



University of Asia Pacific

Department of CSE

Mid-Semester Examination, Fall 2020

Name: Rashik Rahman

Reg ID: 17201012

Year: 4th

Semester: 1st

Course Code: CSE 427

Course Title: Machine Learning

Date: 23.02.2021

"During Examination and upload time I will not take any help from anyone. I will give my exam all by myself."

University of Asia Pacific

Admit Card

Mid-Term Examination of Fall, 2020

Financial Clearance

PAID

Registration No : 17201012

Student Name : Rashik Rahman

Program : Bachelor of Science in Computer Science and Engineering



Sl.NO.	COURSE CODE	COURSE TITLE	CR.HR.	EXAM. SCHEDULE
1	CSE 400	Project / Thesis	3.00	
2	CSE 330	Industrial Training	1.50	
3	CSE 401	Mathematics for computer Science	3.00	
4	CSE 403	Artificial Intelligence and Expert Systems	3.00	
5	CSE 404	Artificial Intelligence and Expert Systems Lab	1.50	
6	CSE 405	Operating Systems	3.00	
7	CSE 406	Operating Systems Lab	1.50	
8	CSE 407	ICTLaw, Policy and Ethics	2.00	
9	CSE 410	Software Development	1.50	
10	CSE 427	Topics of Current Interest	3.00	

Total Credit: 23.00

1. Examinees are not allowed to enter the examination hall after 30 minutes of commencement of examination for mid semester examinations and 60 minutes for semester final examinations.

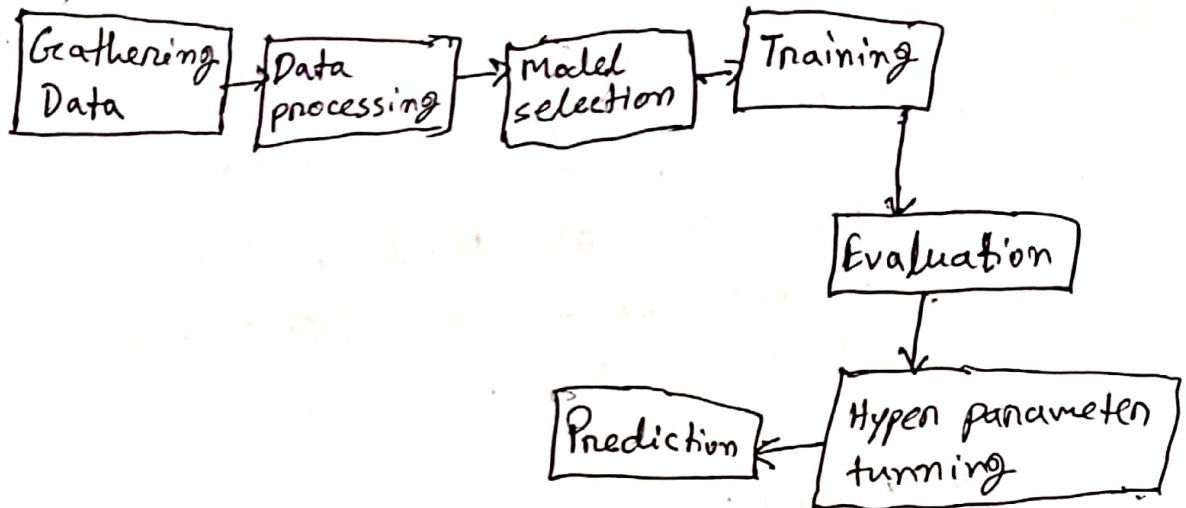
2. No examinees shall be allowed to submit their answer scripts before 50% of the allocated time of examination has elapsed.

3. No examinees would be allowed to go to washroom within the first 60 minutes of final examinations.

4. No student will be allowed to carry any books, bags, extra paper or cellular phone or objectionable items/incriminating paper in the examination hall.
Violators will be subjects to disciplinary action.

This is a system generated Admit Card. No signature is required.

Admit Card Generation Time: 21-Feb-2021 10:23 PM

Answer to the Q.No.1 (a)

- i) Gathering Data: To train a model we need data and this data can be of any kind and needs to be collected. Mostly data are gathered by the means of web scrapping.
- ii) Data processing: In this part it is said that 70% of a machine learning task's time is consumed. Here we need to clean and structure the unstructure data. We structure the data in such a manner that is a good input for model. We also do EDA here.
- iii) Model selection: There are various models out there. No one knows which model performs well for a certain dataset. So we have to select a model that performs well with the given data. We can do this with the help of ~~the~~ fold cross validation, thus pick the best model among all the models.

iv) Training: After selecting the model we train it with provided data.

v) Evaluation: We use different means like accuracy, logloss, RSME, $f1$ score, confusion matrix ~~to etc~~, ROC & AUC curve to evaluate the model.

vi) Hyperparameter tuning: We can use grid search CV to tune the parameters of the selected model thus achieve higher accuracy as we evaluate it again.

vii) Prediction: Finally when the model is finished training & tuning then it is ready to predict for a fresh data.

Answer to the Q.No.1(b)

Machine learning problem are typically two types those are supervised and unsupervised problems. Supervised problem can be further specified in two types these are regression and classification problem. In classification problem the dependent attribute or we can say ~~labels has values~~ predicted values are discrete and we call them labels. Ex:

Score1	Score2	result
29	43	Pass
22	29	Fail
10	47	Fail
31	55	Pass
17	18	Fail
33	54	Pass
32	40	Pass
20	41	Pass

Here we can see that the dependent attribute result has discrete values of just pass and fail. So in a classification problem a model learns from independent ~~data~~ attributes the ~~pro~~ predicts a discrete value.

Answer to the Q. No. 2(a)

The basic ~~ba~~ of naive bayes ~~algorithm~~ algorithm is:

$$P(c|x) = (P(x_1|c) \cdot P(x_2|c) \cdot P(x_3|c) \dots P(x_n|c)) \times \frac{P(c)}{P(x)}$$

①

Calculation of probabilities are:

if a condition B is given and we need to find probability of a condition A happening then we can say,

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

If we breakdown eqn ① then we can understand the followings:

$P(c|x)$ = Posterior probability

$P(x|c)$ = Likelihood

$P(c)$ = ^{class} prior probability

$P(x)$ = Prediction prior probability.

Answer to the Q. No. 2(b)

freq table

chills			running nose			headache			fever		
	Y	N		Y	N		Y	N		Y	N
Y	3	1	Y	4	1	no	1	1	Y	4	1
						mild	2	1			
N	2	2	N	1	2	strong	2	1	N	1	2
							5	3		5	3
	5	3		5	3						

likelihood table

chills			running nose			headache			fever		
	Y	N		Y	N		Y	N		Y	N
Y	3/5	1/3	Y	4/5	1/3	no	1/5	1/3	Y	4/5	1/3
						mild	2/5	1/3			
N	2/5	2/3	N	1/5	2/3	strong	2/5	1/3	N	1/5	2/3

Now for

chills = Y

fever = Y

headache = mild

nose running = N

$$P(\text{Yes}) = \frac{3}{5} \times \frac{4}{5} \times \frac{2}{5} \times \frac{1}{5} = 0.0384$$

$$P(\text{No}) = \frac{1}{3} \times \frac{1}{3} \times \frac{1}{3} \times \frac{2}{3} = 0.025$$

Overall probability:

$$P(\text{Yes}) = 0.0384 \times \frac{5}{8} = 0.024$$

$$P(\text{No}) = 0.025 \times \frac{3}{8} = 0.0094$$

So the answer will be a patient has high probability to have flu.

~~Answer to the Q. No.~~

Answer to the Q. No. 4.

~~Intro~~

Here we use; $-P_1 \log_2(P_1) - P_2 \log_2(P_2) - P_3 \log_2(P_3)$
to calculate gain;

on again of whole dataset,

$$\text{Gain}(\text{dataset}) = -\frac{4}{10} \log_2\left(\frac{4}{10}\right) - \frac{3}{10} \log_2\left(\frac{3}{10}\right) - \frac{3}{10} \log_2\left(\frac{3}{10}\right)$$

$$= -0.4(-1.32) - 0.3(-1.74)$$

$$- 0.3(-1.737)$$

$$= 1.571$$

For Gender:

	Bus	Train	Car	total	$-P_1 \log_2 P_1 - P_2 \log_2 P_2 - P_3 \log_2 P_3$
Male	3	1	1	5	1.371
Female	1	2	2	5	1.522

Info gain of Gender:

$$= \text{Gain}(\text{dataset}) - \left\{ \frac{5}{10} \text{gain}(\text{Male}) + \frac{5}{10} \text{gain}(\text{Female}) \right\}$$

$$= 1.571 - \{ 0.5 \times 1.371 + 0.5 \times 1.522 \}$$

$$= 1.571 - 1.447 = 0.125$$

For canownership

	Bns	Train	Can	total	$-P_1 \log_2 P_1 - P_2 \log_2 P_2 - P_3 \log_2 P_3$
0	2	1	0	3	0.918
1	2	2	1	5	1.522
2	0	0	2	2	0

Info gain of canownership

$$= \text{Gain}(\text{dataset}) - \left\{ \frac{3}{10} \text{gain}(0) + \frac{5}{10} \text{gain}(1) + \frac{2}{10} \text{gain}(2) \right\}$$

$$= 1.571 - \{ 0.3 \times 0.918 + 0.5 \times 1.522 + 0.2 \times 0 \}$$

$$= 0.535$$

17201012
⑧

For travel cost

	Bus	train	Car	total	$-P_1 \log_2 P_1 - P_2 \log_2 P_2 - P_3 \log_2 P_3$
cheap	4	1	0	5	0.722
Standard	0	2	0	2	0
expensive	0	0	3	3	0

∴ Info gain of travel cost

$$= \text{Gain}(\text{dataset}) - \left\{ \frac{5}{10} \text{gain}(\text{cheap}) + \frac{2}{10} \text{gain}(\text{standard}) + \frac{3}{10} \text{gain}(\text{expensive}) \right\}$$

$$= 1.571 - 0.5 \times 0.722$$

$$= 1.21$$

For Income level:

	Bus	train	car	total	$-P_1 \log_2 P_1 - P_2 \log_2 P_2 - P_3 \log_2 P_3$
Low	2	0	0	2	0
Medium	2	2	1	5	1.459
High	0	0	2	2	0

∴ Info gain of Income level:

$$\Rightarrow \text{Gain}(\text{dataset}) - \left\{ \frac{2}{10} \text{gain}(\text{Low}) + \frac{6}{10} \text{gain}(\text{medium}) + \frac{2}{10} \text{gain}(\text{High}) \right\}$$

$$= 1.571 - 0.2 \times 1.459$$

$$= 1.279$$

Summary:

$$\text{Info gain}(\text{Gender}) = 0.125$$

$$\text{Info gain}(\text{Car ownership}) = 0.535$$

$$\text{Info gain}(\text{Travel cost}) = 1.29$$

$$\text{Info gain}(\text{Income level}) = 1.279$$

