

WORKSHEET- 4 MACHINE LEARNING

In Q1 to Q7, only one option is correct, choose the correct option:

1. The value of correlation coefficient will always be:

- A) between 0 and 1
- B) greater than -1
- C) between -1 and 1
- D) between 0 and -1

Ans. C) between -1 and 1

2. Which of the following cannot be used for dimensionality reduction?

- A) Lasso Regularisation
- B) PCA
- C) Recursive feature elimination
- D) Ridge Regularisation

Ans. C) Recursive feature elimination

3. Which of the following is not a kernel in Support Vector Machines?

- A) linear
- B) Radial Basis Function
- C) hyperplane
- D) polynomial

Ans. C) hyperplane

4. Amongst the following, which one is least suitable for a dataset having non-linear decision boundaries?

- A) Logistic Regression
- B) Naïve Bayes Classifier
- C) Decision Tree Classifier
- D) Support Vector Classifier

Ans. A) Logistic Regression

5. In a Linear Regression problem, 'X' is independent variable and 'Y' is dependent variable, where 'X' represents weight in pounds. If you convert the unit of 'X' to kilograms, then new coefficient of 'X' will be?

(1 kilogram = 2.205 pounds)

A) $2.205 \times$ old coefficient of 'X'

B) same as old coefficient of 'X'

C) old coefficient of 'X' $\div 2.205$

D) Cannot be determined

Ans. C) old coefficient of 'X' $\div 2.205$

6. As we increase the number of estimators in ADABOOST Classifier, what happens to the accuracy of the model?

A) remains same

B) increases

C) decreases

D) none of the above

Ans. B) increases

7. Which of the following is not an advantage of using random forest instead of decision trees?

A) Random Forests reduce overfitting

B) Random Forests explain more variance in data than decision trees

C) Random Forests are easy to interpret

D) Random Forests provide a reliable feature importance estimate

Ans. C) Random Forests are easy to interpret

In Q8 to Q10, more than one options are correct, choose all the correct options:

8. Which of the following are correct about Principal Components?

A) Principal Components are calculated using supervised learning techniques

B) Principal Components are calculated using unsupervised learning techniques

C) Principal Components are linear combinations of Linear Variables.

D) All of the above

Ans. B) Principal Components are calculated using unsupervised learning techniques

C) Principal Components are linear combinations of Linear Variables.

9. Which of the following are applications of clustering?

- A) Identifying developed, developing and under-developed countries on the basis of factors like GDP, poverty index, employment rate, population and living index
- B) Identifying loan defaulters in a bank on the basis of previous years' data of loan accounts.
- C) Identifying spam or ham emails
- D) Identifying different segments of disease based on BMI, blood pressure, cholesterol, blood sugar levels.

Ans. A) Identifying developed, developing and under-developed countries on the basis of factors like GDP, poverty index, employment rate, population and living index

- B) Identifying loan defaulters in a bank on the basis of previous years' data of loan accounts.
- D) Identifying different segments of disease based on BMI, blood pressure, cholesterol, blood sugar levels.

10. Which of the following is(are) hyper parameters of a decision tree?

- A) max_depth
- B) max_features
- C) n_estimators
- D) min_samples_leaf

Ans. A) max_depth

- B) max_features
- D) min_samples_leaf

Q11 to Q15 are subjective answer type questions, Answer them briefly.

11. What are outliers? Explain the Inter Quartile Range (IQR) method for outlier detection.

Ans. Outliers badly affect the mean and standard deviation of the dataset. These may statistically give erroneous results. Most machine learning algorithms do not work well in the presence of outlier. So it is desirable to detect and remove outliers. Outliers are highly useful in anomaly detection like fraud detection where the fraud transactions are very different from normal transactions.

IQR is the range between the first and the third quartiles namely Q1 and Q3.

$$\text{IQR} = Q3 - Q1.$$

The data points which fall below $Q1 - 1.5 \text{ IQR}$ or above $Q3 + 1.5 \text{ IQR}$ are outliers.

12. What is the primary difference between bagging and boosting algorithms?

Ans. In Bagging the result is obtained by averaging the responses of the N learners (or majority vote). However, Boosting assigns the second set of weights, this time for the N classifiers, in order to take a weighted average of their estimates. In the Boosting training stage, the algorithm allocates weights to each resulting model. A learner with a good classification result on the training data will be assigned a higher weight than a poor one. So when evaluating a new learner, Boosting needs to keep track of learners' errors, too.

13. What is adjusted R² in linear regression. How is it calculated?

Ans. The adjusted R-squared is a modified version of R-squared that has been adjusted for the number of predictors in the model. The adjusted R-squared increases only if the new term improves the model more than would be expected by chance. It decreases when a predictor improves the model by less than expected by chance. The

adjusted R-squared can be negative, but it's usually not. It is always lower than the R-squared.

The adjusted R^2 is computed as

$$R_{adj}^2 = 1 - \frac{Var(e_i)/(n - k - 1)}{Var(y_i)/(n - 1)} = 1 - \frac{Var(e_i)}{Var(y_i)} \times \frac{n - 1}{n - k - 1}$$

where n is the number of cases used to fit the model and k is the number of predictor variables in the model.

14. What is the difference between standardisation and normalisation?

Ans. Difference between Standardisation and normalization are:

Normalization or Min-Max Scaling is used to transform features to be on a similar scale.

The new point is calculated as:

$$X_{\text{new}} = (X - X_{\text{min}})/(X_{\text{max}} - X_{\text{min}})$$

This scales the range to $[0, 1]$ or sometimes $[-1, 1]$. Geometrically speaking, transformation squishes the n -dimensional data into an n -dimensional unit hypercube. Normalization is useful when there are no outliers as it cannot cope up with them. Usually, we would scale age and not incomes because only a few people have high incomes but the age is close to uniform.

Standardization or Z-Score Normalization is the transformation of features by subtracting from mean and dividing by standard deviation. This is often called as Z-score.

$$X_{\text{new}} = (X - \text{mean})/\text{Std}$$

Standardization can be helpful in cases where the data follows a Gaussian distribution. However, this does not have to be necessarily true. Geometrically speaking, it translates the data to the mean vector of original data to the origin and squishes or expands the points if std is 1 respectively. We can see that we are just changing mean and standard deviation to a standard normal distribution which is still normal thus the shape of the distribution is not affected. Standardization does not get affected by outliers because there is no predefined range of transformed features.

15. What is cross-validation? Describe one advantage and one disadvantage of using cross-validation.

Ans. Cross-validation is a resampling procedure used to evaluate machine learning models on a limited data sample. The procedure has a single parameter called k that refers to the number of groups that a given data sample is to be split into. As such, the procedure is often called k-fold cross-validation.

Advantages of Cross-Validation: Reduces Overfitting: In Cross-Validation, we split the dataset into multiple folds and train the algorithm on different folds. This prevents our model from overfitting the training dataset. So, in this way, the model attains the generalization capabilities which are a good sign of a robust algorithm.

Disadvantages of Cross-Validation: Increases Training Time: Cross-Validation drastically increases the training time. Earlier you had to train your model only on one training set, but with Cross-Validation, you have to train your model on multiple training sets