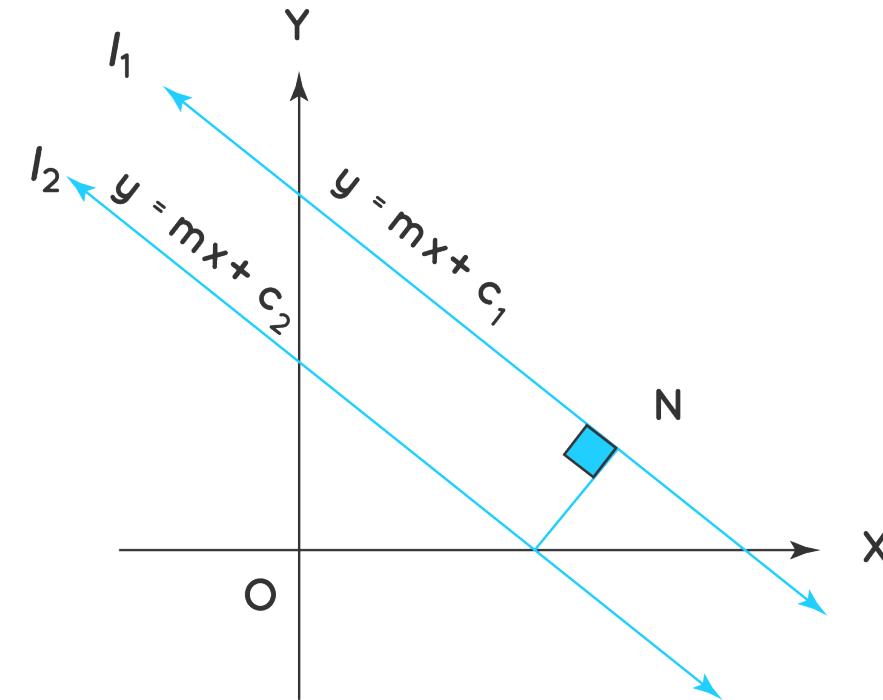
Support Vector Machine

□ How to represent a line?

- ✓ We can use $y=ax+b$
- ✓ We can rewrite it as $[a, -1] \begin{pmatrix} x \\ y \end{pmatrix} + b = 0$

□ The distance between two parallel lines

Distance Between Two Lines



$$d = \frac{|c_2 - c_1|}{\sqrt{1 + m^2}}$$



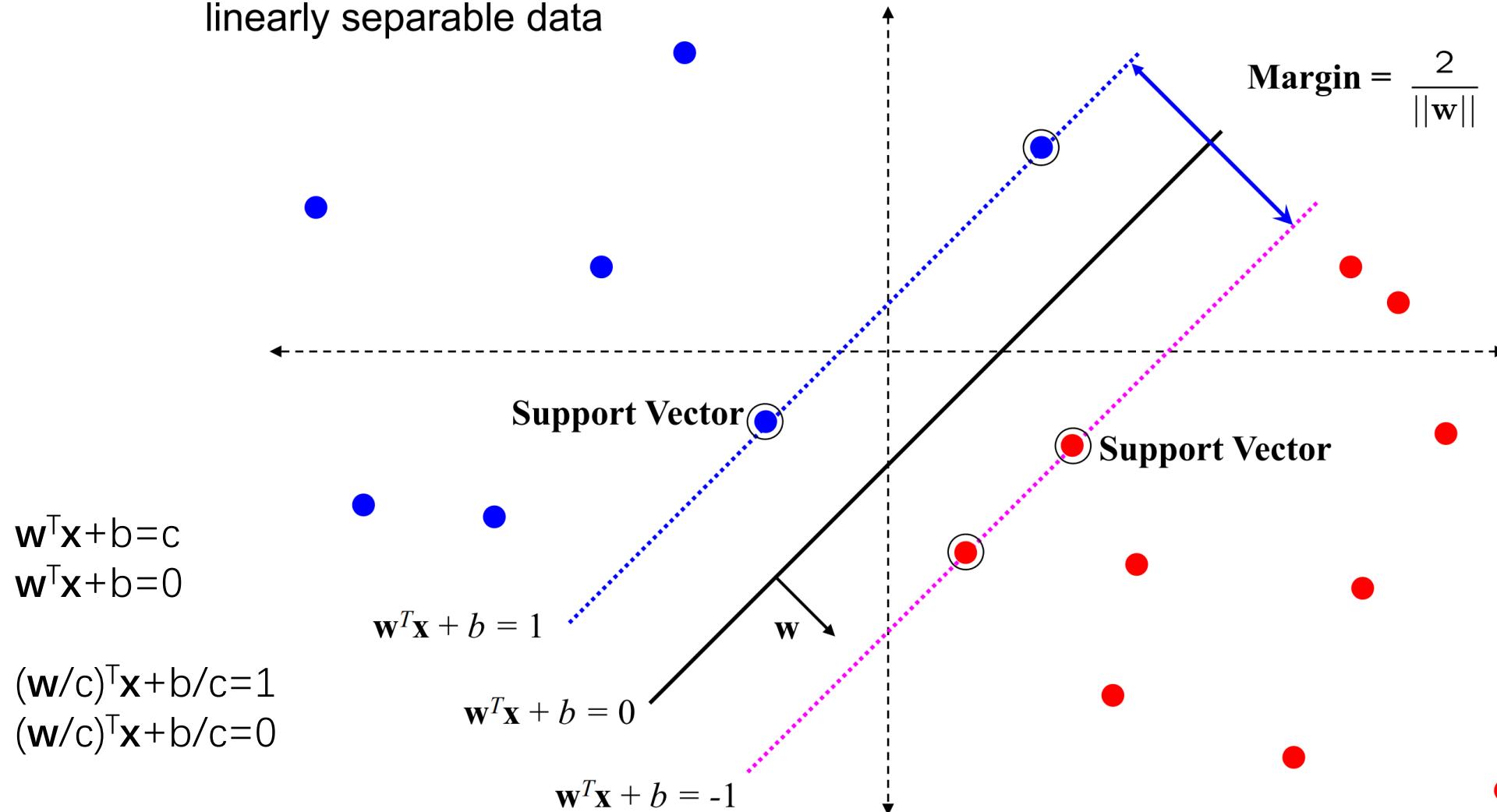
数据智能实验室
DATA INTELLIGENCE LABORATORY



浙江大学
Zhejiang University

Support Vector Machine

linearly separable data



Example from: <https://www.robots.ox.ac.uk/~az/lectures/ml/lect2.pdf>

Support Vector Machine

- Learning the SVM can be formulated as an optimization:

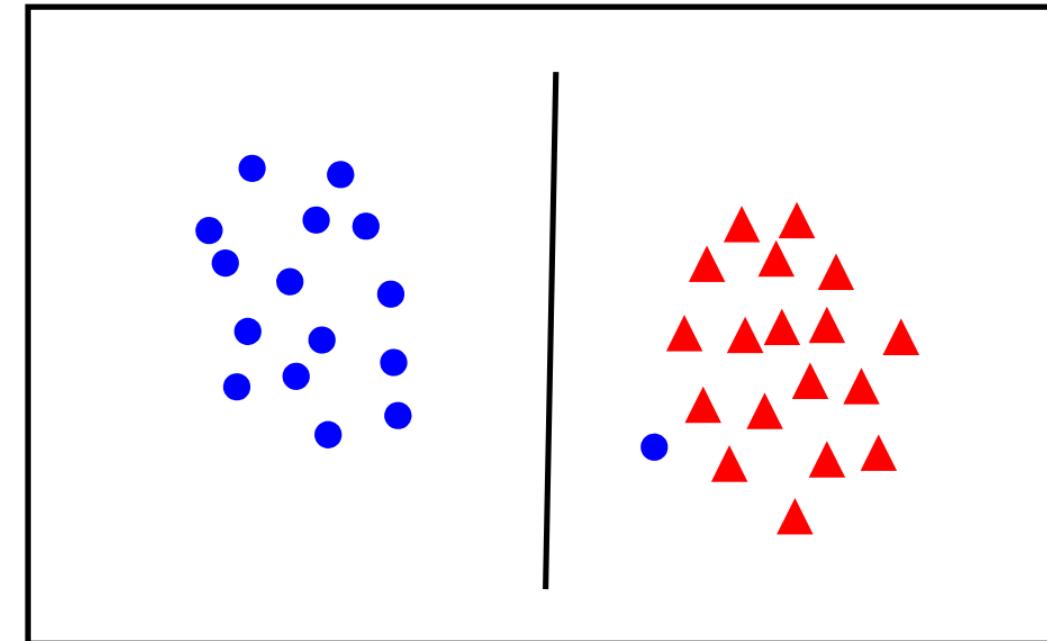
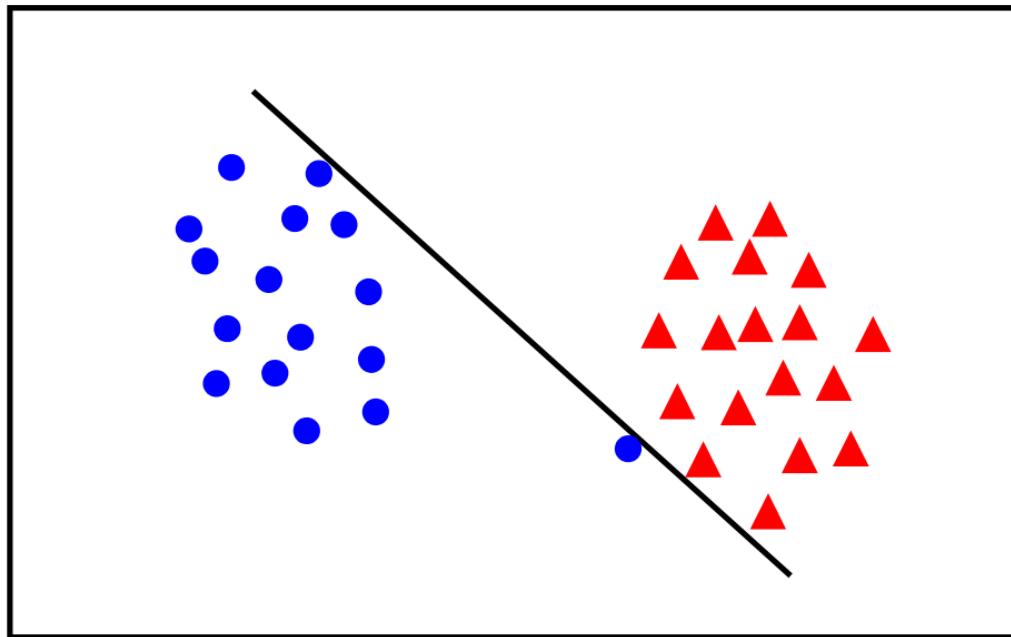
$$\max_{\mathbf{w}} \frac{2}{\|\mathbf{w}\|} \text{ subject to } \mathbf{w}^\top \mathbf{x}_i + b \begin{cases} \geq 1 & \text{if } y_i = +1 \\ \leq -1 & \text{if } y_i = -1 \end{cases} \text{ for } i = 1 \dots N$$

- Or equivalently

$$\min_{\mathbf{w}} \|\mathbf{w}\|^2 \text{ subject to } y_i (\mathbf{w}^\top \mathbf{x}_i + b) \geq 1 \text{ for } i = 1 \dots N$$

- This is a quadratic optimization problem subject to linear constraints and there is a unique minimum

Support Vector Machine



Support Vector Machine (Soft Margin Solution)

The optimization problem becomes

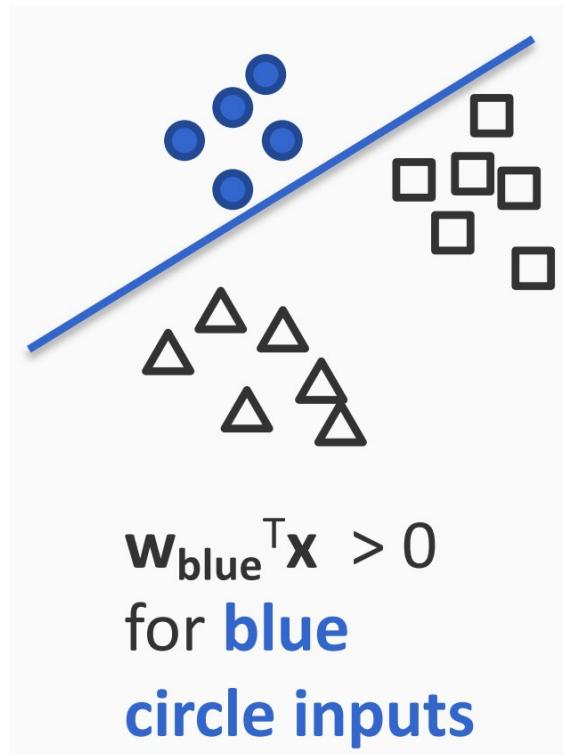
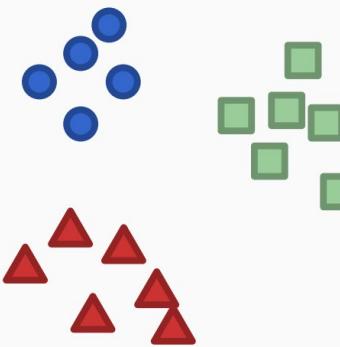
$$\min_{\mathbf{w} \in \mathbb{R}^d, \xi_i \in \mathbb{R}^+} \|\mathbf{w}\|^2 + C \sum_i^N \xi_i$$

subject to

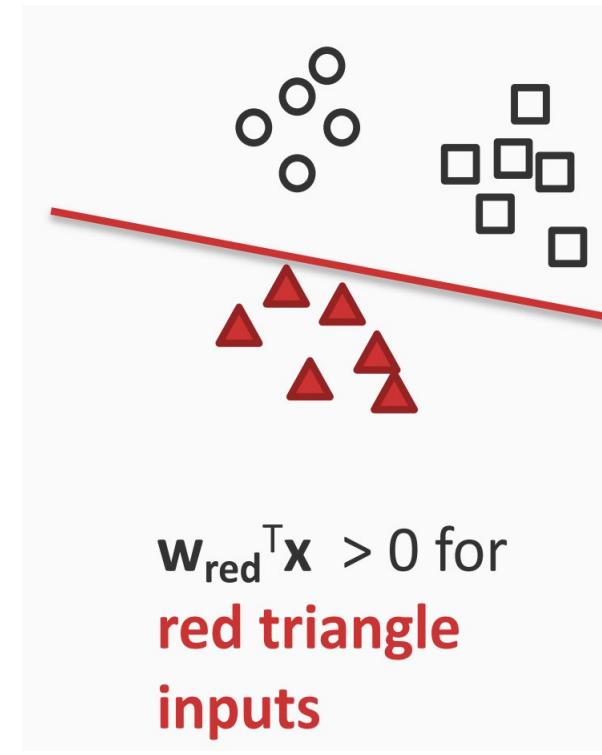
$$y_i (\mathbf{w}^\top \mathbf{x}_i + b) \geq 1 - \xi_i \text{ for } i = 1 \dots N$$

- Every constraint can be satisfied if ξ_i is sufficiently large
- C is a regularization parameter:
 - small C allows constraints to be easily ignored \rightarrow large margin
 - large C makes constraints hard to ignore \rightarrow narrow margin
 - $C = \infty$ enforces all constraints: hard margin

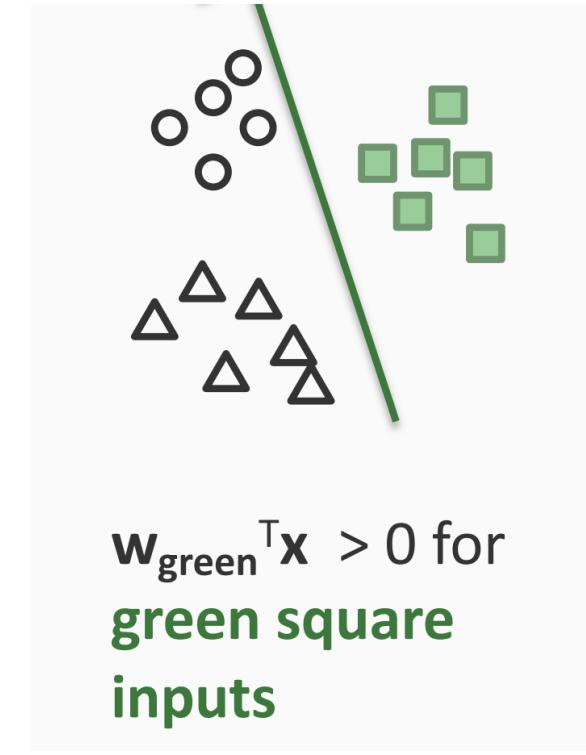
Support Vector Machine



$\mathbf{w}_{\text{blue}}^T \mathbf{x} > 0$
for **blue
circle inputs**



$\mathbf{w}_{\text{red}}^T \mathbf{x} > 0$ for
**red triangle
inputs**



$\mathbf{w}_{\text{green}}^T \mathbf{x} > 0$ for
**green square
inputs**

Support Vector Machine

