Shahjalal University of Science and Technology Department of Computer Science and Engineering

3rdYear 2nd Semester Final Examination, December 2018 (Session: 2015-16)

Course Code: CSE 345

Credits: 2

Course Title: Introduction to Data Science

Total Marks: 50 Time: 2hrs

Group A

[Answer all the questions]

Answer any FIVE Define data

What is categorical variable?

Define population and sample.

- What is selection bias?
- What is an irreducible error?
- What is a scatter plot?

What is 95% confidence interval for linear regression?

Answer any FOUR

4*2.5=10

5*1=5

Write the Data Science Process.

What is Bias? How bias can happen in samples? Differentiate between regression and classification.

Assume there are 4 features (X₁, X₂, X₃, X₄) in a data set, where the interaction of X₂ and X₄ might have some impact on the response. Write the appropriate Linear Regression model for it.

What is a pie chart? Assume there are 100 items: 50 are red, 25 are green and rest are blue. Draw a pie chart for the data

Differentiate between linear regression and logistic regression.

Answer any TWO

2*5=10

Perform KNN regression for the following data and show the predictions for all the X values. Assume k = 2.

X	1	2	3	4	5
Y	1	2	5	4	5

What is F statistics? The true values and predicted values of a model are given below. Calculate the F statistics for the data and interpret the result

X	1	2	3	4	5
Y	6	7	4	3	2
Ŷ	6.8	5.6	4.4	3.2	2

What is MSE? Calculate the MSE for the above model in 4(3b)

Group B

[Answer all the questions]

Answer any FIVE 4.

Define Akaike's Information Criterion. a) What are interacting predictors?

What is overfitting?

What is cross validation?

What is Model Fitness, R2?

What is bootstrapping?

Write a polynomial regression model of degree M.

2+0704

Answer any FOUR

How regression can be done by decision trees?

What is k-fold cross validation?

What is LASSO Regression?

Define Entropy and Information Gain.

Compare between mean and median.

Differentiate between parametric and non-parametric models with example.

Answer any TWO

2*5=10

4*2.5=10

Write the steps of computing principal components of N data with J features. a) For the given data, approximate the values of B_0 and B_1 and write the equation for simple liner regression.

X 1 2 3 4 5 Y 6 7 4 3 2

What is hypothesis testing? Write the steps of testing a hypothesis. (c)

Shahjalal University of Science and Technology Department of Computer Science and Engineering

3rd Year 2nd Semester Final Examination' Dec 2019 (Session: 2016-17) Course Title: Data Science Course Code: CSE 345 Credits: 2

Total Marks: 50 Time: 2 hrs

Group A

[Answer all the questions]

5x1=5Answer any FIVE .1: What is Structured Data? a) Define Quantitative Variable with examples. b) What is a Scatter Plot? c) Define Variance of a Sample. d) What is a Statistical Model? e) What is a Loss Function? f) How to interpret a P-Value? g) If B is a Linear Regression Coefficient and the Standard Error (B) = 1.5, then find the 95% h) Confidence Interval of B. 4x2.5=10Answer any FOUR 2. Describe the process of Data Science. a) Define Data, Distribution of Data, Population, Sample and Bias. b) Name and define different Measures of Centrality; Mention the Computational Complexity of each. c) What are the principles of Data Visualization? d) What is a Histogram? Draw the Histogram of the following data: e) Height(Feet) Frequency 0-2 0 1 2-4 4 4-5

	5-6	8	
	6-8	2	a Disc.
Pre	ediction Pro	blems we disci	issed in this course? Differentiate

What are the Two main types of f) between them.

Answer any TWO

Assume there are 4 predictors (X_1, X_2, X_3, X_4) in a data set, where the interaction of X_2 and X_4 might 3. a) have some impact on the response.

Write the appropriate Linear Regression Model for it.

Write the Final Model if you get the following P-Values after performing hypotheses ii.

testing on the significance of the predictors.

Coefficients	P-Value
B_{θ}	0.00
B_I	0.09
B_2	0.00
B_3	0.03
B_4	0.00
B_5	1.00
	1

What can you tell about the signification of the predictors? iii.

What is F statistics? The true values and predicted values of a model are given below. Calculate the F statistics for the data. And interpret the result. b)

v	1	2	3	4	5
V	1	2	5	4	5
Y(Predicted)	35	3	3	5	14.5
Y (Predicted)	13.5	_			

Perform Linear Regression on the following data to find the equation of the Regression line. c)

Data	Age X	Glucose Level Y
1	43	99
2	21	65
3	25	79
4	42	75

2x5=10

Group B

[Answer all the questions]

Answer any FIVE

5x1=5

- a) What is Logistic Regression?
- b) What is Ensemble Learning?
- c) What is a Random Forest?
- d) Define Entropy and Information Gain.
- e) Write the Bayes' Theorem.
- f) What is a ROC Curve?
- g) What is bootstrapping?
- h) Write a polynomial regression model of degree M.
- 5. Answer any FOUR

4x2.5=10

- a) What is Over-fitting? What are the causes of Over-fitting?
- b) Define the steps to choose the subset of significant predictors using K-fold cross validation.
- e) Why Regularization is used? Define L₁ and L₂ Regularization.
- d) Differentiate between Parametric and Non-Parametric Models with example.
- e) What is the most commonly used Loss Function? Define it and calculate its value from the following prediction.

X	1	2	3	4	5
Y	1	2	5	4	5
Y(Predicted)	3.5	3	3	5	4.5

- f) How to deal with Missing Values in data?
- 6. Answer any TWO

2x5=10

- a) What is Principal Component Analysis? Write the steps of computing principal components of N data with J features.
- b) Suppose you have written a classifier to detect which images in your favorite social network are selfies. You have tested your classifier with some data and got the following predictions.

	Target	Prediction
1	Selfie	Not Selfie
2	Selfie	Selfie
3	Not Selfie	Not Selfie
4	Selfie	Selfie
5	Not Selfie	Selfie
6	Not Selfie	Not Selfie
7	Not Selfie	Not Selfie
8	Selfie	Selfie
9	Selfie	Not Selfie
10	Not Selfie	Not Selfie

Compute the Confusion Matrix for above and calculate Accuracy from it.

c) Briefly describe the process of building a Decision Tree.

Shahjalal University of Science and Technology

Department of Computer Science and Engineering 3^{rd} year 2^{nd} Semester Final Examination—December 2020 (Session 2017-18)

Course No.—CSE 345 Course Title—Data Science

Time—5 Hours

Total Marks#30

(Answer All the Questions)

Credit: **3.00**

Group A

- 1. Determine the following **Five** statements as True or False. If false, write the correct verdict. $5 \times 1 = 5$
 - (a) Tabular format is the most suitable representation of Data.
 - (b) A Scatter Plot displays groups of numerical data through their quartiles.
 - (c) A large variance in data indicates that the values are far from the mean.
 - (d) Data visualization helps us to analyze and explore the data.
 - (e) A statistical model is any algorithm that estimates the underlying function that represents the relationship between dependent and independent variables.
- 2. Answer the following **Two** Questions.

 $2 \times 2.5 = 5$

- (a) Define Data, Population, Sample and Bias as briefly as possible.
- (b) If \boldsymbol{B} is a Linear Regression Coefficient where $\boldsymbol{B}=5$ and the Standard Error $(\boldsymbol{B})=1.5$, then find the 95% Confidence Interval of \boldsymbol{B} .
- 3. Consider the following dataset.

Day	Weather	Temperature	Wind	Play
1	Sunny	Hot	Strong	No
2	Cloudy	Mild	Weak	Yes
3	Sunny	Mild	Weak	Yes
4	Sunny	Mild	Strong	No
5	Rainy	Cool	Weak	No
6	Cloudy	Cool	Weak	Yes
7	Cloudy	Hot	Strong	No

Now, build a **decision tree** according to that dataset which will predict/decide whether you should play or not on a day given the *weather*, *temperature* and *wind* information of that day.

Group B

- 1. Determine the following **Five** statements as True or False. If false, write the correct verdict. $5 \times 1 = 5$
 - (a) Logistic Regression is supervised learning.
 - (b) Ensemble methods use a single model to obtain better prediction.
 - (c) A Random Forest Classifier takes votes from multiple Decision Trees.
 - (d) A model that has a lower AIC (or BIC) is better than other models.
 - (e) A ROC Curve illustrates the trade-off for all possible thresholds chosen for the two types of classification error.
- 2. Answer the following \mathbf{Two} Questions.

 $2 \times 2.5 = 5$

- (a) What is Cross Validation? How to perform K-fold Cross Validation?
- (b) What are the differences between Parametric and Non-Parametric Models?
- 3. Define Entropy and Information Gain. Suppose you have tossed a 4 faced dice 1000 times where 1, 2, 3 and 4 showed up 250, 500, 125, and 125 times respectively. Calculate the Entropy for this dice.

5

5

Shahjalal University of Science and Technology

Department of Computer Science and Engineering 3^{rd} year 2^{nd} Semester Final Examination—December 2020 (Session 2017-18)

Course No.—CSE 345
Course Title—Data Science

Time—5 Hours

Total Marks#30

(Answer All the Questions)

Credit: **3.00**

Group A

- 1. Determine the following **Five** statements as True or False. If false, write the correct verdict. $5 \times 1 = 5$
 - (a) For Quantitative Variables there is no inherent order among the values.
 - (b) A Scatter Plot displays groups of numerical data through their quartiles.
 - (c) Data visualization helps us to analyze and explore the data.
 - (d) "Some samples are more likely to be selected"—this phenomenon is called volunteer bias
 - (e) KNN Regression is a parametric model.
- 2. Answer the following **Two** Questions.

 $2 \times 2.5 = 5$

5

(a) Draw a Pie Chart with the following data:

Height (Feet)	Frequency
0-2	20
2-4	100
4-5	400
5-6	480

- (b) Write down the differences between Regression and Classification.
- 3. Write down the Linear Regression model for the following data, and find the equation of the Regression line.

Data	\mathbf{X}	\mathbf{Y}
1	4	10
2	2	6
3	3	8

Group B

- 1. Determine the following **Five** statements as True or False. If false, write the correct verdict. $5 \times 1 = 5$
 - (a) Logistic Regression is supervised learning.
 - (b) Ensemble methods use a single model to obtain better prediction.
 - (c) A Random Forest Classifier takes votes from multiple Decision Trees.
 - (d) A model that has a lower AIC (or BIC) is better than other models.
 - (e) In hypothesis testing, if F > 1, then we accept the null hypothesis.
- 2. Answer the following **Two** Questions.

 $2 \times 2.5 = 5$

- (a) What is Over-fitting? Write the causes of Over-fitting.
- (b) Name different approaches to impute missing values in a variable.
- 3. Suppose you have written a classifier to separate the images of Roshogolla and Chomchom. You have tested your classifier with some images and got the following predictions. Now make the Confusion matrix for that and calculate accuracy from it.

CSE 345: Introduction to Data Science Time: 35 minutes, Marks: 20

How do you measure node impurity in Decision Tree algorithm with Gini Indexing, Entropy and Information Gain? Discuss with an example.

Data preparation is quite important in data science. What are the steps of data preparation? How do you handle missing values?

Why should we need to address overfitting? How to address overfitting in pre-pruning and post-pruning?

Class Test#02

Time: 30 minutes, Marks: 20

Cluster the following eight points (with (x, y) representing locations) into three clusters:

A1(2, 10), A2(2, 5), A3(8, 4), A4(5, 8), A5(7, 5), A6(6, 4), A7(1, 2), A8(4, 9) Initial cluster centers are: A1(2, 10), A4(5, 8) and A7(1, 2).

The distance function between two points a = (x1, y1) and b = (x2, y2) is defined as-

$$P(a, b) = |x2 - x1| + |y2 - y1|$$

Use K-Means Algorithm to find the three cluster centers after the second iteration.

4

- What is Naïve Bayes algorithm? What advantages do we have of using this algorithm.
- What is Bayesian network? Where can we apply this?
 - What are the differences between KNN and K-mean?

TT #02

Course: Machine Learning (SWE 427) (QT-A)

Marks: 20 (8+12)

Time: 30 mins

1. What are support vectors? How do support vectors help to find the optimal margin of a model?

2. Given the data in the table below, reduce the dimensions from 2 to 1 using the PCA algorithm.

X ₁	X ₂
1	2
3	4
5	6
7	9

TT#01 Course: Machine Learning (**SWE 427**) (*QT-A*)

Marks: 20 Time: **30 mins**

Does gradient descent require a convex cost function to converge? Can we use Mean Squared Error for calculating gradient descent of Logistic Regression to converge to the global optima?
 If not, why?

- 2. Explain the role of the learning rate in gradient descent. What are the potential consequences of setting it too high or too low?
- 3. Create a fictional case study where the improper use of regularization leads to significant model failures.

 What lessons can be learned from this scenario?

 05
- 4. How do filters extract features and how does pooling simplify them in a CNN? Explain in brief. 05

TT#01 Course: Machine Learning (**SWE 427**) (*QT-B*)

Marks: 20 Time: **30 mins**

- 1. How does logistic regression differ from linear regression in terms of the nature of the dependent variable and the type of problems it solves?
- 2. Explain the role of the learning rate in gradient descent. What are the potential consequences of setting it too high or too low?
- Create a fictional case study where the improper use of regularization leads to significant model failures.
 What lessons can be learned from this scenario?
- 4. Provide a detailed mathematical breakdown of the forward propagation process in a simple neural network. **05**

TT#01 Course: Machine Learning (**SWE 427**) (*QT-C*)

Marks: 20 Time: **30 mins**

- 1. How can you visualize the decision boundary for a simple linear classifier? Discuss in brief. 05
- 2. Explain the role of the learning rate in gradient descent. What are the potential consequences of setting it too high or too low?
- Create a fictional case study where the improper use of regularization leads to significant model failures.
 What lessons can be learned from this scenario?
- 4. Evaluate the effectiveness of CNNs in image classification tasks. What problem it solved in deep learning. **05**

Shahjalal University of Science and Technology Snanjaiai University and Communication Technology Institute of Information and Communication Technology

3rd Year 2nd Semester Final Examination' Dec 2019 (Session: 2016-17)
Credits: 2 Comments and Com Course Title: Data Science Course Code: SWE 335 Total Marks: 50 Time: 2 hrs

Group A [Answer all the questions]

Answer any FIVE Define datum and data. What is Residual Standard Error? Define population and sample, What is Messy data? What are the common causes of messiness? What is selection bias? State the complexities of computing mean and median of data. What is 95% confidence interval for linear regression?

4*2.5=10

Answer any FOUR

Write the Data Science Process.

How is data represented and stored? Give examples.

Assume there are 4 features (X_1, X_2, X_3, X_4) in a data set, where the interaction of X_2 and X_4 might have some impact on the response. Write the appropriate Linear Regression model for it.

Write down the imputation methods for missing data. What is cross validation? Describe k-fold cross validation,

2*5=10

Perform KNN regression for the following data and show the predictions for all the X values. Assume 4.45

k = 2.

et of

by

Use K-means to cluster the given data: {20, 3, 9, 10, 9, 3, 1, 8, 5, 3, 24, 2, 14, 7, 8, 23, 6, 12, 18} into 3 groups (use 2 iterations).

What is MSE? Calculate the MSE for the given model. c)

	K	1	2	3	4	5
1	Y	6	7	4	3	2
1	4	6.8	5.6	4.4	32	2

[Answer all the questions]

5*1=5

Answer any FIVE

write down the type of missingness.

What are interacting predictors?

What is overfitting?

What is Variance and Standard Deviation?

What is agglomerative clustering? Give example.

Write a polynomial regression model of degree M. What is Web Scrapping?

4*2.5=10

Answer any FOUR

What are the causes of over-fitting?

Show the difference between Prediction and Estimation.

What is LASSO Regression?

Use hierarchical cluster to cluster the given data {3, 7, 10, 16, 18, 20}. Show each steps.

Differentiate between parametric and non-parametric models with example.

2*5=10

For the given data, approximate the values of B_0 and B_1 and write the equation for simple liner regression.

by

cat

fall

arat

Construct the decision tree to decide what to do in the evening, find the root (use the given dataset).

	Is there a Party?	Lazy?	Activity
Deadline?	1	Yes	Party /
Urgent \	Yes	Yes	Study \.
Urgent >	No		Party
Near · v	Yes	Yes	
None 7	Yes	No	Party
None 7	No	Yes	Pub
	Yes	No	Party /
None Near	No	No	Study \
Near V	No	Yes	TV +
Near A	Yes	Yes	Party .
Urgent	No	No	Study .

What is hypothesis testing? Write the steps of testing a hypothesis. c)

Shahjalal University of Science and Technology Department of Computer Science and Engineering

3rd Year 2nd Semester Final Examination' Dec 2019 (Session: 2016-17) Course Title: Data Science

Course Code: CSE 345 Time: 2 hrs

Total Marks: 50

Group A

[Answer all the questions]

5x1=5

- Answer any FIVE 1.
- What is Structured Data? a)
- Define Quantitative Variable with examples. b)
- What is a Scatter Plot? c)
- d) Define Variance of a Sample.
- e) What is a Statistical Model?
- What is a Loss Function? f)
- If B is a Linear Regression Coefficient and the Standard Error (B) = 1.5, then find the 95% g) h)

Confidence Interval of B.

4x2.5=10

- Answer any FOUR
- Describe the process of Data Science. a)

Define Data, Distribution of Data, Population, Sample and Bias. Name and define different Measures of Centrality; Mention the Computational Complexity of each. b)

c) What are the principles of Data Visualization? d)

What is a Histogram? Draw the Histogram of the following data: e)

1151	Height(Feet)	Frequency
	0-2	0
	2-4	11: 4
	4-5	4
	5-6	8
1	6-8	2 minutes

What are the Two main types of Prediction Problems we discussed in this course? Differentiate between them.

2x5 = 10

Assume there are 4 predictors (X₁, X₂, X₃, X₄) in a data set, where the interaction of X₂ and X₄ might Answer any TWO 3. a) have some impact on the response.

Write the appropriate Linear Regression Model for it.

Write the Final Model if you get the following P-Values after performing hypotheses testing on the significance of the predictors.

	Coefficients	P-Value
	B_0	0.00
V	B_{I}	0.09
	B_2	0.00
	B_3	0:03
Ì	B_4	0.00
1	B_5	1.00

What can you tell about the signification of the predictors? iii.

What is F statistics? The true values and predicted values of a model are given below. Calculate the F b) statistics for the data. And interpret the result.

X	1	2	3	4	5
Y	1	2	.5	4	5
Y(Predicted)	3.5	3	3	5	4.5

Perform Linear Regression on the following data to find the equation of the Regression line. c)

Data	Age X	Glucose Level Y
1	43	99
2	21	65
3.	25 ·	79
4	42	75

CS CamScanner

Group B

[Answer all the questions]

5x1=5

- Answer any FIVE 4.
- What is Logistic Regression? n)
- What is Ensemble Learning? b)
- What is a Random Forest? c)
- Define Entropy and Information Gain.
- Write the Bayes' Theorem. c)
- What is a ROC Curve?
- What is bootstrapping? g)
- Write a polynomial regression model of degree M. h)

4x2.5=10

- Answer any FOUR
- What is Over-fitting? What are the causes of Over-fitting?
- Define the steps to choose the subset of significant predictors using K-fold cross validation. a)
- b) Why Regularization is used? Define L₁ and L₂ Regularization.
- Differentiate between Parametric and Non-Parametric Models with example.
- What is the most commonly used Loss Function? Define it and calculate its value from the following d) c) prediction.

Y	1	2	3	4	5
Y	1	2	5	4	5
Y(Predicted)	3.5	3	3	5	4.5

How to deal with Missing Values in data? I)

2x5 = 10

- What is Principal Component Analysis? Write the steps of computing principal components of N data a)
 - Suppose you have written a classifier to detect which images in your favorite social network are selfies. You have tested your classifier with some data and got the following predictions.

	Target	Prediction
1	Selfie	Not Selfie
2	Selfie	Selfie
3	Not Selfie	Not Selfie
4	Selfie	Selfie
5	Not Selfie	Selfie
6	Not Selfie	Not Selfie
7	Not Selfie	Not Selfie
8	Selfie	Selfie
9	Selfie	Not Selfie
10	Not Selfie	Not Selfie

Compute the Confusion Matrix for above and calculate Accuracy from it.

Briefly describe the process of building a Decision Tree. c)

CS CamScanner

SWE 335 Term Test 1

Introduction to

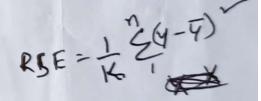
Marks: 25

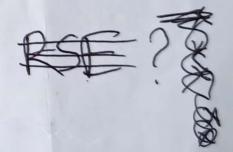
Time: 40 Min

Assume this dataset below is about the unit price of a fabric. There are 10 shades of colors and 10 levels of qualities. Using this, answer the following questions.

	Color	Quality	Price
1	7	5	65
2	3	7	38
3	5	8	51
4	8	1	38
5	9	3	55
6	5	4	43
7	4	0	25
8	2	6	33
9	8	7	71
10	6	4	51

Results of Multiple Regression						
n	10					
k	2					
R-Square	0.850694					
F	22.79061					
p-value	0.000497					





- 1. Write FIVE questions that can be asked about this fabric.
- 2. Draw a scatter plot with Quality and Price and write the findings from it.
- 3. Interpret the regression results given in the above table. What is the hypothesis here? Do you accept it?
- 4. Assume Price is the response variable here. Write the equations for all possible linear regression models using this data up to 2nd order polynomials.
- 5. What is cross validation? Write how a 5-fold cross validation be done with the dataset to find the best model? What is your evaluation criterion?

Shahjalal University of Science and Technology

Institute of Information and Communication Technology (IICT) SWE 3rd Year 2nd Semester Final Exam Dec-2021 (Session: 2018-19)
Course Code: SWE335 Course Title: Introduction to Data Science

Course Code: SWE335

Credits: 2 Time: 2 hrs Total Marks: 50

Group A

[Answer all the questions]

(x) (b) (c) (d) (e) (f) (g)	Answer any FIVE What is Data Science? What is a Statistical Model? What is a Loss Function? What is Supervised Learning? What is Structured data? Give What is a Confidence Interval What is web scrapping? What is Bias-Variance Trade-o	examples.				5x1=5
(x) (b) (c) (d)	Answer any FOUR Assume there are two features a model for it. Define Data, Distribution, Pop Write down the purpose(s) stat What are the things that we was purpose.	ulation, Samp istical model nt to visualize	ole and Bias. ing, with appi e about a data	ropriate example? Name the suita		
E)	What are the things to consider What is Overfitting? Write wh	while evaluate while over the state of the s	ating a model?			
	Answer any TWO i. Perform Linear Regression	n on the follow	wing dataset:			2x5=10 2
b)	ii. Re-estimate the responses iii. Calculate Loss and Fitness i. What is Hypothesis testing? ii. Assume a model $Y = B_0 + B_0$ if you get the following P-Value the predictors?	of the model. Why do we do	140 155 160 200 odel.		is the final model	1 2 1+1 2
ej	iii. Write which predictor(s) arei. Write all possible regression	Coefficien B_0 B_1 B_2 B_3 B_4 conot significate models for the	0.00 0.08 0.03 0.01 1.00			1
p	i. Write all possible regression olynomials and/or an interaction	Y 140 (155 (179 7) 192 7 200 7 215 7	edictors: X1	ata if you consid	der up to 2 nd order	3

ii. What is Cross Validation? Write how K-fold cross validation can be used to select a

suitable model.

Group B[Answer all the questions]

为为分分分分	Answer any FIVE What is Logistic Regularization Define Bayes' Inform What is an ROC curv Define Entropy. What is a Random For Define Eigen Value a How categorical variation	on? nation e? orest? and E	n Criterio	ctor.	computa	ation?			5x1=5
	Answer any FOUR								4x2.5=10
a)	Differentiate between What are the different	regi	ession, o	class	ification	n and clustering	ng.	1-0	
æ)	How regression can b	e do	ne by de	cisio	n trees?	ilid Noll-Para	metric Mode	lS?	
d)	How to avoid Overfit	ting	?						
<i>(</i> 2,	Define Entropy. Suppup 250, 500, 125, and say about this dice?	oose 1 125	you have times re	etos	sed a 4 :	faced dice 100 Calculate the	00 times whe Entropy for t	re 1, 2, 3 and 4 showed his dice. What can you	
A)	What is Imputation?	How	to impu	te m	issing v	alues in data?			
6. ,2)	i. What is Informat	a ca	r and hav	ve th	ie follov	ving models a	ıvailable. Bui	ld a Decision tree with	2x5=10 1 3
			Age	M	ileage	Road Tested	Buy		
		1	Recent	_		Yes	Buy		
		2	Recent		gh	Yes	Buy		
		3	Old Recent	Lo		No No	Don't buy		
	iii. What are the lim			_			Don't buy		1
(اطر	Suppose you have w have tested your class	ritten	a classit	fier t	o label mages a	your pictures on d got the following the Prediction	lowing predic	nd happy faces. You tions:	
				2	Нарру				
				3	Sad	Sad			
				4	Нарру				
				5	Sad	Нарру			
				6	Sad	Sad			
				7 8	Sad	Sad			
				9	Happy Happy				
				10	Sad	Sad			
	i. Make the Confusii. Calculate Accura			y, Pr			your classifie	er.	2
c)	i. What is Principa ii. Write the steps o iii. How can we get	of cor	nputing p	orino	ipal cor	nponents of N	I data with J f	eatures.	3 1 3 1