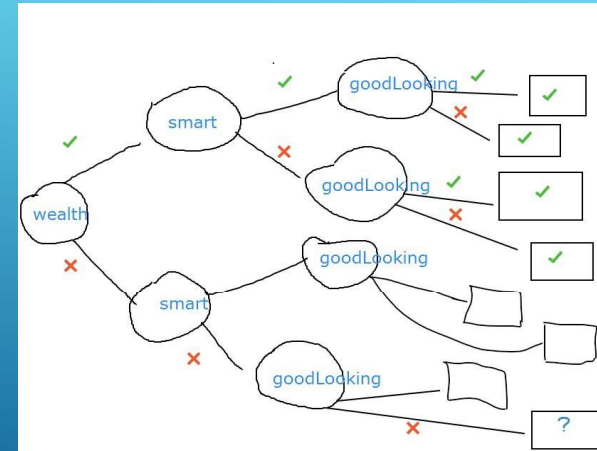


CHAPTER 7 DECISION TREE

PART I: DECISION TREE CLASSIFIERS
PART II: DECISION TREE ENSEMBLES

Associate Professor Yachai Limpiyakorn, Ph.D



GENERAL DECISION MAKING TREE

2

2110773-7 2/2567

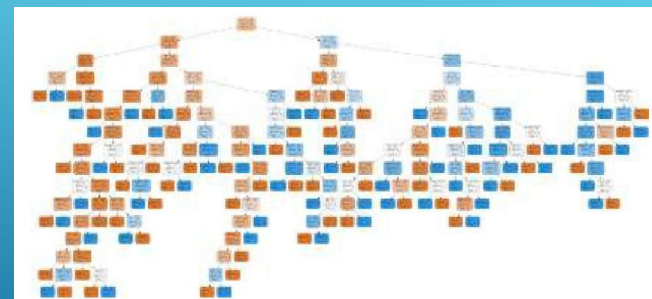
PART 1: DECISION TREE CLASSIFIERS

- A tree-like diagram illustrates all possible decision alternatives and the corresponding outcomes.
- Starting from the root of a tree,
 - ❖ **internal** node represents the basis on which a decision is made;
 - ❖ each **branch** of a node represents how a choice may lead to the next nodes;
 - ❖ terminal node, **leaf**, represents the outcome produced.
 - ❖ Paths from root to leaves represent **classification rules**

3

2110773-7 2/2567

1. Tree Construction



2. Tree Pruning

DECISION TREE LEARNING

4

2110773-7 2/2567

ประเภท

คุณสมบัติ

คำ

Name	Hair	Height	Weight	Lotion	Result
Sarah	blonde	average	light	no	sunburned
Dana	blonde	tall	average	yes	none
Alex	brown	short	average	yes	none
Annie	blonde	short	average	no	sunburned
Emily	red	average	heavy	no	sunburned
Pete	brown	tall	heavy	no	none
John	brown	average	heavy	no	none
Katie	blonde	short	light	yes	none

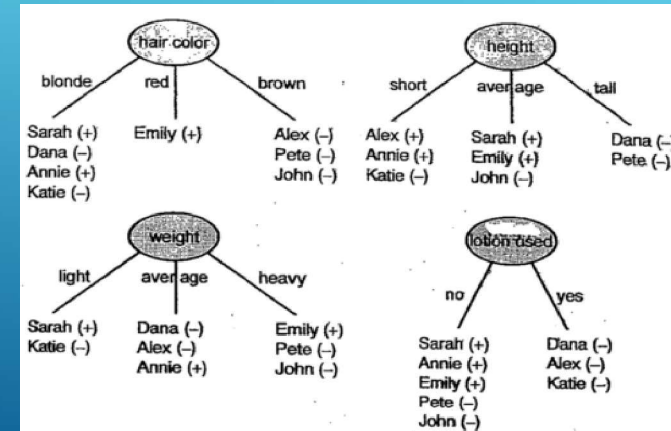
- สร้างต้นไม้จำนวนเท่ากับค่าคุณลักษณะ และให้แต่ละคุณลักษณะเป็นโหนดราก
- สำหรับต้นไม้แต่ละต้น แบ่งตัวอย่างสอนไปแต่ละกิ่งตามค่าคุณลักษณะของโหนดราก
- สำหรับต้นไม้แต่ละต้น ให้สร้างโหนดภายในด้วยคุณลักษณะที่เหลือ
- ดังตัวอย่างการสร้างต้นไม้ที่มี "hair color" เป็นโหนดราก
- งานทำเช่นนี้ไปเรื่อยๆ จนกระทั่งสร้างต้นไม้ตัดสินใจที่เป็นไปได้ครบหมดทุกต้น

5

2110773-7 2/2567

OPTIMAL DECISION TREE

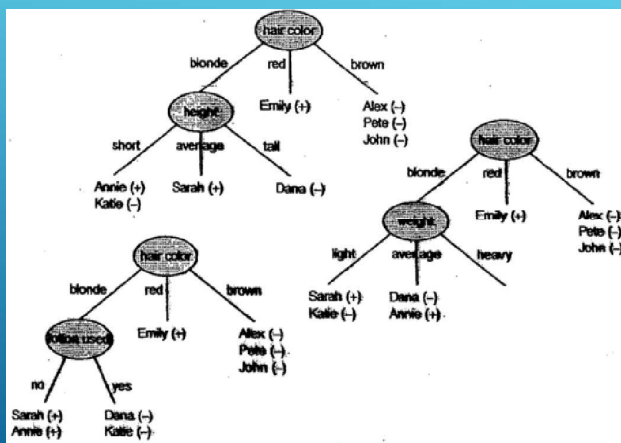
สร้างต้นไม้ที่ประกอบด้วยโหนดราก 4 แบบที่เป็นไปได้



6

2110773-7 2/2567

โหนดราก HAIR COLOR และโหนดภายใน 3 แบบที่เป็นไปได้



7

2110773-7 2/2567

Rain	Windy	Xtreme windy	Decision
0	0	0	Nothing
0	0	1	Nothing
0	1	0	Nothing
0	1	1	Nothing
1	0	0	Umbrella
1	0	1	Umbrella
1	1	0	Rain jacket
1	1	1	Stay home

TREE CONSTRUCTION

8

2110773-7 2/2567



Occam's Razor

สมมติฐานที่สั้นกว่าที่สามารถอธิบายข้อมูลได้เหมือนกัน
จะเป็นสมมติฐานที่ดีกว่า

the simplest explanation is usually the best one

~raining → don't bring anything

raining and ~windy → use an umbrella

raining and ~extremely windy → wear a rain jacket

raining and extremely windy → stay home

<https://medium.com/@ml.at.berkeley/machine-learning-crash-course-part-5-decision-trees-and-ensemble-models-dcc5a36af8cd>

9

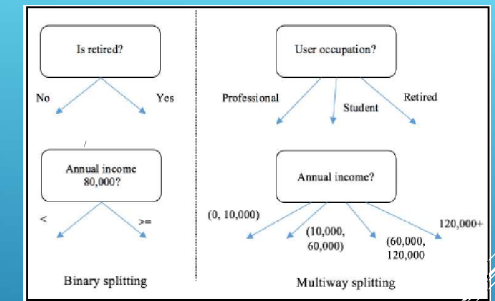
2110773-7 2/2567

Select the best attribute using Attribute Selection Measures(ASM) to split the records.

Make that attribute a decision node and breaks the dataset into smaller subsets.

Starts tree building by repeating this process recursively for each child until one of the condition will match:

- All the tuples belong to the same class.
- There are no more remaining attributes.
- There are no more instances.



BASIC IDEA OF DECISION TREE ALGORITHM (TOP-DOWN CONSTRUCTION)

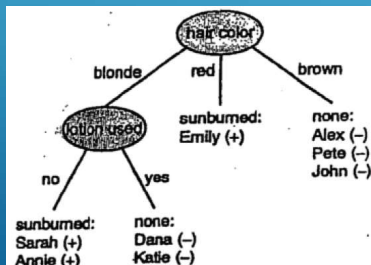
10

2110773-7 2/2567

Multi-way Split

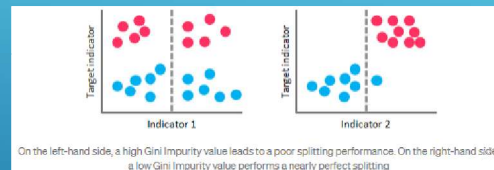
► Information Gain – ID3 [Ross Quinlan]

► Gain ratio – C4.5 [Ross Quinlan]



Binary Split

► GINI – CART (Classification and Regression Tree)



SPLIT MEASURE/ ASM

11

2110773-7 2/2567

กำหนด message M ประกอบด้วยค่าที่เป็นไปได้

$\{m_1, m_2, \dots, m_n\}$ และ

ความน่าจะเป็นที่จะเกิดค่า $m_i = P(m_i)$ จะได้ว่า

จำนวนบิตน้อยที่สุดที่ใช้ encode m_i แต่ละตัว ที่ให้

ค่าเฉลี่ยจำนวนบิตที่น้อยที่สุด คือ

$$\text{Optimal code length } (m_i) = -\log_2 P(m_i)$$

ค่าสารสนเทศของ M หรือค่าเอนโทรปีของ M เขียนแทนด้วย $I(M)$ คำนวณโดย

$$I(M) = \sum_i^n -P(m_i) \log_2 P(m_i)$$

ENTROPY/ INFORMATION

Message	Probability	Standard Code	Optimal Code
A	0.5	00	0
B	0.25	01	10
C	0.125	10	110
D	0.125	11	111
Average Encoding Length		2 bits	1.75 bits

Average Encoding Length of Optimal Code is calculated by
 $= (-0.5 \log_2 0.5) + (-0.25 \log_2 0.25) + (-0.125 \log_2 0.125) + (-0.125 \log_2 0.125)$
 $= (0.5 \times 1) + (0.25 \times 2) + (0.125 \times 3) + (0.125 \times 3) = 1.75 \text{ bits}$

12

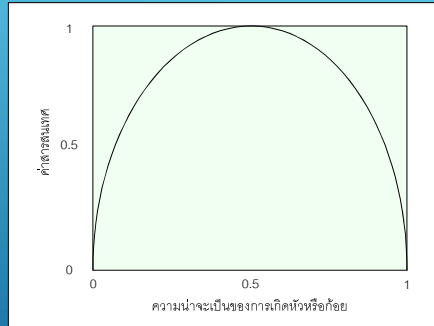
2110773-7 2/2567

INFORMATION/ ENTROPY OF COIN FLIP

$$\log_2 1 = 0$$

$$I(\text{การโยนหัวโยนก้อย}) = -P(\text{หัว}) \log_2(P(\text{หัว})) - P(\text{ก้อย}) \log_2(P(\text{ก้อย}))$$

- ▶ M1=HHHHHHHH
 $I(M1) = (-1 \log_2 1) + (-0 \log_2 0) = 0$
- ▶ M2=TTTTTTT
 $I(M2) = (-0 \log_2 0) + (-1 \log_2 1) = 0$
- ▶ M3=HHHHTTTT
 $I(M3) = (-0.5 \log_2 0.5) + (-0.5 \log_2 0.5) = 1$



Lower entropy implies a purer dataset. In a perfect case where the dataset contains only one class, the entropy is $-(1 \cdot \log_2 1 + 0) = 0$

13

2110773-7 2/2567

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	yes
31...40	high	yes	fair	yes
>40	medium	no	excellent	no

ID3 [J.R. QUINLAN]

Id3Estimator

14

2110773-7 2/2567

ID3 CONSTRUCTION (1)

$$S = [9+, 5-]$$

$$I(S) = \frac{-9}{14} \log_2 \frac{9}{14} + \left(\frac{-5}{14} \log_2 \frac{5}{14} \right) = 0.940$$

Age

LE 30 GT 40

S1=[2+, 3-] [31..40] S3=[3+, 2-]

$$I(S1) = \frac{-2}{5} \log_2 \frac{2}{5} + \left(\frac{-3}{5} \log_2 \frac{3}{5} \right)$$

$$S2 = [4+, 0-]$$

$$I(S2) = \frac{-4}{4} \log_2 \frac{4}{4} + (-0 \log_2 0)$$

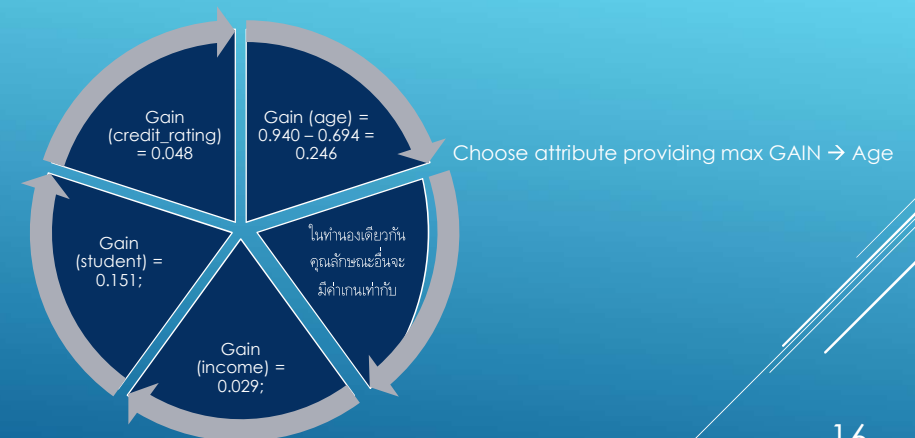
$$I(S3) = \frac{-3}{5} \log_2 \frac{3}{5} + \left(\frac{-2}{5} \log_2 \frac{2}{5} \right)$$

$$I_{\text{age}}(T) = \frac{5}{14} \left(\frac{-2}{5} \log_2 \frac{2}{5} - \frac{3}{5} \log_2 \frac{3}{5} \right) + \frac{4}{14} (0) + \frac{5}{14} \left(\frac{-3}{5} \log_2 \frac{3}{5} - \frac{2}{5} \log_2 \frac{2}{5} \right) = 0.694$$

15

2110773-7 2/2567

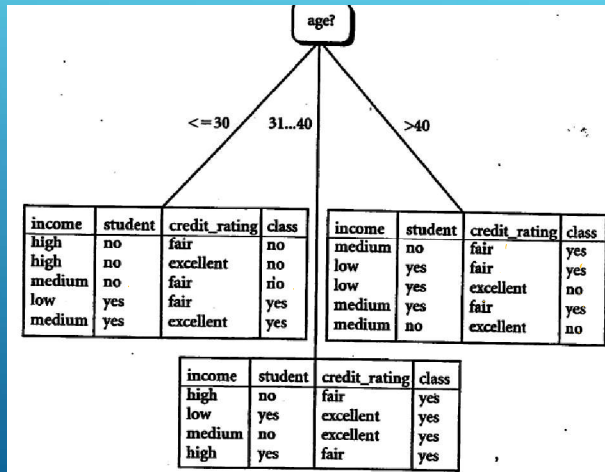
ID3 CONSTRUCTION (2)



16

2110773-7 2/2567

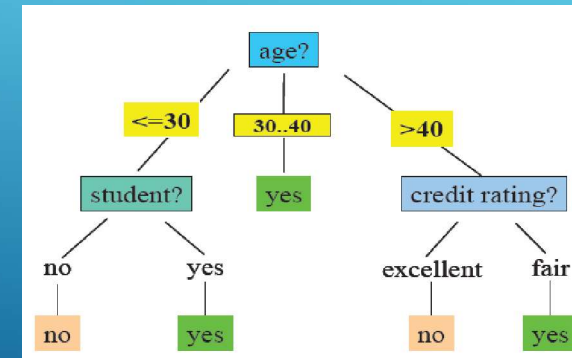
ID3 CONSTRUCTION (3)



17

2110773-7 2/2567

OUTPUT OF ID3 LEARNING (MULTI-WAY SPLIT)



18

2110773-7 2/2567

Gini Impurity

- Lower Gini indicates a purer dataset
- For a dataset with K classes, suppose data from class k ($1 \leq k \leq K$), take up a fraction f_k ($0 \leq f_k \leq 1$) of the entire dataset:

$$GiniImpurity = 1 - \sum_{k=1}^K f_k^2$$

- To evaluate quality of a split, add up the Gini of all resulting subgroups, combining the proportions of each subgroup as corresponding weight factors.
- The smaller the weighted sum of Gini Impurity, the better the split.

User gender	Interested in tech	Click	Group by gender
M	True	1	Group 1
F	False	0	Group 2
F	True	1	Group 2
M	False	0	Group 1
M	False	1	Group 1

User gender	Interested in tech	Click	Group by interest
M	True	1	Group 1
F	False	0	Group 2
F	True	1	Group 1
M	False	0	Group 2
M	False	1	Group 2

#1 split based on gender

#2 split based on interest in tech

Weighted Gini Impurity of #1 split based on gender

$$\#1 \text{ Gini Impurity} = \frac{3}{5} \left[1 - \left(\frac{2^2}{3} + \frac{1^2}{3} \right) \right] + \frac{2}{5} \left[1 - \left(\frac{1^2}{2} + \frac{1^2}{2} \right) \right] = 0.467$$

Weighted Gini Impurity of #2 split based on tech interest

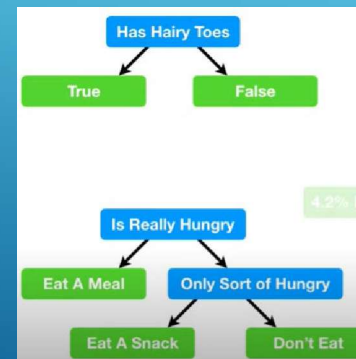
$$\#2 \text{ Gini Impurity} = \frac{2}{5} \left[1 - (1^2 + 0^2) \right] + \frac{3}{5} \left[1 - \left(\frac{1^2}{3} + \frac{2^2}{3} \right) \right] = 0.267$$

19

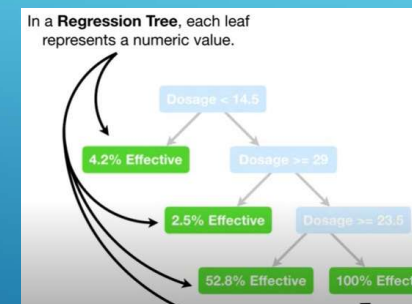
2110773-7 2/2567

SPLIT METRICS: GINI INDEX

CART (Classification and Regression Tree)



Classification Tree



Regression Tree

20

2110773-7 2/2567

- ▶ Most ML learning models in Python work with numerical data
- ▶ Three approaches to manage categorical data:
 - ▶ Drop categorical variables if NOT relevant
 - ▶ Label encoding or ranking in case of ordinal variables
 - ▶ One-Hot encoding

Label	Encoded Label
Africa	1
Asia	2
Europe	3
South America	4
North America	5
Other	6

Label encoding

	is_africa	is_asia	is_europe	is_sam	is_nam
Africa	1	0	0	0	0
Asia	0	1	0	0	0
Europe	0	0	1	0	0
South America	0	0	0	1	0
North America	0	0	0	0	1
Other	0	0	0	0	0

One-hot encoding

DATA PREPROCESSING

21

2110773-7 2/2567

ข้อดี

- ▶ ช่วยให้การเรียนรู้ง่ายขึ้น แทนที่ไม่ได้จะเรียนรู้ Pattern เปลี่ยนมาเรียนรู้จากสวิตช์
- ▶ ความหมายของข้อมูลแบบ Nominal จะตรงขึ้น Europe 3 ไม่ได้ใกล้เคียง S.America 4 มากกว่า N.America 5
- ▶ สามารถ Dot Product กับ Matrix ที่ต้องการ

ข้อเสีย

- ▶ ความหมายของลำดับข้อมูลแบบ Ordinal จะหายไป เนื่องจากทุก Category แตกต่างกันไปหมด
- ▶ ถ้าข้อมูลมี Value หลากหลายมาก เช่น มีสีเสื้อ 10,000 สี จะทำให้มีปัญหาเรื่องเนื้อที่/Memory ที่เก็บค่า 0 เป็นส่วนใหญ่ เรียกว่า Sparse Matrix
- ▶ การเพิ่ม Category ใหม่ ยับย่อยอยู่ตลอด จะทำให้มีปัญหา เช่น เพิ่มสีเสื้อใหม่

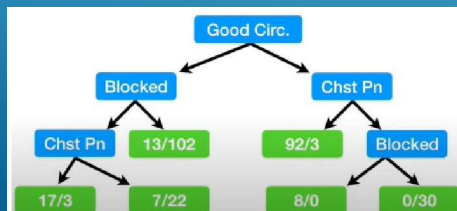
ONE-HOT ENCODING

22

2110773-7 2/2567

- ▶ Always produces **binary** splits
- ▶ **Gini index**. A Gini score of 0 indicates perfect purity and a score of 1 indicates maximum impurity.
- ▶ CART should be allowed to go till 7–8 tree depth in accordance with the nature of producing tall and skinny trees.
- ▶ Splitting stops when CART detects no further gain can be made, or some pre-set stopping rules are met. (Alternatively, the data are split as much as possible and then the tree is later pruned).
- ▶ The optimal Tree is identified by evaluating the performance of every Tree through test set; or performing k-fold cross-validation.

CART ALGORITHM



23

2110773-7 2/2567

- ▶ $Gini(interest, tech) = weighted_impurity([1, 1, 0], [0, 0, 0, 1]) = 0.405$
- ▶ $Gini(interest, Fashion) = weighted_impurity([0, 0], [1, 0, 1, 0, 1]) = 0.343$
- ▶ $Gini(interest, Sports) = weighted_impurity([0, 1], [1, 0, 0, 1, 0]) = 0.486$
- ▶ $Gini(occupation, professional) = weighted_impurity([0, 0, 1, 0], [1, 0, 1]) = 0.405$
- ▶ $Gini(occupation, student) = weighted_impurity([1, 0, 0, 1], [0, 0, 1]) = 0.476$
- ▶ $Gini(occupation, retired) = weighted_impurity([1, 0, 0, 0, 1, 1], [0]) = 0.429$

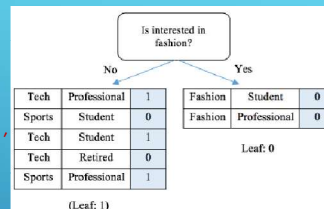
User interest	User occupation	Click
Tech	Professional	1
Fashion	Student	0
Fashion	Professional	0
Sports	Student	0
Tech	Student	1
Tech	Retired	0
Sports	Professional	1

IMPLEMENTING A CART TREE (1)

24

2110773-7 2/2567

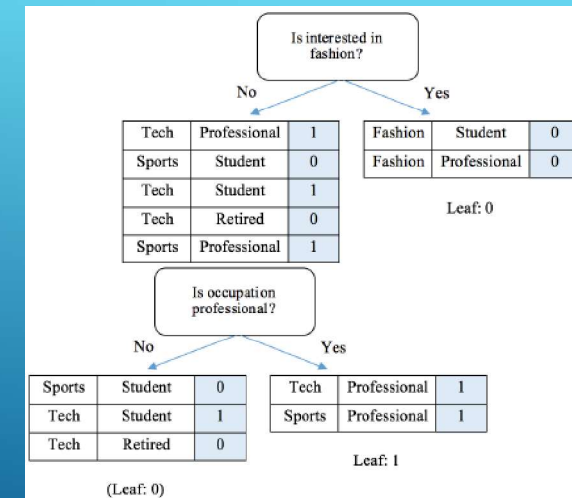
- ▶ $\text{Gini}(\text{interest, tech}) = \text{weighted_impurity}([0, 1], [1, 1, 0]) = 0.467$
- ▶ $\text{Gini}(\text{interest, Sports}) = \text{weighted_impurity}([1, 1, 0], [0, 1]) = 0.467$
- ▶ $\text{Gini}(\text{occupation, professional}) = \text{weighted_impurity}([0, 1, 0], [1, 1]) = 0.267$
- ▶ $\text{Gini}(\text{occupation, student}) = \text{weighted_impurity}([1, 0, 1], [0, 1]) = 0.467$
- ▶ $\text{Gini}(\text{occupation, retired}) = \text{weighted_impurity}([1, 0, 1, 1], [0]) = 0.300$



IMPLEMENTING A CART TREE (2)

25

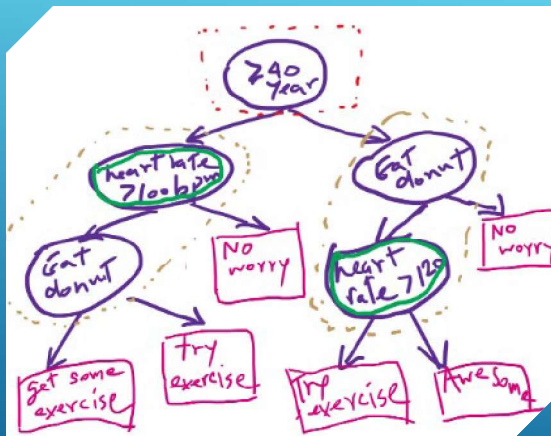
2110773-7 2/2567



IMPLEMENTING A CART TREE (3)

26

2110773-7 2/2567



EXPLORATION

27

2110773-7 2/2567

Weight	Heart Disease	Weight	Heart Disease
220	Yes	Lowest 155	No
180	Yes	180	Yes
225	Yes	190	No
190	No	220	Yes
155	No	Highest 225	Yes

Step 1) Sort the patients by weight, lowest to highest.

NUMERIC VARIABLE: WHAT'S THE BEST WEIGHT USED TO DIVIDE THE PATIENT?

28

2110773-7 2/2567

Step 2) Calculate the average weight for all adjacent patients.

Step 3) Calculate the impurity values for each average weight.

Weight	Heart Disease
155	No
167.5	Yes
180	Yes
185	No
190	Yes
205	Yes
222.5	Yes
225	Yes

Weight	Heart Disease	Gini impurity
155	No	
167.5	Yes	Gini impurity = 0.3
180	Yes	
185	No	Gini impurity = 0.47
190	No	
205	Yes	Gini impurity = 0.27
220	Yes	
222.5	Yes	Gini impurity = 0.4
225	Yes	

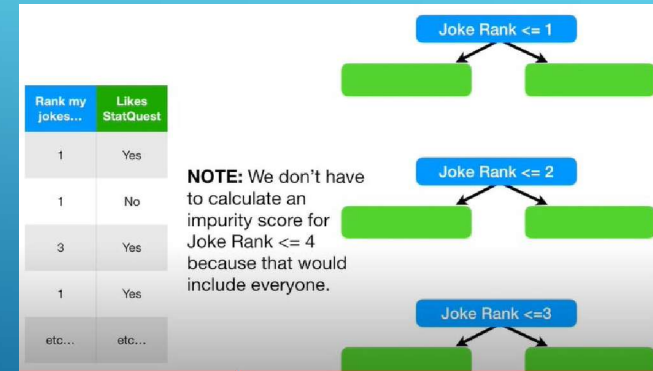
The lowest impurity occurs when we separate using **weight < 205...**

...so this is the cutoff and impurity value we will use when we compare weight to chest pain or blocked arteries.

NUMERIC VARIABLE: WHAT'S THE BEST WEIGHT USED TO DIVIDE THE PATIENT?

29

2110773-7 2/2567



ORDINAL VARIABLE

30

2110773-7 2/2567

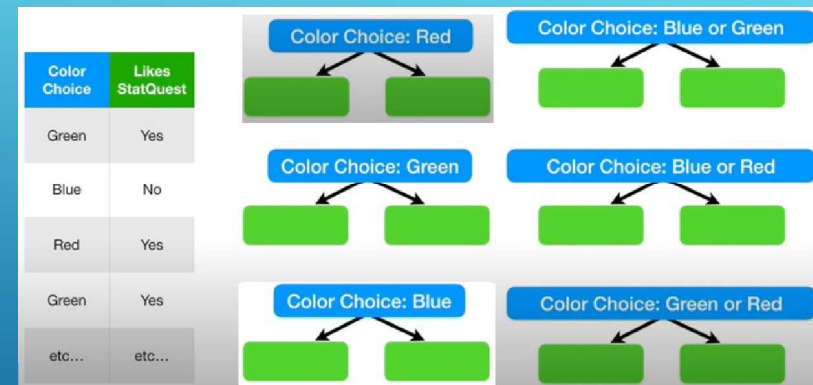
Color Choice	Likes StatQuest
Green	Yes
Blue	No
Red	Yes
Green	Yes
etc...	etc...

When there are **multiple choices**, like "color choice can be blue, green or red", you calculate an impurity score for each one as well as each possible combination.

NOMINAL VARIABLE (1)

31

2110773-7 2/2567



NOTE: We don't have to calculate an impurity score for "Color Choice: Blue or Green or Red" since that...

NOMINAL VARIABLE (2)

32

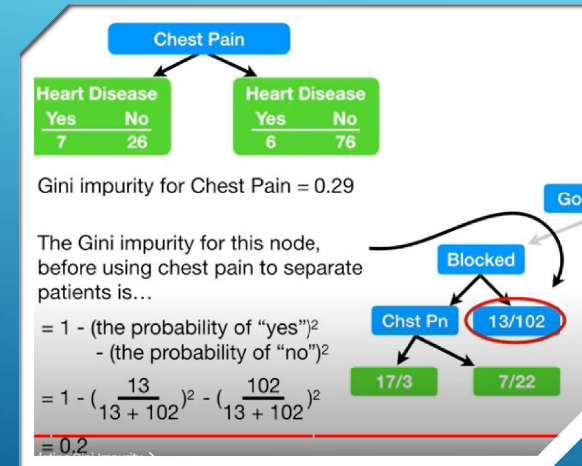
2110773-7 2/2567

- ▶ Pruning is a technique used to deal with overfitting, that reduces the size of DTs by removing sections of the Tree that provide little predictive or classification power.
- ▶ The goal is to reduce complexity and gain better accuracy by reducing the effects of overfitting and removing sections of the DT that may be based on noisy or erroneous data.
- ▶ There are two different strategies to perform pruning on DTs:
 - Pre-prune: When you stop growing DT branches when information becomes unreliable.
 - Post-prune: When you take a fully grown DT and then remove leaf nodes only if it results in a better model performance. This way, you stop removing nodes when no further improvements can be made.

TREE PRUNING

33

2110773-7 2/2567



STOP SPLITTING WHEN NO FURTHER GAIN CAN BE MADE

Chest Pain	Good Blood Circulation	Blocked Arteries	Heart Disease
No	No	No	No
Yes	Yes	Yes	Yes
Yes	Yes	No	No
Yes	No	???	Yes
etc...	etc...	etc...	etc...

34

2110773-7 2/2567

- ▶ Optimization of DT classifier performed by only pre-pruning using maximum depth of DT.
- ▶ **max_depth : int or None, (default=None) or Maximum Depth of a Tree:** If None, nodes are expanded until all the leaves contain less than min_samples_split samples. The higher value of maximum depth causes overfitting, and a lower value causes underfitting.

DECISION TREE (CLASSIFICATION) USING SCIKIT-LEARN

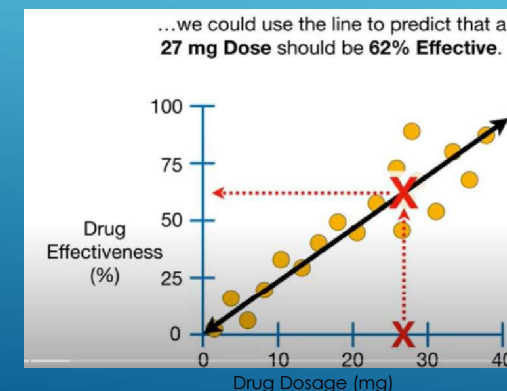
35

<https://www.datacamp.com/community/tutorials/decision-tree-classification-python>

2110773-7 2/2567

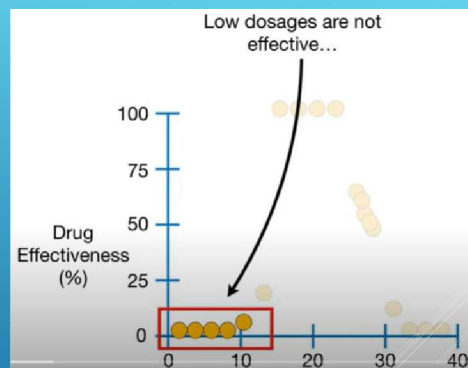
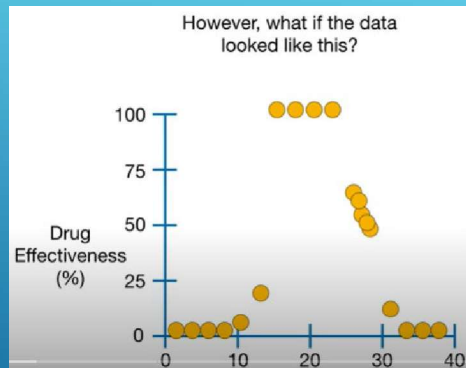
LINEAR REGRESSION

Easily fit a line to the data, the higher the dose, the more effective the drug...



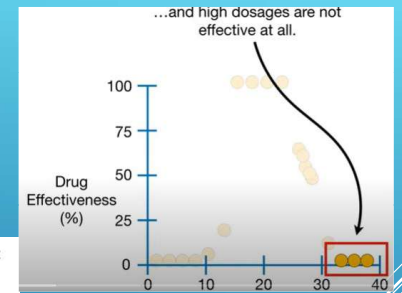
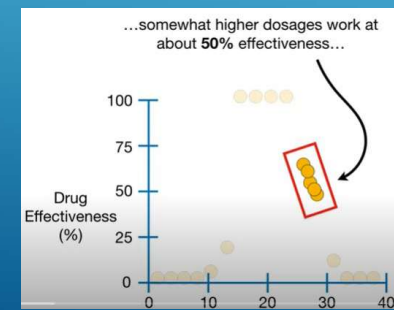
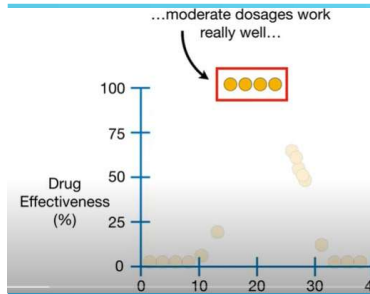
36

2110773-7 2/2567



37

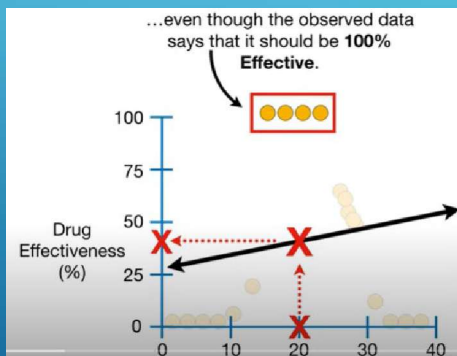
2110773-7 2/2567



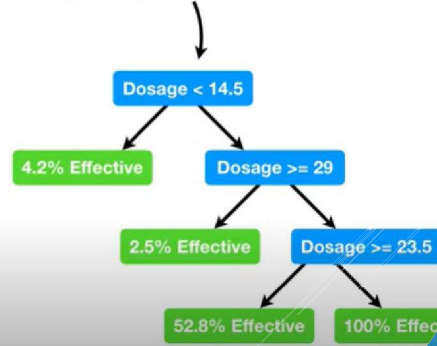
38

2110773-7 2/2567

For example, if someone told us they were taking a 20 mg Dose...

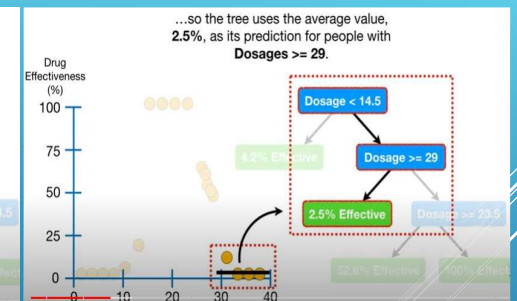
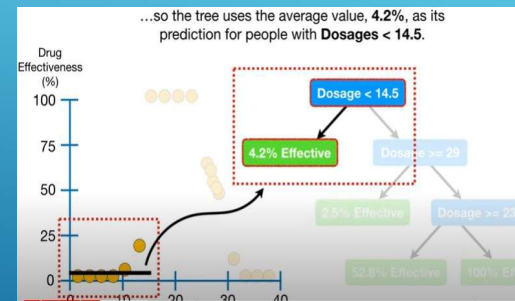


One option is to use a **Regression Tree**.



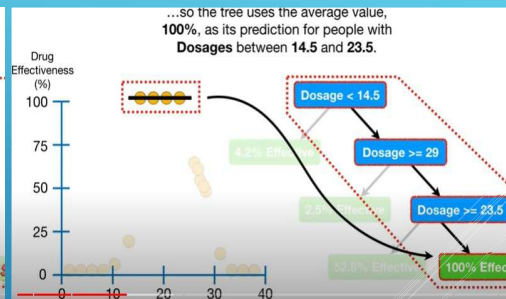
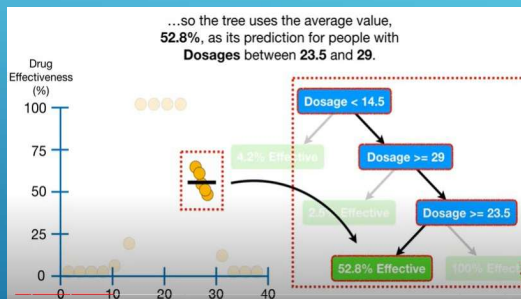
39

2110773-7 2/2567



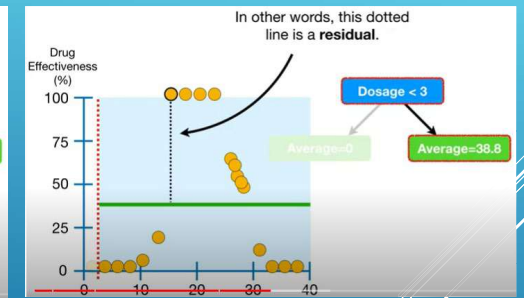
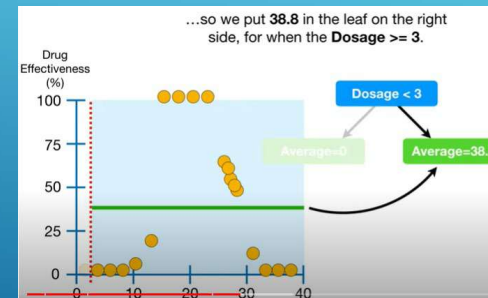
40

2110773-7 2/2567



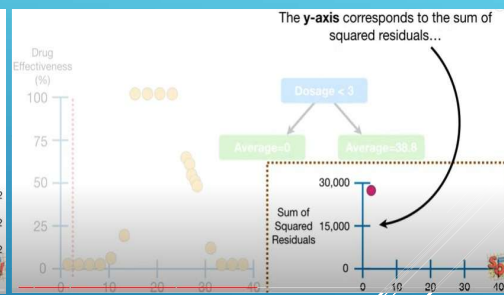
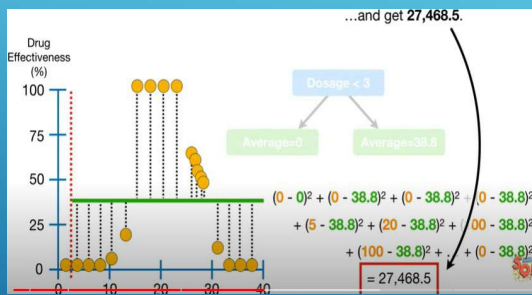
41

2110773-7 2/2567



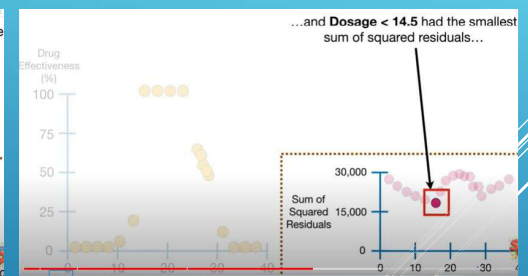
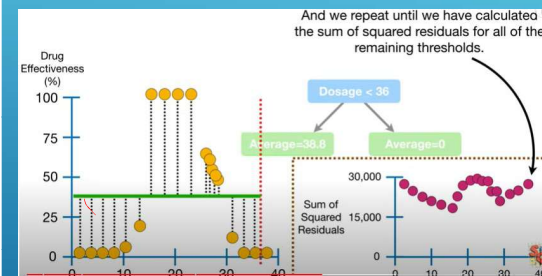
42

2110773-7 2/2567



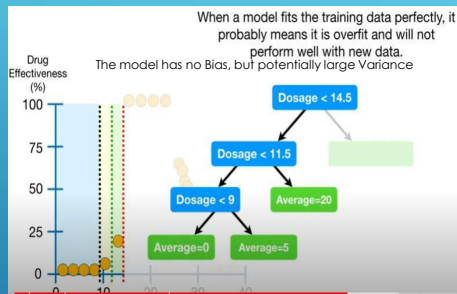
43

2110773-7 2/2567



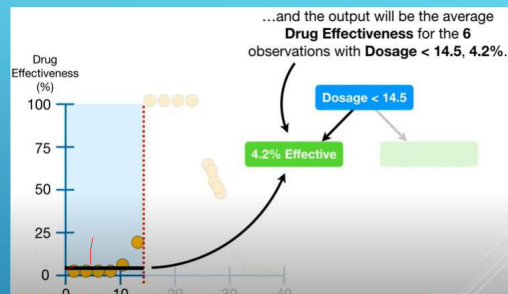
44

2110773-7 2/2567



The simplest is to only split observations when there are more than some minimum number.

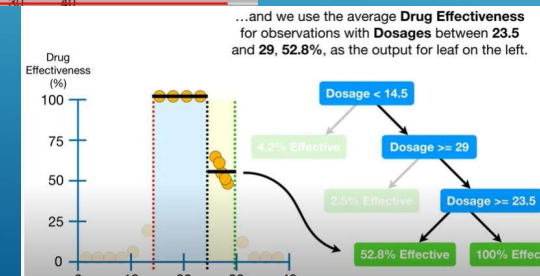
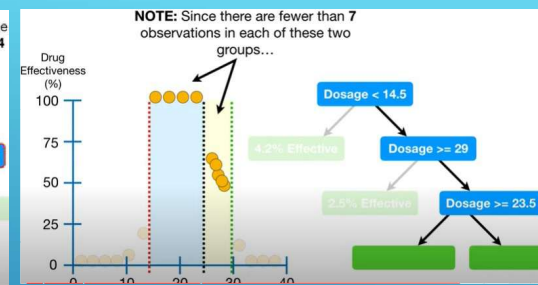
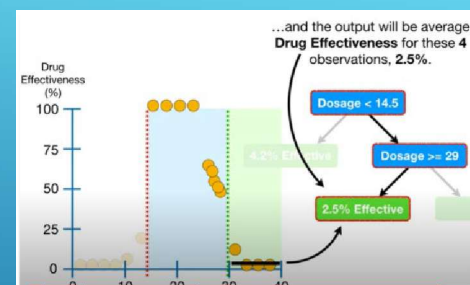
Typically, the minimum number of observations to allow for a split is 20.



However, since this example doesn't have many observations, I set the minimum to 7.

45

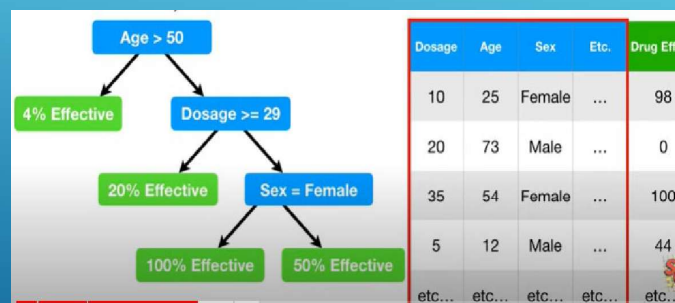
2110773-7 2/2567



46

2110773-7 2/2567

- ▶ A Regression tree looks for splits that minimize the Least Square Deviation (LSD), sometimes referred as "variance reduction", that implies the variance within the node.



REGRESSION TREE

38

2110773-7 2/2567

Pro

- ▶ easy to interpret and visualize
- ▶ easily capture Non-linear patterns with non-parametric nature of the algorithm.
- ▶ requires fewer data preprocessing, no need to normalize features
- ▶ can be applied for variable selection

Con

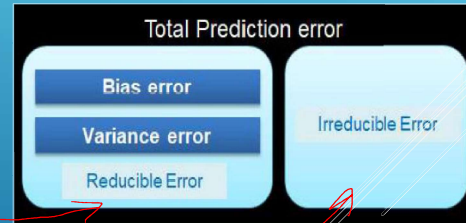
- ▶ Sensitive to noisy data. It can overfit noisy data.
- ▶ Biased with imbalanced dataset, balance out the dataset before creating DT is recommended
- ▶ small variation(or variance) in data can result in different DT. This can be reduced by bagging and boosting algorithms.

DT CLASSIFIER

48

2110773-7 2/2567

- ▶ Every model has both bias and variance error components in addition to white noise.
- ▶ The ideal model will have both low bias and low variance.
- ▶ Unfortunately, bias and variance are inversely related to each other; while trying to reduce one component, the other component of the model will increase.
- ▶ The true art lies in creating a good fit by balancing both.



$$E(y_0 - \hat{f}(x_0))^2 = \text{Var}(\hat{f}(x_0)) + [\text{Bias}(\hat{f}(x_0))]^2 + \text{Var}(\epsilon)$$

49

2110773-7 2/2567

Bias-error

- ▶ Difference between predicted and actual data points caused by **oversimplified** model or unable to capture underlying pattern of data.
- ▶ It misses how the features in the training data set relate to the expected output.
- ▶ A model with high bias is too simple and has low number of predictors.

Variance-error

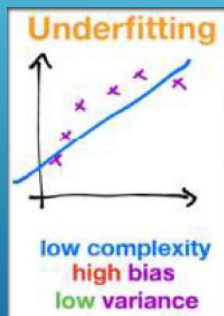
- ▶ High variance error of a model implies that it is highly sensitive to small fluctuations. This model flounders outside of its comfort zone(training data)
- ▶ Any model which has very large number of predictors will end up being a very **complex model** which will deliver very accurate predictions for the training data that it has seen already but this complexity makes the generalization of this model to unseen data very difficult, i.e. a high variance model. Thus, this model will perform very poorly on test data.

REDUCIBLE ERROR/ INADVERTENT MISTAKES

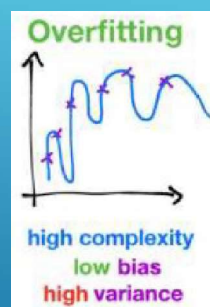
50

2110773-7 2/2567

Bias-error



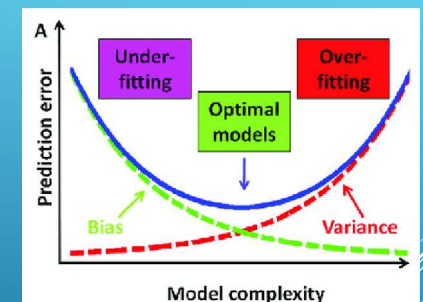
Variance-error



51

2110773-7 2/2567

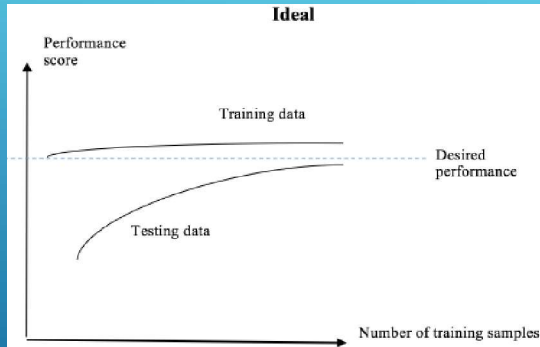
- ▶ On the one hand, we want our algorithm to model the training data very closely, otherwise we'll miss relevant features and interesting trends.
- ▶ However, on the other hand we don't want our model to fit too closely, and risk overinterpreting every outlier and irregularity.
- ▶ **High-Bias:** Suggests more assumptions about the form of the target function.
- ▶ **High Variance:** Suggests large changes to the estimate of the target function with changes to the training dataset.



BIAS-VARIANCE DILEMMA

52

2110773-7 2/2567

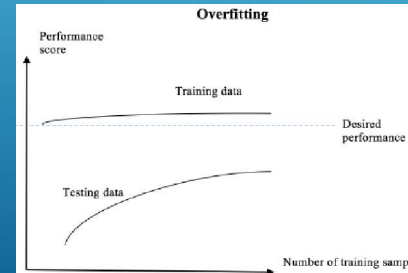


- ▶ A **learning curve** is usually used to evaluate the bias and variance of a model.
- ▶ For a model that fits well on the training samples, the performance of training samples should be above desire. Ideally, as the number of training samples increases, the model performance on testing samples improves; eventually the performance on testing samples becomes close to that on training samples.

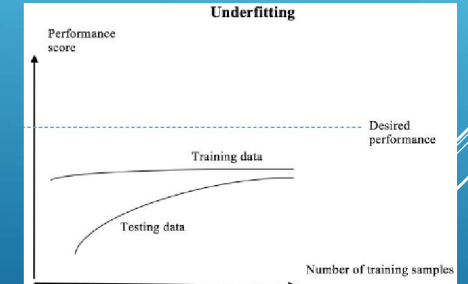
DIAGNOSING OVERFITTING AND UNDERFITTING

53

2110773-7 2/2567



When the performance on testing samples converges at a value far from the performance on training samples, overfitting can be concluded. In this case, the model fails to generalize to instances that are not seen.



For a model that does not even fit well on the training samples, underfitting is easily spotted: both performances on training and testing samples are below desire in the learning curve.

54

2110773-7 2/2567