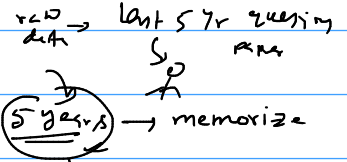
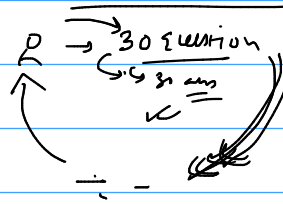


CHP 7

- Bias
- Variance
- Underfitting
- Overfitting
- Bias-variance tradeoff.

Last Year Question p.p. 4



Overfitting:

Raw data

useful information + noise

- Model is not generalized.

- training set accuracy: $\approx 100\%$

identification test set accuracy \ll training accuracy

variance: for an unknown i/p \approx Am.

reason: - len number of sample data

error \rightarrow variance

10 people
Exam $\rightarrow \approx 100\%$
1 Q₁ \approx Ans₁
2 Q₁ \approx Ans₂
3 Q₁ \approx Ans₃

model
Training data
overfit
predict (training data set)
high
But unknown data set
accuracy low

underfitting: not learn too much.

- training set accuracy is low
- test set accuracy is also low.

- error introduced because it underfit is called bias.

- reason:

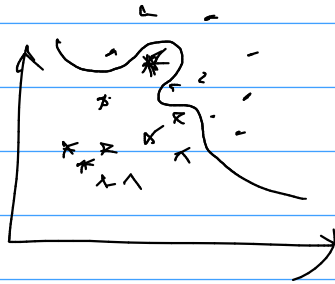
- few number of parameters.
- low training iteration

underfit \uparrow overfit
Bias-variance tradeoff

H/W Q₁ write 4 diff of bias vs variance

Q₂ write 4 diff of Underfit and overfitting of a model.

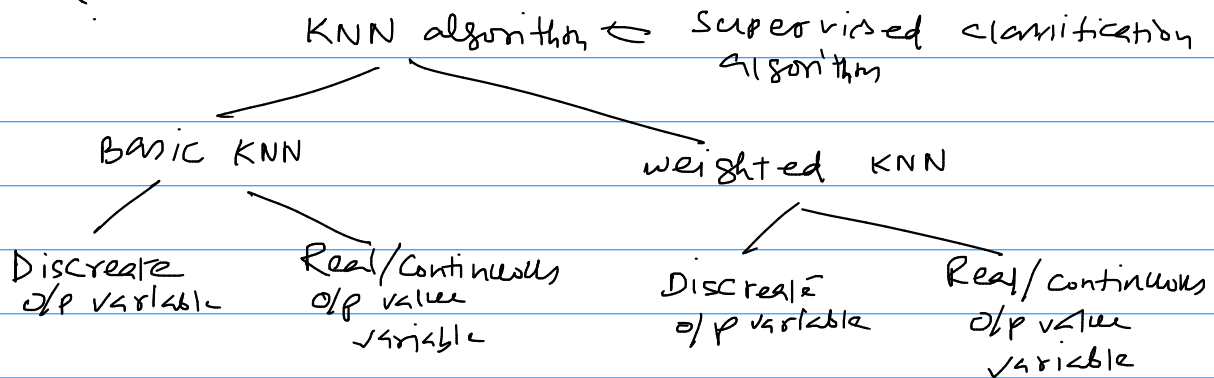
Q₃: What do you mean by bias-variance tradeoff



Ch 10 KNN algorithm

↓
K-nearest neighbour algorithm

↓
?



BASIC KNN (Discrete)

eg:

	Height	Weight	target	distance	nearest neighbour
K=3 given	150	50	M	8.06	
	155	55	M	2.77	1
	160	60	L	6.71	3
	161	59	L	6.40	2
	158	65	L	11.05	
unknown	157	54	?		

among K=3 nearest neighbour
M → 1

L = 2

note:

Distance: Euclidean distance

So
157 54 L

$$E(x, y) = \sqrt{\sum_{i=1}^n (x_i - y_i)^2}$$

Basic KNN (Continuous/Real data)

Euclidean distance

Solution Part

Q8

K=3

height	weight	target	Distance	nearest-neighbour
150	50	1.5	8.06	
155	55	1.2	2.24	1
160	60	1.8	6.71	3
161	59	2.1	6.40	2
158	65	1.7	11.05	
157	54	?		

$$\text{So, O/p for } (157, 54) = \frac{1.2 + 1.8 + 2.1}{3} = 1.7$$

- How to decide value of K

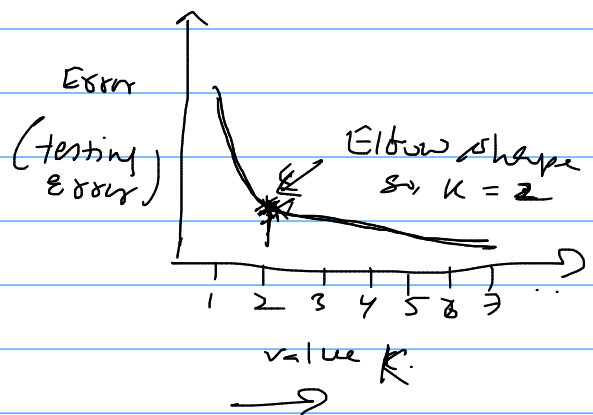
Ans: Approach 1 $K = \sqrt{n}$
 \nwarrow Sample size

Approach 2 K \in odd number (to avoid any tie)

Approach 3: Elbow method.

- Consider odd K
near 15 Elbow shape

	X	Y	Target
Training data	x_1	y_1	M
	x_2	y_2	L
	x_3	\vdots	M
	\vdots	\vdots	\vdots
	x_{50}	y_{50}	M
Testing data	x_{51}	y_{51}	L
	x_{52}	y_{52}	L
	\vdots	\vdots	M
	x_{60}	y_{60}	L



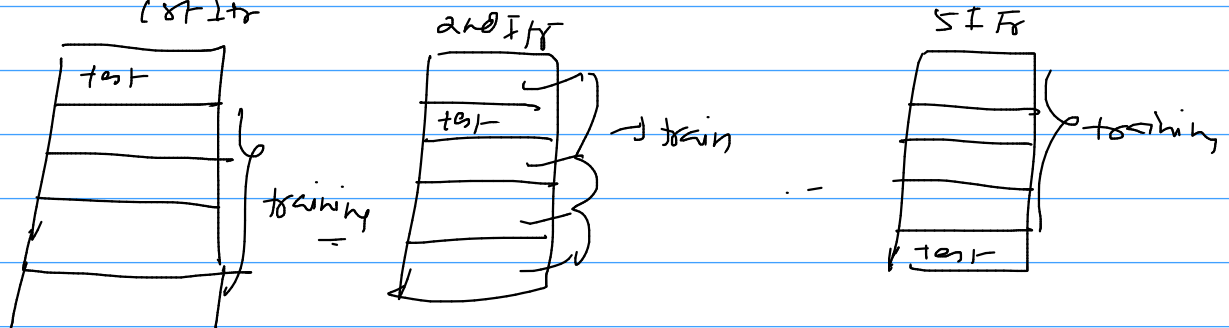
K=2

① - find o/p label for all testing data points

② - find test accuracy. (By comparing given o/p label for test data and model predicted o/p for test data)

$$\text{Error} = (100 - \text{Accuracy}) \%$$

Approach: use 5-fold cross validation



$$\text{test accuracy} = \frac{\sum_{i=1}^n (\text{test accuracy Iter}_i)}{5}$$

$$\text{Error} = (100 - \text{test accuracy}) \%$$

- plot the elbow curve

- Basic KNN algorithm for discrete valued target function is given below -

Training algorithm:

- For each training example $\langle x, f(x) \rangle$, add the example to the list training example.

Classification algorithm:

- Given a query instance x_q to be classified.

- Let x_1, x_2, \dots, x_k denotes the k instances from training example that are nearest to x_q .

- Return

$$\hat{f}(x_q) \leftarrow \underset{v \in V}{\text{argmax}} \sum_{i=1}^k \delta(v, f(x_i))$$

where $\delta(a, b) = 1$ if $a = b$ otherwise $\delta(a, b) = 0$
 V : set of output/target labels

DATA set

Solve

height weight target- dis neighbor

height	weight	target	dis	neighbor
150	50	M	8.06	
155	55	M	2.24	1
160	60	L	6.71	3
161	59	L	6.40	2
158	65	L	11.05	

157 54 ? (L)

input \rightarrow o/r

$\langle x, f(x) \rangle \leftarrow$ training set

$x: (150, 50)$

$f(x) = M$

$x_q: \langle 157, 54 \rangle$

for $k = 3$

$x_1: (155, 55)$

$x_2: (161, 59)$

$x_3: (160, 60)$

Predicted $\hat{f}(x_q)$: o/r of x_q

V : set of target o/r

so, $V \Rightarrow \{M, L\}$

BASIC KNN Algorithm

Supervised Algorithm

Classification Algorithm

two line

output is discrete label

output is continuous label

$$d(x_i, x_j) = \sqrt{\sum_{k=1}^n (x_{ik} - x_{jk})^2}$$

height	weight	target	distance	nearest
150	50	Medium	8.06	
155	55	Medium	2.24	1
160	60	Large	6.71	3
161	59	Large	6.40	2
158	65	Large	11.05	

157 54 \rightarrow Large

k : number of nearest neighbors considered

output level (V) = {Medium, Large}

$$\hat{f}(x_q) = \underset{v \in V}{\text{argmax}} \sum_{i=1}^k \delta(v, f(x_i))$$

$f(\text{Medium}) = f(\text{Medium, Medium}) + f(\text{Medium, Large})$

$f(\text{Large}) = f(\text{Large, Medium}) + f(\text{Large, Large})$

argmax

$$\delta(M, M) + \delta(M, L) + \delta(M, L)$$

$$= 1$$

$$\delta(L, M) + \delta(L, L) + \delta(L, L)$$

$$= 2$$

$$= L$$

BASIC KNN when output is continuous/real value:

Training Algorithm

- For each training example $\langle x, f(x) \rangle$, add the example to the list training-examples

Classification Algorithm

- Given a query instance x_q to be classified
- let x_1, x_2, \dots, x_k denote the k -nearest instance of x_q .
- Return $f(x_q) \leftarrow \frac{\sum_{i=1}^k f(x_i)}{k}$

BASIC KNN when output level is continuous.

Height	Weight	target	Distance	Nearest	Consider, $k=3$
150	50	1.5	8.06		
155	55	1.2	2.39	1	
160	60	1.8	6.71	3	
161	59	2.1	6.40	2	
158	65	1.7	11.05		
157	54	1.7			

$$f(157, 54) = \frac{1.2 + 2.1 + 1.8}{3} = \frac{5.1}{3} = 1.7$$

$$f(\langle 157, 54 \rangle) =$$

Discrete o/r

Distance Weighted KNN Algorithm

- Supervised algorithm
- Classification algorithm
- output type
 - Discrete
 - real valued

for discrete output value

$$f(x_q) \leftarrow \underset{v \in V}{\operatorname{argmax}} \sum_{i=1}^K w_i S(v, f(x_i))$$

where, $w_i = \frac{1}{d(x_q, x_i)^2}$

Height	Weight	target	distance	Nearest	$w_i = \frac{1}{(\text{distance})^2}$
150	50	M	8.06		
155	55	M	2.29	1	$\frac{1}{(2.29)^2} \approx 0.45$
160	60	L	6.71	3	0.15
161	59	L	6.40	2	0.16
158	65	L	11.05		
157	54	M			

$$f((157, 54), M) = 0.45 f(M, M) + 0.15 f(M, L) + 0.16 f(M, L)$$

$$= 0.45 \times 1$$

$$f((157, 54), L) = 0.45 f(L, M) + 0.15 f(L, L) + 0.16 f(L, L)$$

$$= 0.45 \times 0 + 0.15 \times 1 + 0.16 \times 1$$

$$= 0.31$$

Distance Weighted KNN for real valued output data

$$f(x_q) \leftarrow \frac{\sum_{i=1}^K w_i f(x_i)}{\sum_{i=1}^K w_i}$$

where, $w_i = \frac{1}{d(x_q, x_i)^2}$

→ continuous

BASIC KNN when output level is continuous.

Height	Weight	target	Distance	Nearest	Consider, K=3
150	50	1.5	8.06		
155	55	1.2	2.29	1	
160	60	1.8	6.71	3	
161	59	2.1	6.40	2	
158	65	1.7	11.05		
157	54	M			

$$f((157, 54)) = \frac{1.2 + 1.8 + 2.1}{3} = \frac{5.1}{3} = 1.7$$

$w_i = \frac{1}{(\text{distance})^2}$

$$\frac{1}{(2.29)^2} \approx 0.45$$

$$\frac{1}{(6.71)^2} \approx 0.15$$

$$\frac{1}{(6.40)^2} \approx 0.16$$

$$f(157, 54) = 0.45 \times 1.2 + 0.15 \times 1.8 + 0.16 \times 2.1$$

$$= 0.45 + 0.15 + 0.16$$

$$= 1.507$$

$$\sqrt{(157 - 150)^2 + (54 - 50)^2}$$

✓ NOTE :

(x)
- KNN algorithm is based on distance. So, if the i/p are in different scale then it might produce incorrect result-

↓ Soⁿ is make all inputs are in same scale

- Z-normalization \Leftarrow
or
- min-max standardization \Leftarrow (in exam for standardization you can prefer either of this two method).

↓ then apply KNN algorithm.