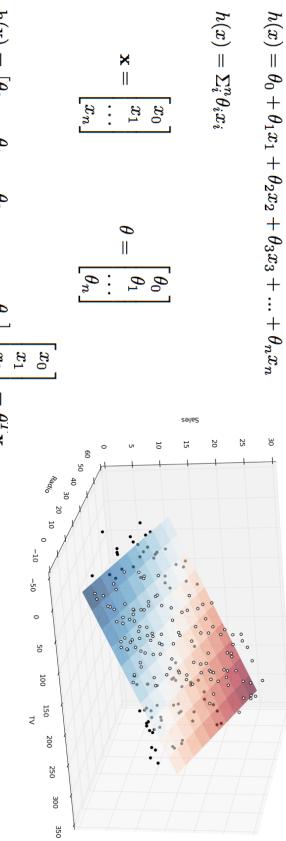


2

Linear Regression Model



$$h(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 + \dots + \theta_n x_n$$

$$h(x) = \sum_i^n \theta_i x_i$$

$$\mathbf{x} = \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix}$$

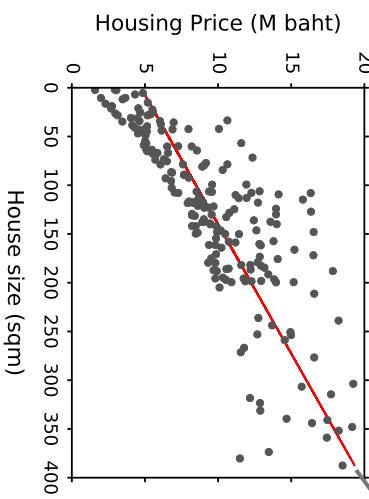
$$\boldsymbol{\theta} = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix}$$

$$\boldsymbol{\theta}^T \mathbf{x} = \boldsymbol{\theta}^T \mathbf{x}$$

Linear Regression Model

$$h(x) = \theta_0 + \theta_1 x$$

parameters
of the model



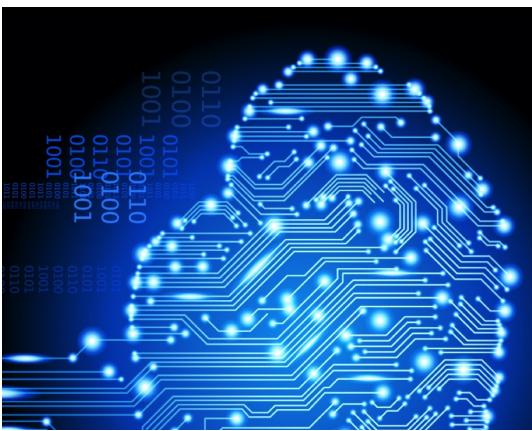
$$\theta_0 = 0.038$$

$$\theta_1 = 4.717$$

Regression, Classification & Clustering

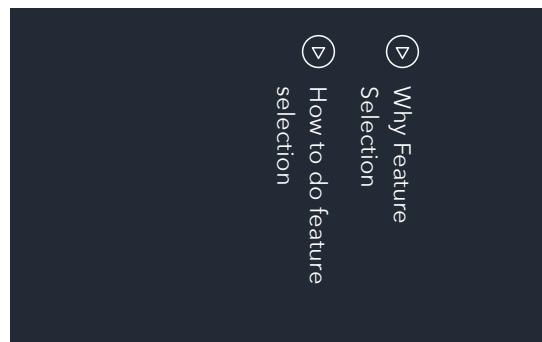
.....

Instructor: Warasinee Chaiangangnon, PhD



LINEAR & LOGISTIC REGRESSION

- ④ Why Feature Selection
- ④ How to do feature selection



3

4

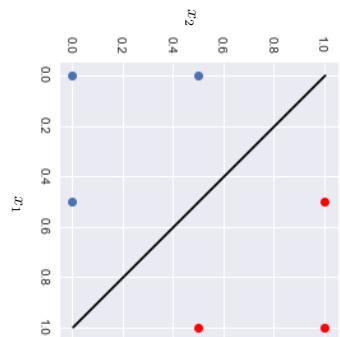
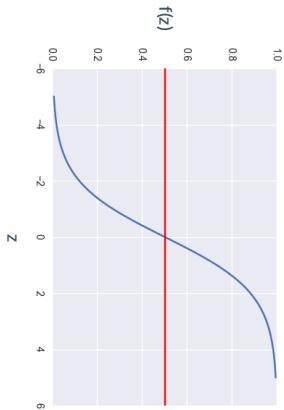
Logistic Regression Model

The Model

$$h(x) = f(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots)$$

$$f(z) = \frac{1}{1 + e^{-z}}$$

$$\begin{aligned} y' &= 1, \text{ if } f(z) > 0.5 \text{ or } z > 0 \\ y' &= 0, \text{ if } f(z) < 0.5 \text{ or } z < 0 \end{aligned}$$



model: $h(x) = f(\theta_0 + \theta_1 x_1 + \theta_2 x_2) = f(z)$
 parameter:
 $\theta_0 = -1$
 $\theta_1 = 1$
 $\theta_2 = 1$

x1	x2	z	y'
1	1		
0	0		
0.5	0		
0	0.5		

Regression Performance

- R-squared: how close the data are to the fitted regression line.
- R-squared = the percentage of the response variable variation that is explained by a linear model.

R-squared = Explained variation / Total variation

- R-squared = 0 (model explains none of the variability of data).
- R-squared = 1 (model explains all of the variability of data).

$$\begin{aligned} \rho_e &= \frac{1}{L} \sum_{l=1}^L \int_{\Omega_l} \frac{\Delta \Phi}{2\pi} = \frac{\Delta x}{2\pi} = \frac{x_2 - x_1}{2\pi} \\ k &= \frac{\Delta x}{\sqrt{2\pi}} \\ \rho_{eB} &= \frac{4\pi r^2}{k^2} X_L = \frac{U_m}{4\pi r^2} \frac{\varepsilon_m \varepsilon_B}{4\pi r^2} \\ \rho_{EB} &= \frac{|E_{PA} - E_{PB}|}{|E_{PA}|} = |P - P_B| \frac{1}{T} = \frac{4\pi r_1 r_2}{(m_1 + m_2)} \\ \rho_{EB} &= k \frac{Q}{\rho} \frac{M_m}{M_B} = \frac{Q}{N_A} \frac{M_m}{N_B} \\ \rho_{EB} &= N_A m_B = \frac{Q}{N_A} \frac{M_m}{N_B} \\ \ell_t &= \ell_0 (1 + d \Delta t) \quad I = \frac{U_e}{R + R_i} \cdot Q \cdot g \cdot L \cdot \frac{S_{ind}}{S_{ind}} \end{aligned}$$

LOGISTIC REGRESSION QUIZ



7

R-squared = Explained variation / Total variation

- R-squared = 0 (model explains none of the variability of data).
- R-squared = 1 (model explains all of the variability of data).

$$\begin{aligned} E &= \frac{1}{2} \hbar k \beta_n \quad \beta = \frac{\Delta I_c}{I_c} \quad \phi_e = \frac{\beta}{\pi \omega} \\ E &= \frac{1}{2} (\vec{E} \times \vec{B}) \quad \Delta I_c = \frac{\hbar k^2}{\pi \omega} \quad \phi = \frac{\beta}{\pi \omega} \\ E_f &= \frac{1}{2} \hbar k^2 \Gamma_{PC} = \frac{1}{2} \hbar k^2 \Gamma_{AU} \\ Q_f &= \frac{2m}{1} \frac{1}{\sigma_{\text{eff}}} \frac{\rho_{PC}}{\rho_{AU}} \quad M = \frac{Q_f}{\rho_{AU}} \end{aligned}$$

8

Logistic Regression Model

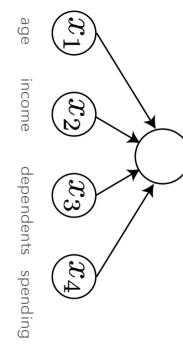
A simple example

10

Let's look at a simple logistic regression procedure.

$$h = f(\sum_j w_j x_j)$$

x1	x2	x3	x4	history
40	50	0	30	1
25	40	2	35	1
18	10	0	12	0
34	22	1	10	1



w1	0.7
w2	0.6
w3	-0.1
w4	-0.2

A simple example

11

Then fit the logistic regression to the data.

Suppose after fitting, here are the weight numbers.

$$h = f(\sum_j w_j x_j)$$

x1	x2	x3	x4	history
40	50	0	30	1
25	40	2	35	1
18	10	0	12	0
0.44	0.63	0	0.6	1
0.28	0.50	0.5	0.7	1
0.20	0.13	0	0.24	0
0.38	0.28	0.25	0.2	1

Logistic Regression Model

9

A simple example

12

We first need to do preprocessing, such as normalization and standardization.

model: $h(x) = f(\theta_0 + \theta_1 x_1 + \theta_2 x_2) = f(z)$

parameter:

$$\theta_0 = -1$$

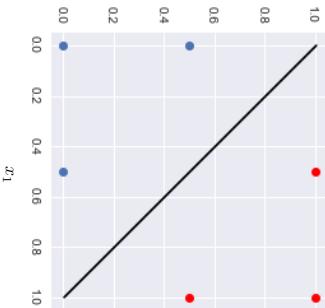
$$\theta_1 = 1$$

$$\theta_2 = 1$$

Logistic regression model draws a boundary at $z=0$

$$1 * x_1 + 1 * x_2 - 1 = 0$$

$$x_2 = -x_1 + 1$$



What the output means

14

Classification results

	X	W	X*W
age	0.44	0.7	0.31
income	0.63	0.6	0.38
dependent	0.00	-0.1	0.00
spending	0.60	-0.2	-0.12
		sum:	0.57
h:			0.64

h indicates the probability of customer being good



$$\begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ y_4 \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \\ 0 \\ 1 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}$$

h = 0.64
64% chance that he will be good
36% chance that he will be bad

$$E(W) = -\frac{1}{m} \sum_{c=1}^k \sum_{i=1}^m [y_c^i \log(h(x^i)_c) + (1 - y_c^i) \log(1 - h(x^i)_c)]$$

A simple example

13

Let us make prediction for a single customer...:

$$h = f(\sum_j w_j x_j)$$

	X	W	X*W
age	0.44	0.7	0.31
income	0.63	0.6	0.38
dependent	0.00	-0.1	0.00
spending	0.60	-0.2	-0.12
		sum:	0.57
h:			0.64

Regression Cost Function

Regression

For single output unit:

$$E(W) = \frac{1}{2m} \sum_{i=1}^n [y^i - h(x^i)]^2$$

For multiple output units:

$$E(W) = \frac{1}{2m} \sum_{c=1}^k \sum_{i=1}^m [y_c^i - h(x^i)_c]^2$$

Classification Cost Function

14

Classification problem with 2 classes : Binomial Log Loss

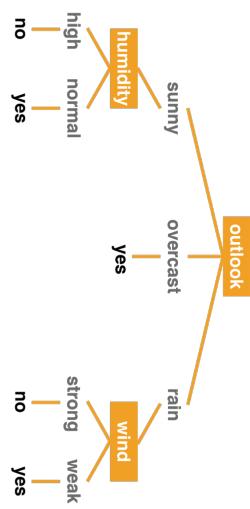
$$E(W) = -\frac{1}{m} \sum_{i=1}^m [y^i \log(h(x^i)) + (1 - y^i) \log(1 - h(x^i))]$$

Classification problem with more than 2 classes : Multinomial Log Loss

Decision Tree

- What is decision tree
- How it is constructed

DECISION TREE



- Here's an example of a decision tree for classification. Is Josh going to play tennis today?

Classification Performance



- Decision tree is a very old and simple idea. Today, it is perhaps the most popular machine learning algorithm due to the following reasons.

- Log Loss
- Accuracy
- Precision
- Recall
- F1-Measure

		Predicted class	
		<i>P</i>	<i>N</i>
Actual Class	<i>P</i>	True Positives (TP)	False Negatives (FN)
	<i>N</i>	False Positives (FP)	True Negatives (TN)

$REC = TPR = \frac{TP}{P} = \frac{TP}{TP + FN}$

$F_1 = 2 \cdot \frac{PRE \cdot REC}{PRE + REC}$



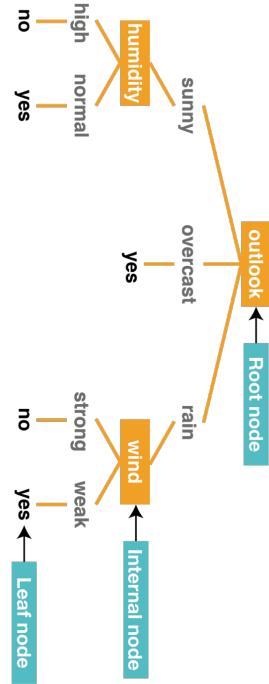
Decision Tree



- Decision tree is a very old and simple idea. Today, it is perhaps the most popular machine learning algorithm due to the following reasons.
 - 1. Easy to understand
 - 2. Easy to implement (not a lot of parameters to tweak)
 - 3. Can be used for both classification and regression problems

Decision Tree

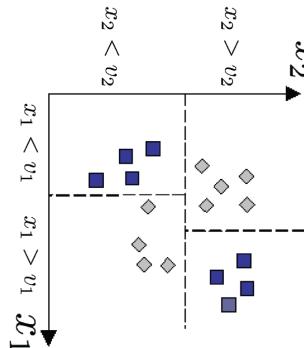
Root Node, Internal Node, Leaf Node



22

Decision Boundary

Decision trees divide the space into axis-parallel rectangles and label each rectangle with class membership.

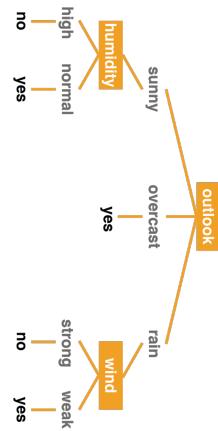


24

Decision Tree

21

- The y variable is {yes, no} There are 4 features:
 - x_1 = Outlook, {sunny, overcast, rain}
 - x_2 = Temperature, {hot, warm, cold}
 - x_3 = Humidity, {high, normal}
 - x_4 = Wind, {strong, weak}



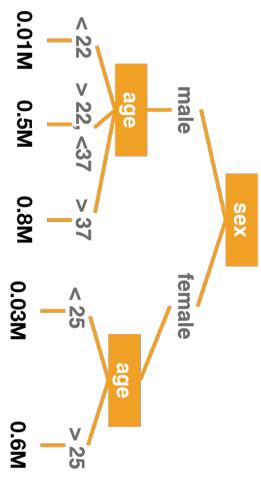
Decision Tree & Continuous Variable

23

Both x (features) and y (predicted value) can be continuous variables, for example...

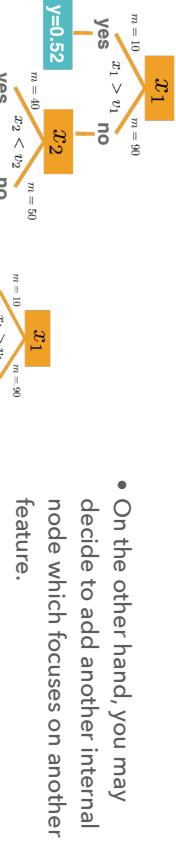
Predict amount of money in saving account of single people (Y) from the following attributes:

- x_1 = age (15-65 years old)
- x_2 = gender (male or female)
- x_3 = Humidity, {high, normal}
- x_4 = Wind, {strong, weak}



Constructing Decision Tree

26



- On the other hand, you may decide to add another internal node which focuses on another feature.
- Repeat the process until we account for all samples.



Constructing Decision Tree

25

One can construct a decision tree by a constructive search. Suppose we have a dataset with $m = 100$ samples.

- Decide on the feature at the root node and a boundary to create branches.



- Decide whether to add another internal node. If there's no use in creating another internal node, just add a leaf node and make a prediction.

Pseudo-Algorithm

28

```
Function: Build subtree  
Require: node  $n$ , data at the node  $D$   
 $(n_L, n_R, D_L, D_R) = \text{FindBestSplit}(D)$ 
```

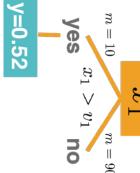
```
if StoppingCriteria( $D_L$ ) then  
    FindPrediction( $D_L$ )  
else  
    BuildSubtree( $n_L, D_L$ )  
end if  
  
if StoppingCriteria( $D_R$ ) then  
    FindPrediction( $D_R$ )  
else  
    BuildSubtree( $n_R, D_R$ )  
end if
```

Constructing Decision Tree

27

- At a given node, you always have to make a decision. Do I continue adding a node to the tree?

- if yes, which feature do I use to split the tree, at what threshold?
- if no, how do I make a prediction?



Information Gain

30

- Information gain is defined as:

$$IG(Y|X) = H(Y) - H(Y|X)$$

- The difference between the uncertainty of Y when we don't know X and when we know X.
- This is the information about Y that we gained from having known X.

Find the best split

29

Information Gain

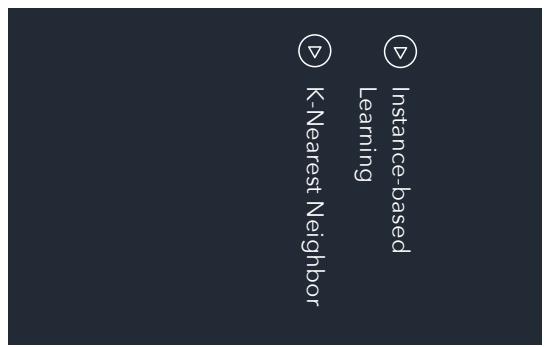
31

- Intuitively, the best split will send all the 'yes' samples to one side and 'no' samples to the other side.
- To determine the best split, we consider how much information about a given feature x tells us about the value of y.
- The best split involves the feature and associated threshold that has the most 'information gain'.

Stopping Criteria

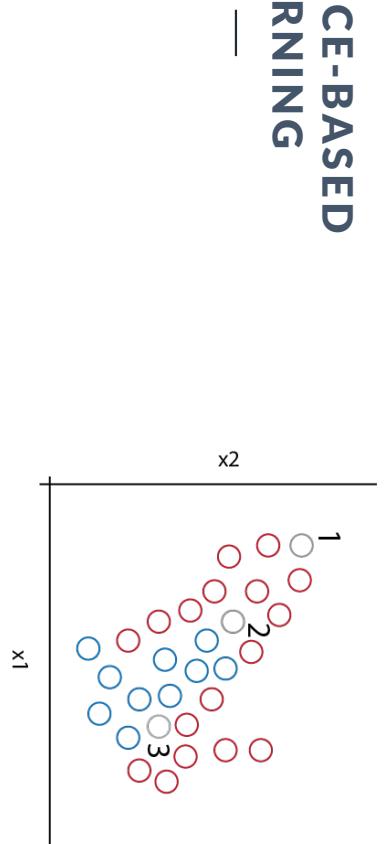
32

- When to decide not to make a branch?
1. When the leaf is pure, i.e. the variance of Y is small
 2. When the number of samples in the leaf is too small



INSTANCE-BASED LEARNING

Would you classify point 1, 2, 3 as blue or red. Fill in the table.



Pt	BLUE or RED
1	
2	
3	

34

Intuition Quiz

36

Making Prediction

33

- For a classification problem:
 - Predict most common y of the examples in the leaf.
 - For a regression problem:
 - Predict the average y of the examples in the leaf or build a linear regression model on the examples in the leaf.

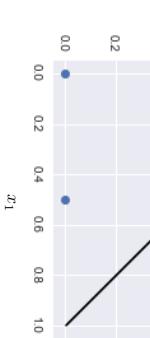
INTUITION QUIZ

35

Logistic Regression Model

38

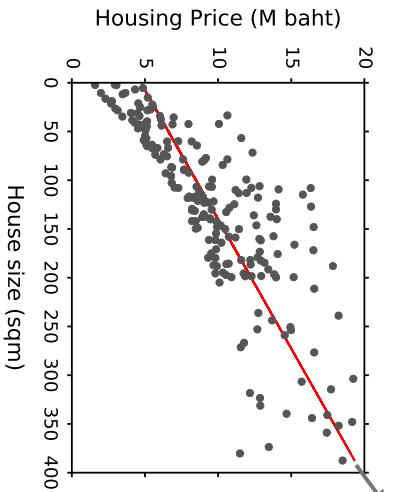
model: $h(x) = f(\theta_0 + \theta_1x_1 + \theta_2x_2) = f(z)$
 parameter:
 $\theta_0 = -1$
 $\theta_1 = 1$
 $\theta_2 = 1$



x1	x2	z	y'
1	1		
0	0		
0.5	0		
0	0.5		

Linear Regression Model

37



$$h(x) = \theta_0 + \theta_1 x$$

parameters
of the model

$$\theta_0 = 0.038$$

$$\theta_1 = 4.717$$

IBL: How Decision is Made

39

- Your source of knowledge is the similarity between two different data points. So you use similarity to make decisions such as classification and regression.
- You make decisions about one data point based on neighboring points.

Instance-based Learning

40

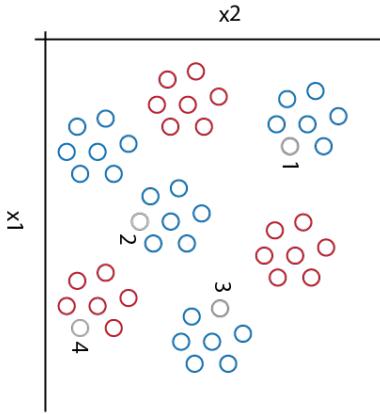
- Lazy algorithm: when you see your training set, you do nothing, just store them in the memory.
- When new sample comes you compare the new sample with the existing samples in the memory.
- Examples of algorithms in this family: nearest neighbor, kernel machines.

IBL - K-Nearest Neighbor Methods

- K-Nearest neighbor: locate k nearest neighbors around x' .

- For classification problem, let k neighbors vote for the right label of x' .

- For regression problem, average the y values of all neighbors and predict that y as the label of x' .



IBL - Nearest Neighbor Methods

- Nearest neighbor:
when you see a new data point (x'), locate the nearest data point (x) and predict the label of x' to be the same as label of x .

$$41) \quad \rho = \frac{1}{\sqrt{\pi} r^2} \int_{-\infty}^{\Delta x} \frac{dx}{2\pi} = \frac{\Delta x}{2\pi} = \frac{x_2 - x_1}{2\pi}$$

$$\frac{\rho}{\rho_{E_k}} = \frac{|E_{PA} - E_{PB}|}{|E_{PA} - E_{PB}|} = |k_B - k_A| / \Gamma = \frac{4\pi n_1 n_2}{(m_1 + m_2)}$$

$$m_1 = N_1 m_0 = \frac{Q}{N_1} M_0$$

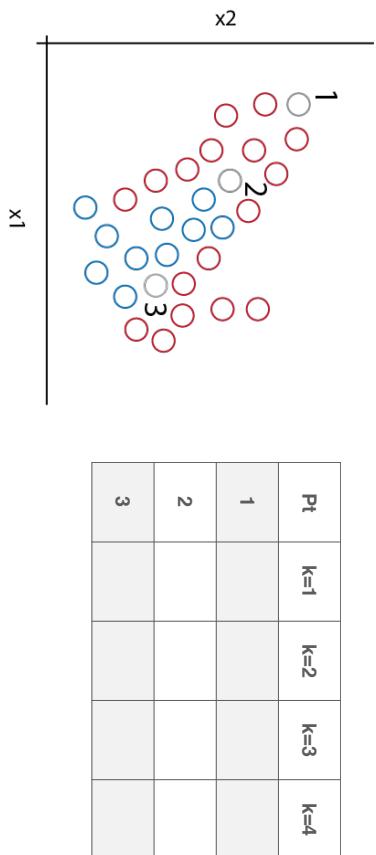
$$3) \quad \ell_0 = \lambda_0 (1 + d \Delta t) \quad I = \frac{U_0}{R + R_i} \frac{S_{Ind}}{S_{Ind}}$$

$$42) \quad F = M C$$

K-NEAREST NEIGHBOR QUIZ



43)



IBL - K-Nearest Neighbor Quiz

Use K-Nearest Neighbor Rule to classify point 1, 2, and 3 with different values of k.

44)

Distance and Similarity Metrics

46

- To determine whether two points are close, we use distance metrics.
- Distance metrics are the numerical value that tells you whether two points are close (low value) or far apart (high value).
- There are several ways to define distance metrics, such as euclidean distance, minkowski distance.
- Similarity metrics are the numerical value that tells you whether two points are close (very similar - high value) or far apart (very dissimilar - low value).
- Distance and similarity metrics are important in many ML models such as 'Support Vector Machine', 'K-Nearest Neighbor', 'K-Mean Clustering'

Pros and Cons

45

- Pros:
 - Training takes no time
 - Complex decision boundary is possible
 - Information is not lost
 - Cons:
 - Query is slow (the more data the slower)
 - Storage space is huge
 - Easily fooled by irrelevant attributes

Distance Metrics for Boolean Features

48

- Jaccard Distance

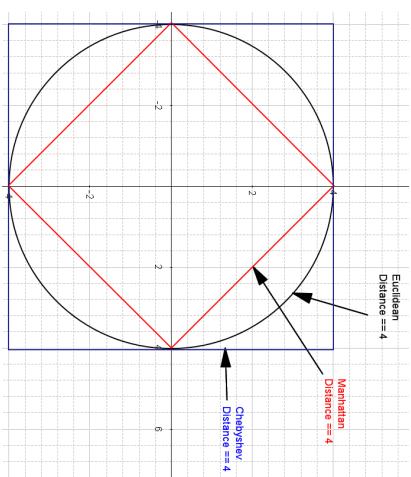
Feature	Me	My Dad
Man Barber	F	T
Toyota	T	T
MK	T	T
Water Park	T	F
Temple	F	T
Bar	F	F

N=6
 NTF=2
 NFT=1
 NFT=2
 NFF=1
 NNEQ : number of non-equal dimensions
 NNEQ = NTF + NFT = 3
 NNZ : number of nonzero dimensions
 NNZ = NTF + NFT + NTT = 5
 NNEQ / NNZ = 3/5 = 0.6

Distance Metrics for Real Value Features

47

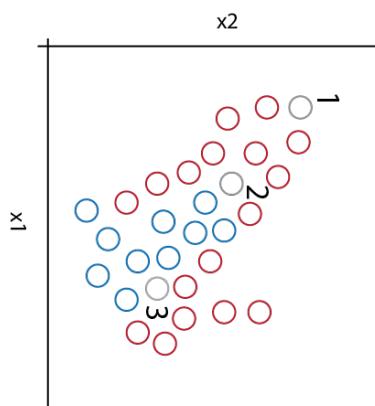
- Euclidean Distance
- $\sqrt{\sum((x - x')^2)}$
- Manhattan Distance
- $\sum(|x - x'|)$



How to Avoid Overfitting

50

- Use k as an overfitting control.
 - If k is one, you are very susceptible to noise (overfitting).
 - If k is high, you are averaging over really large regions, you lose resolution (under fitting).



- Unsupervised Learning
- K-Means Algorithm
- Clustering Performance

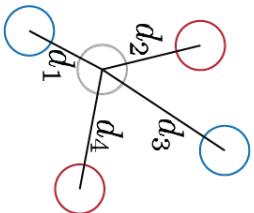
K-MEANS CLUSTERING

Using Distances as Weights

49

- Neighbors who are closer to the target data point should get more say in the voting process.

$$y' = \frac{w_1 y_1 + w_2 y_2 + w_3 y_3 + w_4 y_4}{w_1 + w_2 + w_3 + w_4}$$



$$w_i = \frac{1}{d_i}$$

How to Avoid Overfitting

51

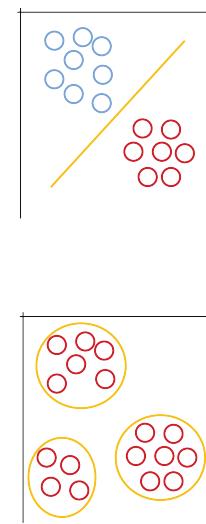
- Remove noisy instances prior to using nearest neighbor algorithm. Remove x if all nearest neighbors of x are in the opposite class.
- Form prototypes. If you observe lots and lots of very similar samples, lump them into a prototype by finding an average over all dimensions.

$$y' = \frac{w_1 y_1 + w_2 y_2 + w_3 y_3 + w_4 y_4}{w_1 + w_2 + w_3 + w_4}$$

52

Unsupervised Learning

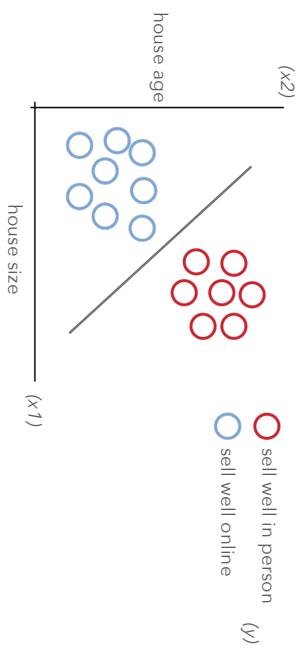
In **supervised learning**, your data come with labels indicating what class corresponds to each sample. Sometimes, data do not come with categorical labels, but you can tell that there is a grouping structure.



Supervised vs. Unsupervised Learning

- Classification: predicting discrete output from input

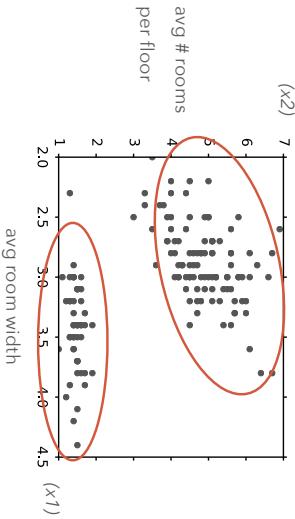
- Example: predicting whether a house would sell well in person or online



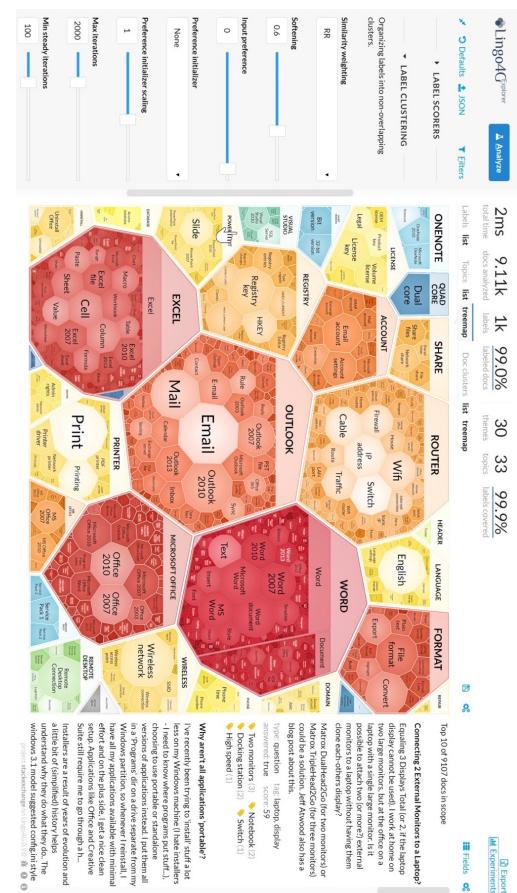
Unsupervised Learning

- Clustering: grouping unlabeled input by similarities

- Example: grouping houses in the database



55

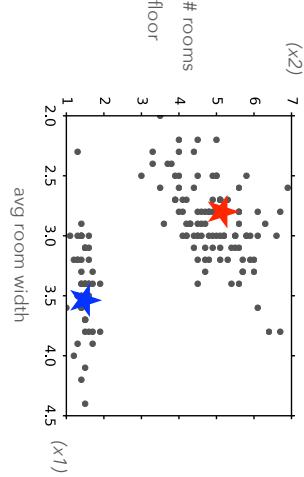


K-Means Clustering

58

- K-means clustering is the most popular clustering algorithm.

- K-means clustering aims to partition n observations into k clusters in which each observation belongs to the cluster with the nearest mean, serving as a prototype of the cluster.



Supervised vs. Unsupervised Learning

57

- Supervised learning problem: "given all history of feature data points (X) and their labels (Y), find a model that computes label (y') for a given sample data point (x'). The prediction must minimize a cost function (minimize errors)."
- Clustering problem: "given all feature data points (X), find their labels (Y). The label Y must minimize a given objective function."

- The red and blue stars are called "centroids".
 - The red star is the centroid of top cluster.
 - The blue star is the centroid of the bottom cluster.
- The location of the centroid is at the mean of all points in the cluster.

	X	Y	x'	y'	Predict	Cost Function
Supervised	Yes	Yes	No	y'	Error	
Unsupervised	Yes	No	Yes	No	Y, y'	Objective F.

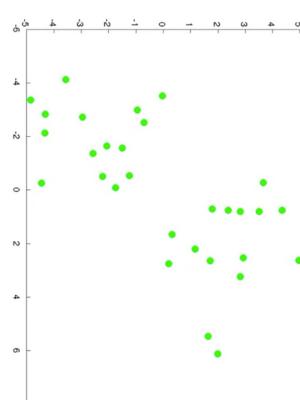
K-Means Clustering

59

K-Means Algorithm

60

First, we have a bunch of unlabeled data points. We decide that we are going to find two clusters in this data.



- The red and blue stars are called "centroids".
 - The red star is the centroid of top cluster.
 - The blue star is the centroid of the bottom cluster.
- The location of the centroid is at the mean of all points in the cluster.

- The red and blue stars are called "centroids".
 - The red star is the centroid of top cluster.
 - The blue star is the centroid of the bottom cluster.
- The location of the centroid is at the mean of all points in the cluster.

- The red and blue stars are called "centroids".
 - The red star is the centroid of top cluster.
 - The blue star is the centroid of the bottom cluster.
- The location of the centroid is at the mean of all points in the cluster.

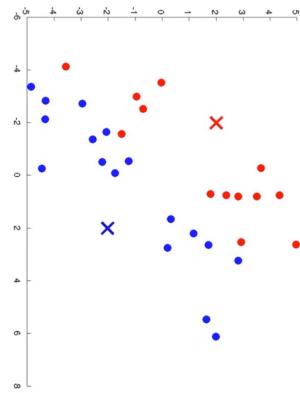
- The red and blue stars are called "centroids".
 - The red star is the centroid of top cluster.
 - The blue star is the centroid of the bottom cluster.
- The location of the centroid is at the mean of all points in the cluster.

- The red and blue stars are called "centroids".
 - The red star is the centroid of top cluster.
 - The blue star is the centroid of the bottom cluster.
- The location of the centroid is at the mean of all points in the cluster.

K-Means Algorithm

62

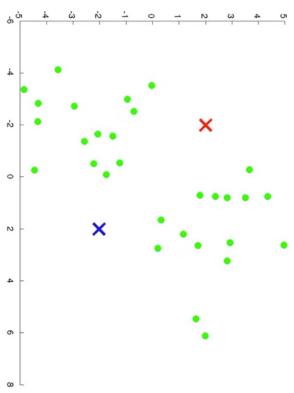
Then, we do **cluster assignment** step, where we go through every sample and determine whether each dot is closer to red or blue centroid. Label the sample to red/blue accordingly.



K-Means Algorithm

63

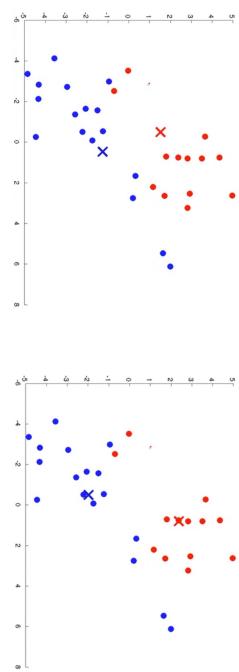
The first step is to pick two random locations to be our cluster centroids.



K-Means Algorithm

64

After the assignment step, we do **centroid movement** step where we move the red and blue centroids to the means of our clusters.



K-Means Algorithm

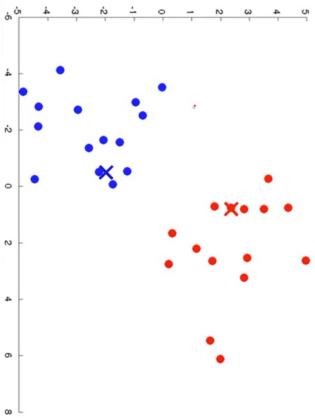
64

Repeat the **cluster assignment** step and **centroid movement** step, alternately.

K-Means Algorithm

66

- To summarize, k-means algorithm has two steps:
 - In **cluster assignment step**, we fixed the centroids and label each data point to belong to the nearest cluster centroid.
 - In **centroid movement step**, we fixed data point labels and move each centroid to the mean of its data points.
- Iterating over these two steps and you will achieve the best clusters.

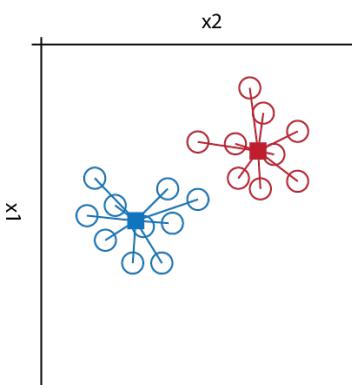


K-Means Algorithm

65

The algorithm converges to the solution when cluster centroids are not changed anymore.

- What K-Means algorithm is doing?
- The algorithm minimize the total distances between all data points and the centroids of the clusters they belong to.



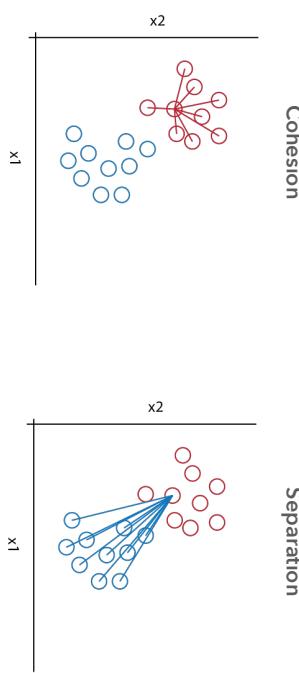
K-Means Algorithm

67

Clustering Performance

68

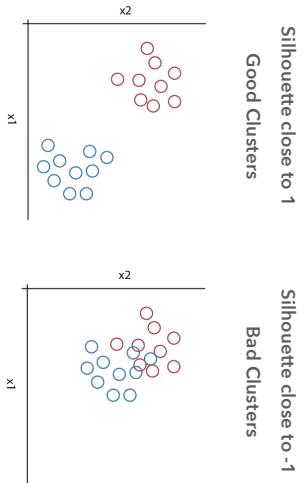
- How good is your clustering results? It is not as simple as counting right and wrong predictions. Performance of clustering algorithm is often judged based on how well you algorithm separates out the data into several clusters.



Clustering Performance

70

- Silhouette Coefficient = $1 - (a/b)$
- If $b >> a$, SC is close to 1.
- If $a << b$, SC is close to -1.
- Completeness: all members of a given class are assigned to the same cluster.
- Both homogeneity and completeness are bounded below by 0.0 and above by 1.0 (higher is better)
- V-Measure Score: the harmonic mean between homogeneity and completeness.



$$V = 2 * (\text{homogeneity} * \text{completeness}) / (\text{homogeneity} + \text{completeness})$$

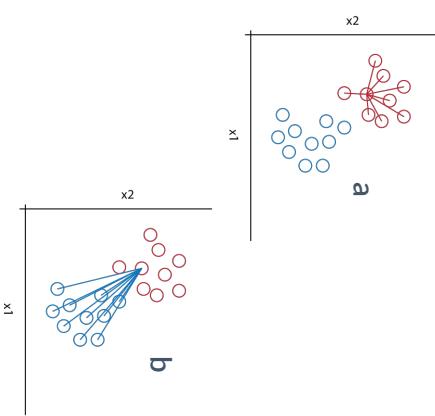
Clustering Performance

69

- Silhouette Coefficient is one such measure. It's defined by two separate scores:
 - a: mean distance between a sample and all other points in the same class.
 - b: mean distance between a sample and all other points in the next nearest cluster.
- Silhouette Coefficient

$$= \frac{(b-a)}{\max(b,a)}$$

$$= 1 - \frac{(a/b)}{} - \text{if } b > a$$



Clustering Performance

71

- Other methods of clustering performance include Calinski-Harabaz Index (dispersion from centroid), percentage of explained variance (between-group variance / total variance), and many others.
- We can calculate these performance metrics using any distance metrics we learned about (Euclidean, Manhattan, Jaccard).
- Typically, it would be best if we could find some ground truth to validate our clusters. For example, if we cluster 1M news articles, we could get people to label 5000 news articles to establish ground truth for validating our clusters.

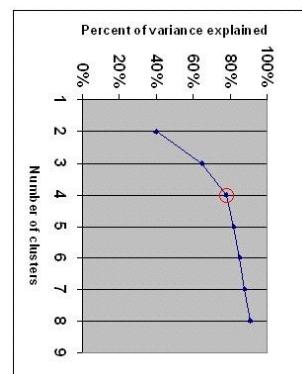
Clustering Performance

72

Choosing the Right K

74

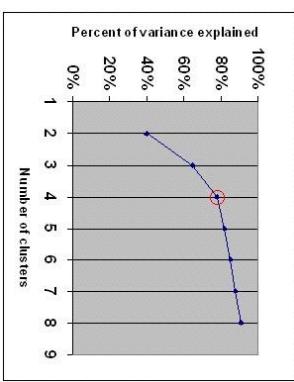
- Pick the k at the elbow point. At this point, more clusters do not necessarily mean higher performance.
- Similarly we use metrics like Information Criterion (the famous one being Bayesian - BIC, and Akaike - AIC). These are measures that find a good balance between the variance explained and the complexity of the model.



Choosing the Right K

73

- How to determine the number of clusters in the dataset?
- The most common way is called elbow methods.
- You compute performance metrics such as Silhouette Coefficient or Percentage of Variance (in the figure), while varying k.



Choosing the Right K

75

- Pick the k at the elbow point. At this point, more clusters do not necessarily mean higher performance.
- Similarly we use metrics like Information Criterion (the famous one being Bayesian - BIC, and Akaike - AIC). These are measures that find a good balance between the variance explained and the complexity of the model.

