

**G.K. GUJAR MEMORIAL CHARITABLE TRUST'S
DR. ASHOK GUJAR TECHNICAL INSTITUTE'S
DR. DAULATRAO AHER COLLEGE OF ENGINEERING, KARAD**

UNIT TEST EXAMINATION - I / II / III

Name of Student Aavej Vahid Patel
 Class : Batch No - Dg Division :
 Roll No.
 Subject : Assessment 3 Date : 14/2/24
 Signature of Supervisor :

Que. No.	1	2	3	4	5	6	7	8	9	10	Total Marks
Marks Obtained											(53) 100

(Sign. of Subject Teacher)

Q. 1]

Sec. A .

→ list :- ① lists are mutable; you can modify add or remove element after creation.

② List defined using square brackets
 eg. list = [1, 2, 3]

Tuple :- ① Tuple are immutable; once created, their element cannot be changed

② Tuple use parentheses
 eg. Tuple = (1, 2, 3)

• Example use cases.

1) Use list when mutability matter

eg- managing a dynamic collection of user input
 ud = ['John', 25]

ud = 26 ≠ modify age

2) Use tuple for Immutable data

coordinates = (10, 20)

choose list when you need a collection that can be modified & tuple when you want a fixed unchangeable set of value.

Q.2. → ~~Write a Python~~
def factorial (n):
 If n == 0 or n == 1;
 return 1.
 else:
 return * factorial (n-1)

✓ = Example usage:
number 5
result = factorial (number)
print ("The factorial of {} number is: {}".format(number, result))

Q.3]
→ List comprehension is a concise way to
create lists in Python. It provides a more readable
& compact syntax for generating lists by specifying
the element to include & the condition
to meet. The basic structure is
[expression for item in iterable if condition]

eg :-
Here an example that generates a list of
squares for even numbers between 0 & 9 using
list comprehension.

Squares = [x**2 for x in range(10) if x%2==0]
print squares

The result a list [0, 4, 16, 36, 64]

Q4]

→ 1] Numpy :-

- Numpy is a powerful library for numerical computing in python. It provides support for large, multi-dimensional arrays & matrices, along with a collection of mathematical functions to operate on these array efficiently.

2] Pandas :

- Pandas is a data manipulation & analysis library that provides easy-to-use data structures, such as Dataframes, for working with structured data.
- It simplifies task like cleaning, filtering & aggregating data.
- widely used in data analysis & preparation, handling missing data, & performing operations on structured .

3] Matplotlib :-

- Matplotlib is a 2D plotting library for creating static, animated & interactive visualization in python. It allows user to create a wide variety of charts & graphs to visualize data.
- Ideal for data visualization in scientific computing, statistical analysis

Sec 'B'

Q.1]

→ Supervised learning :-

- Supervised learning is a type of machine learning where the algorithm is trained on a labeled dataset. In this approach, the model learns to map input data to corresponding output labels.

Example :-

- Task - Image classification [distinguishing between cats & dogs]
- Dataset - Collection of images with labels
- Training - The algorithm learns the patterns & features associated with each class
- Testing - Once trained, the model can predict the labels of new images.

3/

Unsupervised learning :

- Unsupervised learning involves training a machine learning algorithm on an unlabeled dataset.

Example :

- Task - clustering customer preferences in an e-commerce platform.
- Dataset - Purchasing history data without label categories.
- Algorithm - The model identifies natural grouping or cluster based on similarities in purchasing behaviour.
- Output :- Clusters or segments of customers with similar preference emerge from analysis.

Q.3]

→ Steps :-

① Data collection:-

Gather relevant data for the problem at hand. The quality & quantity of data significantly impact model performance.

② Data cleaning & processing:-

Clean & preprocess the data to handle missing values, outliers & inconsistencies. Convert data into a suitable format for training the model.

③ Feature Engineering:-

Create new features or transform existing ones to enhance the model's ability to capture patterns in the data. This step helps improve model performance.

④ Data splitting:-

Divide the dataset into training & testing sets.

⑤ Model selection:-

Choose a machine learning model. The selection depends upon nature & characteristics of data.

⑥ Model training:-

Train the selected model using the training data.

⑦ Model evaluation:-

Assess the model's performance.

⑧ Model deployment:-

Integrate model into production environment.

⑨ Monitoring & maintenance.

Q. 4]

→ Cross-validation :-

- cross-validation is resampling technique used in a machine learning to assess the performance of a model & to reduce the risks of overfitting or underfitting.
- It involves partitioning the dataset into subsets & training the model on these subsets & evaluating it on remaining subsets.

Importance of cross-validation

1] Better performance Estimation :- Cross-validation provides a more reliable estimate of how well a model will generalize to unseen data compared to a single train-test split.

2] Reduced overfitting or underfitting Risk :-
By evaluating the model on multiple subsets cross-validation helps in identifying if model is overfitting & underfitting.

example:

```
from sklearn.model_selection import KFold  
from sklearn.model_selection import cross_val_score  
from sklearn.ensemble import RandomForestClassifier
```

Example using KFold cross-validation with Random Forest classifier

model = RandomForestClassifier(),

kf = KFold (n_splits=5, shuffle=True,
random_state=42)

Performance cross-validation & calculate accuracy
accuracy_scores = cross_val_score(model, X, y,
cv=kf, scoring='accuracy')

Display average accuracy
print("Average Accuracy", accuracy_scores
mean())

Q5]

→ Regression Classification

1] The output variable has to be real value or continuous in nature 1] The data output variable in SML problem

2] Regression algorithm helps to map the input value & the continuous output variable which is discrete in nature

3] Regression algorithms are only used for data that is continuous.

4] linear regression, decision trees, neutral networks are common used for tasks eg.

predicting house prices based on features like square footage number of bedroom & location. Identifying whether an email is spam or not based on its content & characteristics.

Q. 6] →

The k-nearest neighbors (KNN) algorithm is a simple & versatile supervised machine learning algorithm used for both classification & regression tasks. It makes prediction based on the majority class or average of the k-nearest data points in the feature space.

How KNN works :

1] Training -

The algorithm stores the entire training dataset in memory.

2] Prediction (Classification) :-

- for a new data point, the algorithm identifies the k-nearest neighbors from the training set based on a distance metric.

3] Prediction :- (Regression)

for regression tasks the algorithm calculates the average of the target values of the k-nearest neighbors & assigns it to the new data point.

Main parameter

(1) Number of Neighbors (k) :-

(2) Distance Metric :-

Measure similarity or dissimilarity between data points.

(3) Weighting of Neighbors :-

- Defines the contribution of each neighbor to the prediction.

Section 'C'

- Q.1 → 1] Mean :- It is the average of a set of values calculated by summing all values & dividing by total number of values.

$$\text{Mean} = \bar{x} = \frac{\sum x_i}{n}$$

- 2] Median :- It is the middle value of sorted list of values. If the number of the middle value - The median is most appropriate when the data is skewed or contains outliers. It provides a better representation of the central tendency in such cases.

- 3] Mode :- It is the value that appears most frequently in a dataset. The mode is most appropriate for categorical or nominal data, where calculating the mean & median may not be meaningful. It is also useful for identifying the most common value in a dataset.

- Q.2, 1) In a normal distribution, approximately 95% of data falls within 2 standard deviation of the mean. This is known as 95% rule or the empirical rule.
- 2) If a data point is 2 standard deviation above the mean it falls within the top 2.5% of distribution (since approximately 95% of data is within 2 standard deviation of the mean). 5% of the data outside

This range divided equally b/w the upper & lower limit. The percentage of the data below 1.5 points in normal distribution is approx 100% - 2.5% = 97.5%.

Q.3)

\rightarrow The P-value in hypothesis testing represents the probability of observing a test statistic as extreme as or more extreme than one observed in sample data, under the assumption that the null hypothesis is true.

Q.3) Formulate hypothesis: In hypothesis testing we start with a null hypothesis (H_0), which represents the default assumption or the status quo, & alternative hypothesis (H_1 or H_a) which represents what we are trying to find evidence for.

i) Select a significance level: Before conducting the hypothesis test a significance level (α) is chosen, typically set at 0.05 or 0.01. This represents the threshold for rejecting null hypothesis.

Q.5.

$\rightarrow z = \frac{x - \mu}{\sigma}$ - x is value of random variable

μ is mean of the distribution & σ is the standard deviation of the distribution.

$$z = \frac{58 - 50}{8} = \frac{8}{8} = 1$$

So the Z-score for the value of $x = 58$ is 1.

number to the prediction.

Section 'D'

Q1) Overfitting is a common problem in machine learning where a model learns the learning data too well capturing noise & random function in the data instead of underlying pattern. As a result the model performs well on the training data but fails to generalize to new unseen data.

Causes :- i) Model complexity ii) Insufficient data
iii) Noise in the Data.

~~Q2) Effect :- i) Poor generalization ii) High variance~~
~~Mitigation Techniques for overfitting :-~~
i) Simplifying the model ii) cross validation
iii) Regularization iv) feature selection
v) Data Augmentation vi) Early stopping
vii) Ensemble methods

Q2] →

Support vector machine (SVM)

Definition :- SVM is a supervised machine learning algorithm used for classification & regression tasks. It works by finding the hyperplane that best separate different classes in the feature space

Role of the kernel in SVM

• Linear SVM :-

- In its simplest form, SVM uses a linear kernel & the decision boundary is a straight line. Linear means that each feature is multiplied by a weight and then summed up.

- The linear kernel is effective when the data is already linearly separable in the feature space.

• Non-linear SVM :-

- For non-linearly separable data SVM can be combined with non-linear kernels to map the input features into a higher-dimensional space.

- This transformation allows for finding a hyperplane in the transformed space that corresponds to a non-linear decision boundary in original space.