

USER USER Collaborative filtering

Collaborative filtering:

A hasn't watched All the movies
eg. in this case.

$A=B \Rightarrow$ Similar Ratings

$A+C \Rightarrow$ const dissimilar cases

DATE	HM	HD	HP3	DW	SW	2/3
A	4			5	1	
B	5	5	4			
C				2	4	5
D		3				3

Jaccard Similarity

$$\text{Sim}(A, B) = \frac{|R_A \cap R_B|}{|R_A \cup R_B|} = \frac{1}{5}$$

$$\text{Sim}(A, C) = \frac{2}{4}$$

$$\therefore \text{Sim}(A, B) < \text{Sim}(A, C)$$

just sees A, B have watched lesser movies than A, C have in Common & not if they both liked it.
Problem

Cosine Similarity

> Insert unknown values : 0

$$\text{Sim}(A, B) = \frac{\sum A \cdot B}{\sqrt{A^2} \sqrt{B^2}} = \frac{4 \cdot 5}{\sqrt{4^2 + 5^2 + 1^2} \sqrt{5^2 + 5^2 + 4^2}} = 0.38$$

$$\text{Sim}(A, C) = 0.32 \leftarrow \text{Cosine Lbl both}$$

$\therefore \text{Sim}(A, B) > \text{Sim}(A, C)$ but not by much.

problem: treats missing ratings as least

Centered Cosine

Pearson Correlation

Rating Normalization

$$A : (4+5+1) / 3 = 10/3$$

$$B : (5+5+4) / 3 = 14/3$$

$$C : (2+4+5) / 3 = 11/3$$

$$D : (3+3) / 2 = 6/2$$

Subtract this value from each cell.

$$\text{i.e. } 4 - \frac{10}{3} = \frac{2}{3}, \quad 5 - \frac{10}{3} = \frac{5}{3}, \quad 1 - \frac{10}{3} = -\frac{7}{3}$$

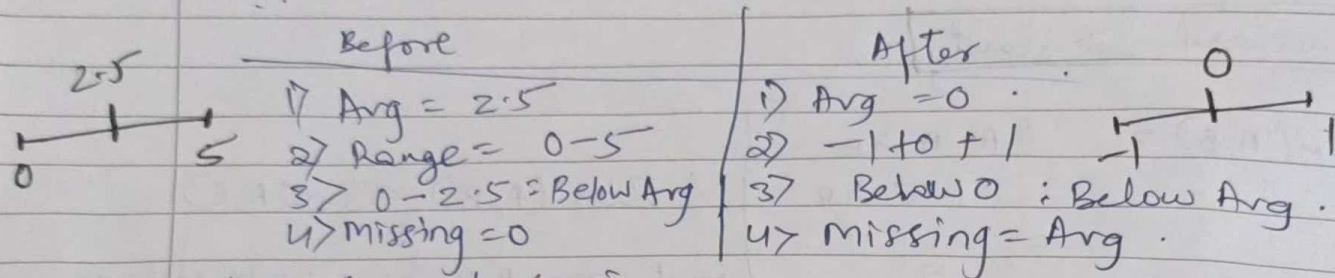
classmate

Go use Euclidean dist as well

②

	HP1	HP2	HP3	TW	SW1	SW2	DATE	Sum
A	2/3			5/3	-7/3			$\rightarrow \Sigma = 0$
B	1/3	1/3	-2/3					$\Sigma = 0$
C				-5/3	1/3	4/3		
D		0					0	

this makes Avg. Rating \forall users = 0



Now Compute Cosines

$$\text{Sim}(A, B) = \frac{2}{3} * \frac{1}{3}$$

$$\sqrt{\left(\frac{2}{3}\right)^2 + \left(\frac{5}{3}\right)^2 + \left(\frac{7}{3}\right)^2} \cdot \sqrt{\left(\frac{1}{3}\right)^2 + \left(\frac{1}{3}\right)^2 + \left(\frac{-2}{3}\right)^2}$$

$$= \frac{2/9}{140/729} \approx 0.092$$

$$\text{Sim}(A, C) = -0.56$$

$$\therefore \text{Sim}(A, B) > \text{Sim}(A, C)$$

ITEM ITEM Collaborative filtering

(More useful, outperforms user-user CF)

\therefore items are simpler than users

\rightarrow items : fixed set of genres.

\rightarrow users : kind of like fuzzy logic Nature.

Because you watched Crime patrol, you might also like Sardhaan India.

precision @ k

Recall @ k



#K Recommended

#Actual Relevant

classmate

Estimate Rating done by Users for movie

DATE

	Users											
	1	2	3	4	5	6	7	8	9	10	11	12
1	1											
2			3									
3			5	4	?	5			5		4	
4	2	4					4					
5		2	4	1	2		3		4	2	1	3
6	1		4	3	5			4		3	5	
			3		4	2					2	
					3			2			4	5

	1	2	3	4	5	6	7	8	9	10	11	12
1	-2.6	0	-0.6	0	?	1.4	0	0	1.4	0	0.4	0
2	0	0	1.84	0.84	?	0.84	0	0	0	-1.16	-2.16	-0.16
3	-1	1	0	-2	-1	0	0	0	1	0	2	0
4	0	-1.4	0.6	0	1.6	0	0	0.6	0	0	-1.4	0
5	0	0	0.7	-0.3	0.7	-1.3	0	0	0	0	-1.3	1.7
6	-1.6	0	0.4	0	0.4	0	0	-0.6	0	0	1.4	0

$\text{Sim}(M_1, \text{other Movies Rated by User 5}) = ?$ 2 movies highest similar.

$$\text{Sim}(M_1, M_3) = \frac{2.6 \times 1 + 1.4 \times 1 + 0.4 \times 3}{\sqrt{1+1+9}} = 0.41$$

$$\text{Sim}(M_1, M_4) = -0.10$$

$$\text{Sim}(M_1, M_5) = -0.31$$

$$\text{Sim}(M_1, M_6) = -0.59$$

$$\text{Sim}(M_1, M_1) = 1$$

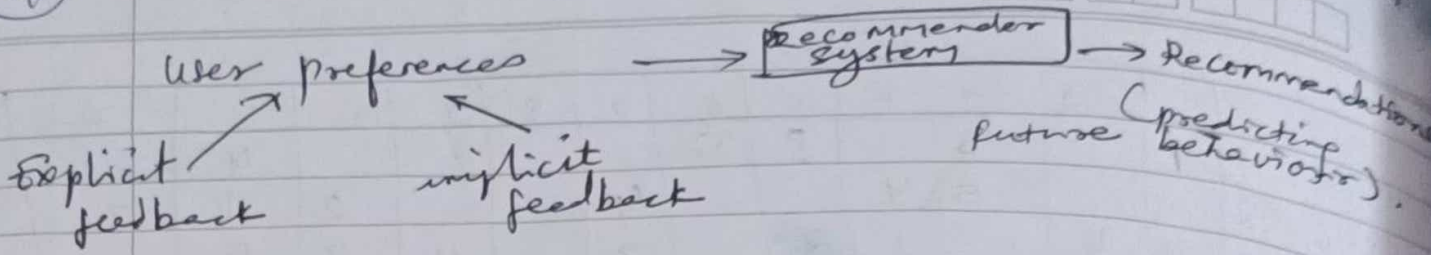
$$\text{Sim}(M_1, M_2) = -0.18$$

$$\therefore \text{Ratings for } [?] \Rightarrow \frac{0.41 \times 2 + 0.59 \times 3}{0.41 + 0.59} = 2.6$$

$$\therefore [?] = 2.6$$

(4)

DATE



Collaborative:

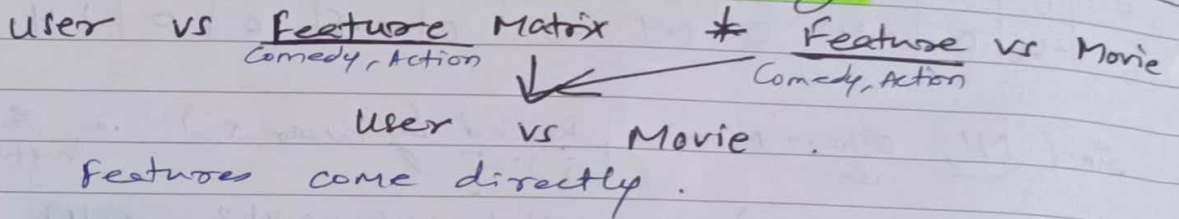
Similar users like similar things.
User x item matrix

Content Based:

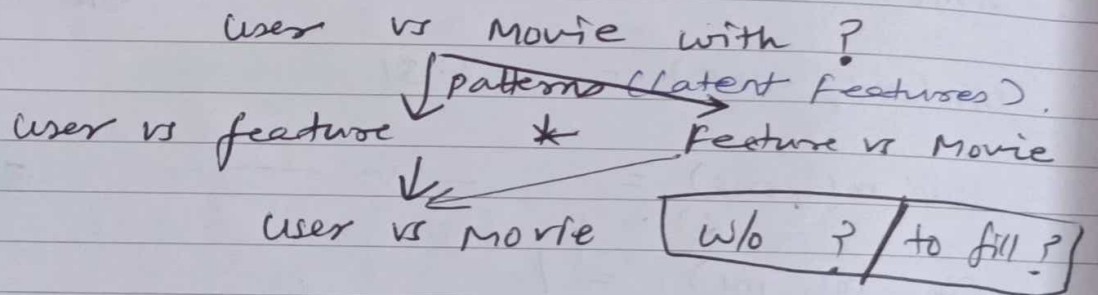
Considers item/item features
item x item
movie genre,
Year of Release,
Cast,
director,
prod. house.

user x user features
age, gender,
Spoken language

1. Content Based Filtering



2. Collaborative Based Filtering



you keep guessing values for both matrices until you get closest Matrix Multipl. product same as from where u started.

Applications of Recommend. systems: News/Songs/...
Eg. Synthetic Control. What'll be effect of "Gun Control" policies if Implemented? you check for countries you that already Implemented em

- experience .
- You, as a reader or student of a particular course, learn something and then become equipped or capable of carrying out various tasks based on what you learnt. You might have heard about **Bloom's Taxonomy**. Bloom's Taxonomy is a classification of the different objectives and skills that learners could achieve out of a particular learning or a course. Fig. 1.1.1 outlines Bloom's Taxonomy.

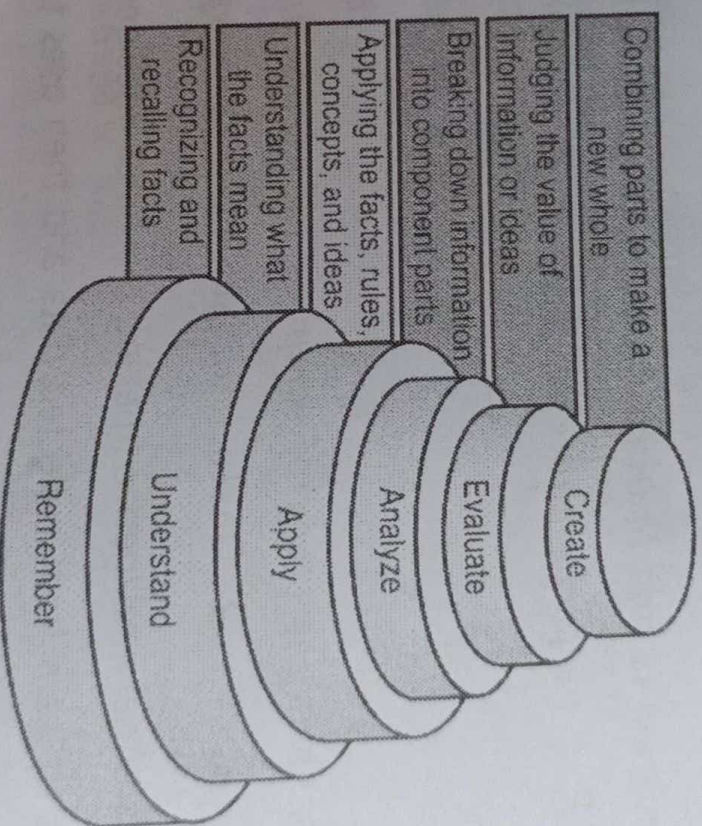


Fig. 1.1.1 : Bloom's Taxonomy

Table 1.1.2

Comparison Attributes	Structured Data	Semi-structured Data	Unstructured Data
Volume of Data	Low	Medium	High
Processing Complexity	Low	Medium	High
Data generated by	Humans and Machines	Machines	Humans
Data usually stored in	Relational Databases	Textual files	Binary files
Patterns and Schema	Fixed	Flexible	Random
Specialised Tools	Not required	Not required	Required

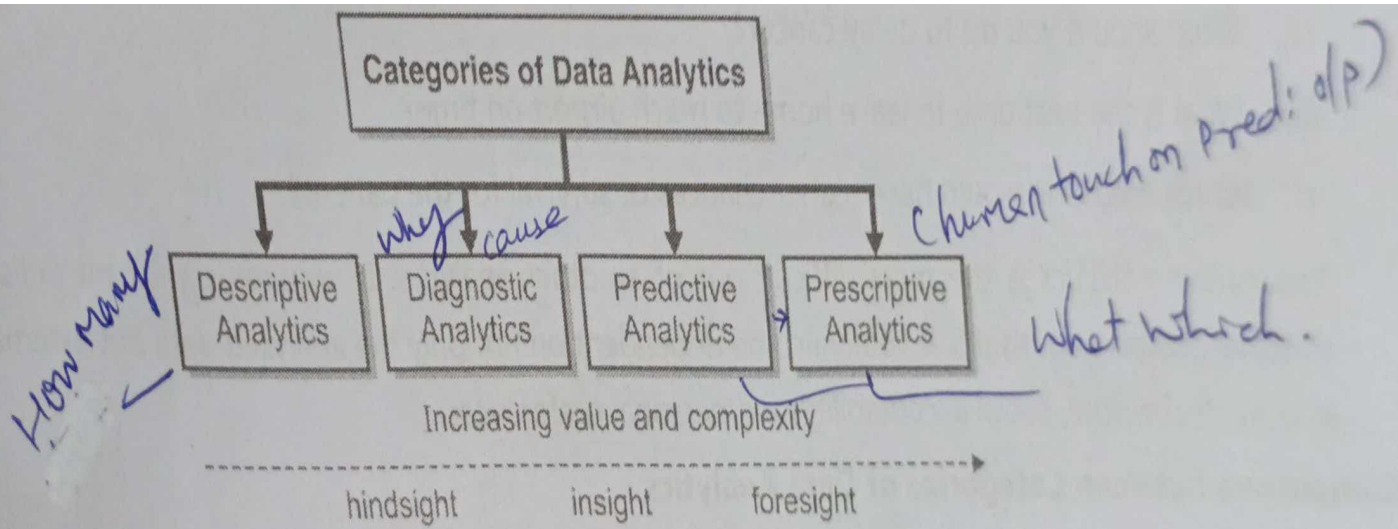


Fig. 1.1.5 : Categories of Data Analytics

Table 1.1.3

Comparison Attribute	Descriptive Analytics	Diagnostic Analytics	Predictive Analytics	Prescriptive Analytics
Complexity	Least	Medium	High	Highest
Time requirement to produce results	Low	Medium	High	Very High
Value of results	Short Term	Medium Term	Long Term	Very long term
Data enrichment level	Data	Information	Knowledge	Wisdom
Analytics Frequency	Most Common	Frequent	Not often	Rare

Table 1.2.1

Comparison Attribute	Supervised Learning	Unsupervised Learning
Training Dataset Contains	Both input and output	Only input
Used for	Classification and Prediction	Finding patterns and understanding data
Training Data	Is Labelled	Not labelled
Number of targets	Known beforehand	Not known
Feedback from user	Provided	Not provided
Complexity	High	Low

Table 1.2.2 *subset of AI*

Comparison Attribute	Machine Learning	Artificial Intelligence
Focus	Learn from data	Solve complex problems
Complexity	Low	High
Scope	Narrow	Broad
Human interaction	Minimal	High

Comparison Attribute	Machine Learning	Data Mining
Building a trained model	Required	Not required
Human effort required	Only for building model	For extracting information from data
Use of specific algorithms	Frequent	Rare
Accuracy	High	Low
Tasks carried out by	Machines	Humans
Self-learning	Yes	No

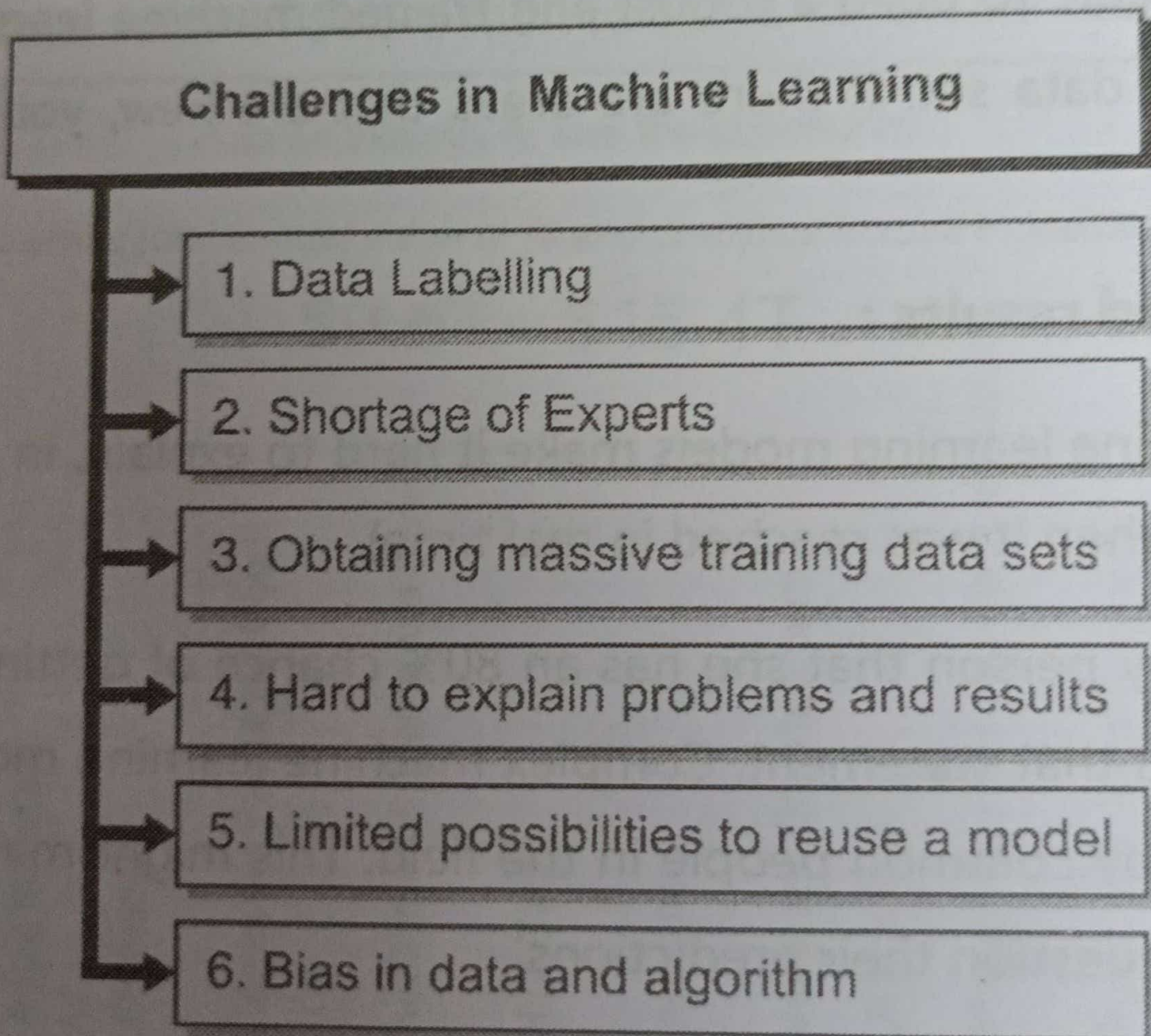
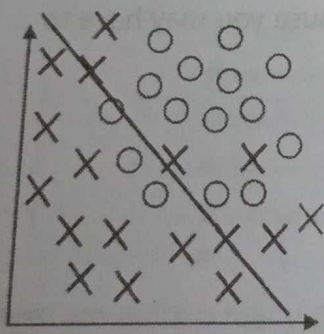
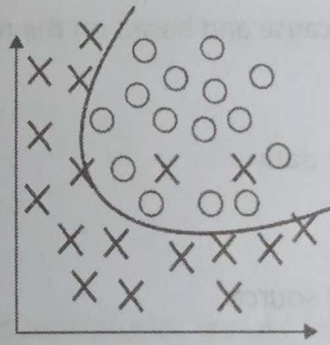


Fig. 1.3.1 : Challenges in Machine Learning

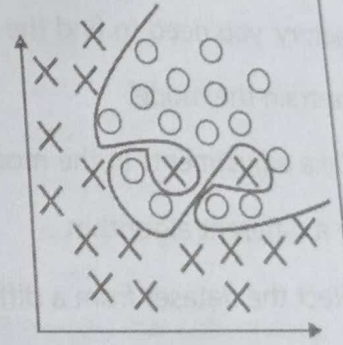
everything u read != what u learn
not everything is correct



Under-fitting
(too simple to
explain the variance)



Appropriate-fitting

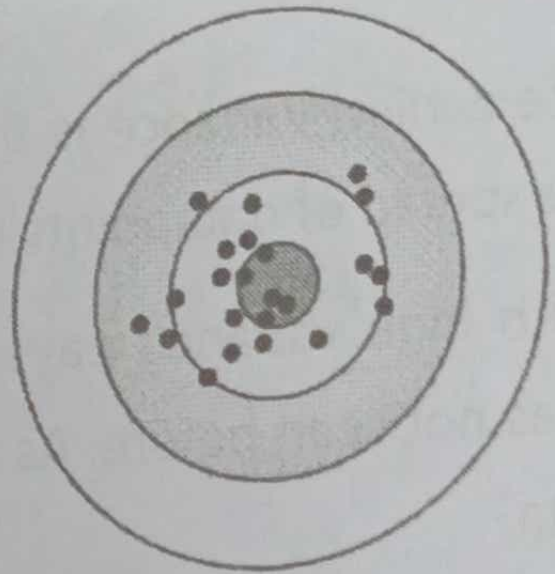
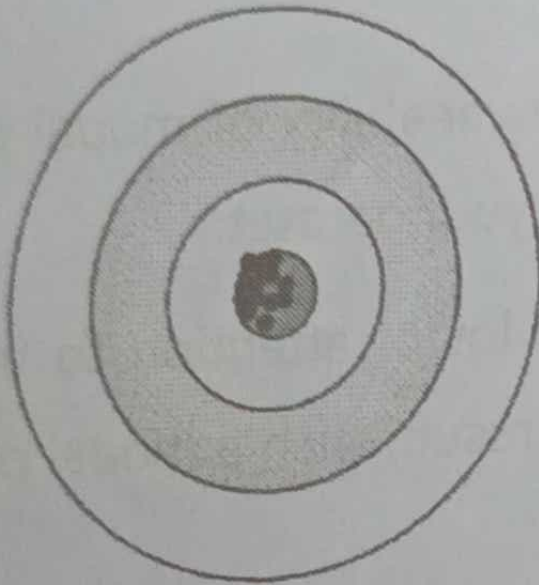


Over-fitting
(forcefitting-too
good to be true)

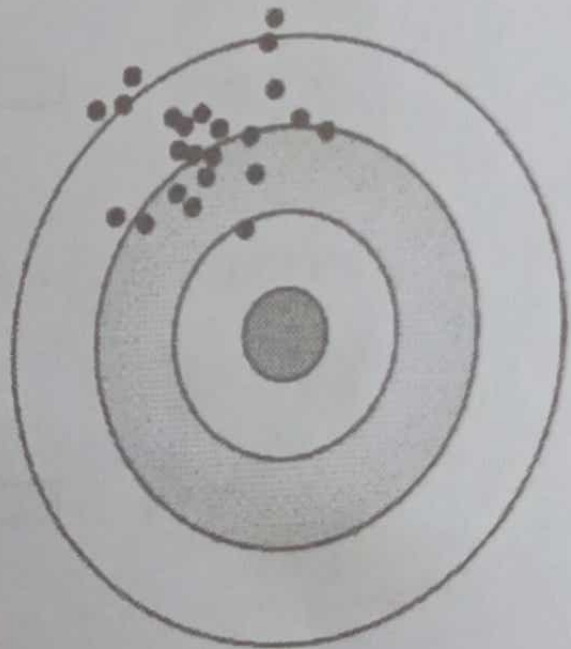
Low Variance

High Variance

Low Bias



High Bias



Algorithm Name	Bias	Variance
Linear Regression	High	Low
Decision Tree	Low	High
Bagging	Low	High
Random Forest	Low	High

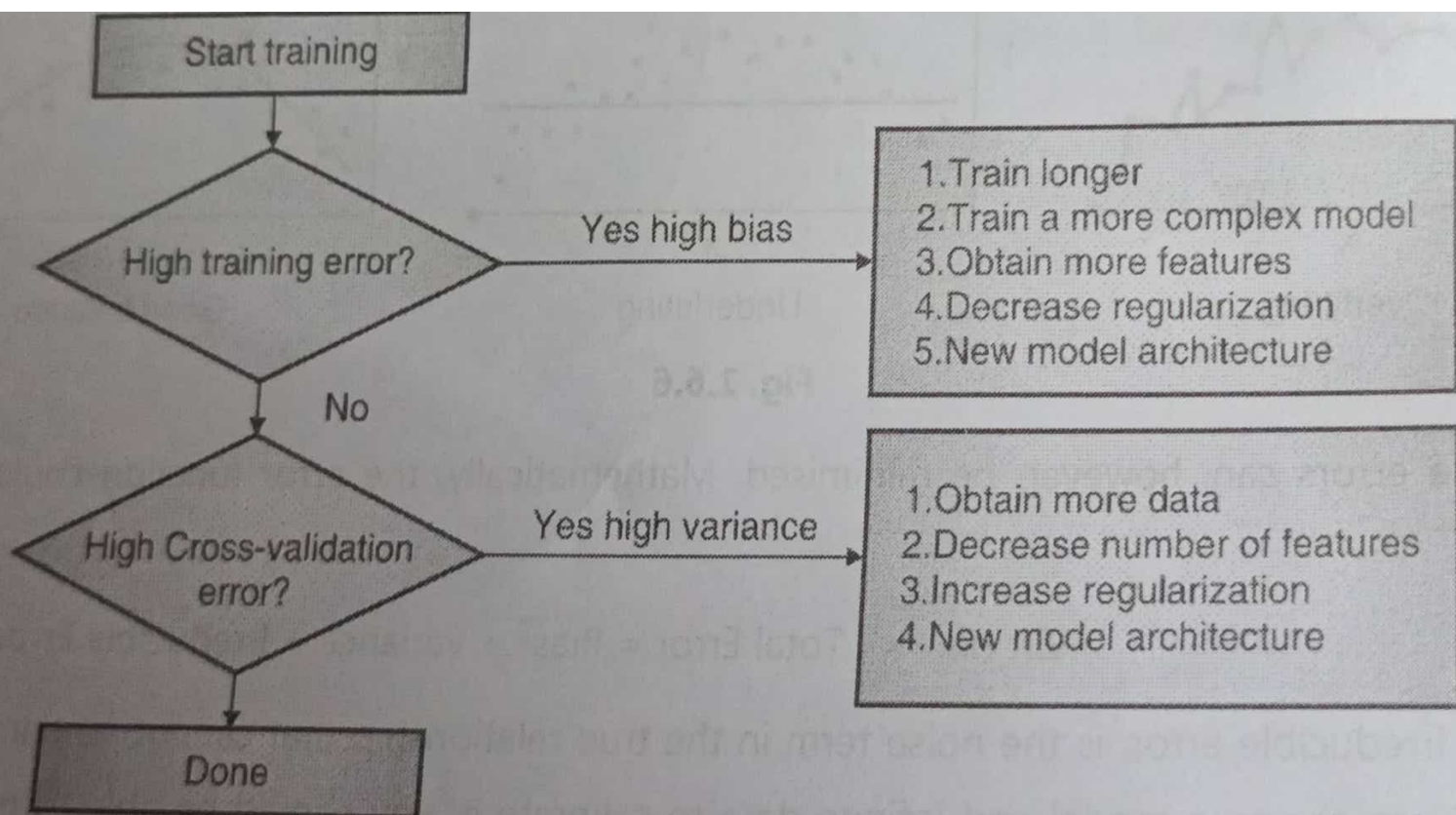


Fig. 1.6.8

1.6.5 Characteristics (Detection) of a High Bias Model

Characteristics of a high bias model are as following.

1. Failure to capture proper data trends
 2. ✓ Low training accuracy
 3. Potential towards underfitting
 4. More generalised or overly simplified
 5. ✓ High error rate
-

(Copyright No. L-98904/2021)

1.6.6 Characteristics (Detection) of a High Variance Model

Characteristics of a high variance model are as following.

1. Noise in the data set
2. Low testing accuracy
3. Potential towards overfitting
4. Complex models
5. Trying to put all data points as close as possible

<u>Simple linear regression</u>	<u>multiple linear regression</u>
1 independent variable (X) 1 dependent variable (Y)	n independent variable (X ₁ , X ₂ ,...) 1 dependent variable (Y)
$Y = mX + c$	$y = b_0 + b_1 \cdot x_1 + b_2 \cdot x_2 + \dots + b_n \cdot x_n$
only one relationship (X,Y)	more than 1 relationship. between (X ₁ ,Y) ; (X ₂ ,Y)... between dependent & independent & X ₁ ,X ₂ ; X ₁ ,X ₃ between 2 independent.
<ul style="list-style-type: none"> ➤ Y increases by factor of m if X increases by 1 ➤ Y = c or intercept if X=0 	<ul style="list-style-type: none"> ➤ Y increases by factor of b₁ if X₁ increases by 1 keeping X₂ constant ➤ Y increases by factor of b₂ if X₂ increases by 1 keeping X₁ constant ➤ Y = a or b₀ if all independent variables are 0

- Adding more independent variables to a multiple regression procedure does not mean the regression will be “better” or offer better predictions; in fact it can make things worse. This is called OVERFITTING.
- The addition of more independent variables creates more relationships among them. So not only are the independent variables potentially related to the dependent variable, they are also potentially *related to each other*. When this happens, it is called MULTICOLLINEARITY.
- The ideal is for all of the independent variables to be correlated with the dependent variable but NOT with each other.

Multiple Regression Model

$$y = \underbrace{\beta_0 + \beta_1 x_1 + \beta_2 x_2 + \cdots \beta_p x_p}_{\text{linear parameters}} + \underbrace{\epsilon}_{\text{error}}$$

Multiple Regression Equation

$$E(y) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \cdots \beta_p x_p$$

error term assumed to be zero

SVM		
Comparison Attribute	Logistic Regression	
Good for	Linear classification	Both linear and non-linear classification
Decision boundary	Multiple	One (best one)
Approach	Statistical	Geometrical
Errors	Comparatively higher	Comparatively Lower



- As in case of SVM, consider the two decision boundaries and a hyperplane. Your objective is to consider the points that are within the decision boundary line. The best fit line (or regression line) is the hyperplane that has maximum number of points. Assume that the decision boundaries are at any distance, say 'a', from the hyperplane. So, these are the lines that you draw at distance '+a' and '-a' from the hyperplane. Based on SVM, the equation of the hyperplane is as following.

$$Y = wx - b,$$

- The equations of decision boundaries are as following.

$$wx - b = a$$

$$wx - b = -a$$

- Thus, any hyperplane that satisfies SVR should satisfy $-a \leq wx - b \leq a$
- Your goal is to decide a decision boundary at 'a' distance from the original hyperplane such that data closest to the hyperplane or the support vectors are within that boundary line.

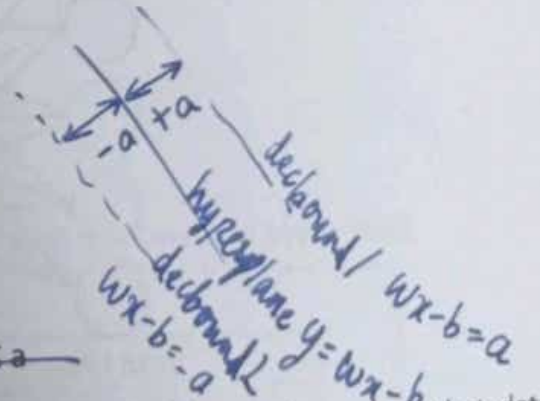


Table 4.6.1

Comparison Attribute	One vs One (OvO)	One vs Rest (OvR)
Speed	Slower than OvR	Faster than OvO
Computation Complexity	High	Low
Suitable for	Algorithms that don't scale	Algorithms that scale
No. of binary datasets or models for C classes	$\frac{C \times (C - 1)}{2}$	C
Interpretability	Low	High
Used	Less commonly	More Commonly

Table 3.1.1

Comparison Attribute	Bagging	Boosting
Weak models work	Parallely	Sequentially
All models have	Equal weight	Adjusted weights
Reduces	Variance	Bias
Overfitting	Addressed	Not addressed
Merge predictions of	Same type	Different types

C4.5 GAIN RATIO= GAIN/ SPLIT INFO

CART .. GINI INDEX

CHANGE IN PTS WHEN DATA PTS CHANGE

LINK

- The distance matrix is, $\text{AVG}[\text{dist}(P3, P6), P1]$
 $\text{dist}((P3, P6), P1) = \frac{1}{2} (\text{dist}(P3, P1) + \text{dist}(P6, P1))$
 $= \frac{1}{2} (0.22 + 0.23)$
 $= \frac{1}{2} (0.45)$
 $= 0.23$

www.anuradhabhatia.com

Technical Agglomerative Clustering [Average Link]

6 yr ago Data Warehouse and Mining ...more

Anuradha Bhatia 10.5K

Subscribe

Like



Share



Remix



Comments 51

Add a comment...



Single Linkage	This is the distance between members of the two clusters.
Complete Linkage	This is the distance between that are farthest apart.
Average Linkage	This method involves looking at all distances between all pairs of members and averaging all of these distances. This is also known as UPGMA - Unweighted Pair Grouping Method.

www.anuradhabhatia.com

Technical Agglomerative Clustering [Average Link]

6 yr ago Data Warehouse and Mining ...more

Anuradha Bhatia 10.5K

Subscribe

10K



Share



Remix



Comments 51

Add a comment...



