

Machine Learning-4

Answers:

1. C
2. D
3. C
4. A
5. C
6. B
7. D
8. B, C
9. A, B, C, D
10. A, D
11. Outliers are those data points that are significantly different from the rest of the dataset. They are often abnormal observations that skew the data distribution, and arise due to inconsistent data entry, or erroneous observations. IQR is used to measure variability by dividing a data set into quartiles. The data is sorted in ascending order and split into 4 equal parts. Q1, Q2, Q3 called first, second and third quartiles are the values which separate the 4 equal parts. Most commonly used method to detect outliers is visualization. We can use the IQR method of identifying outliers to set up a “fence” outside of Q1 and Q3. Any values that fall outside of this fence are considered outliers. To build this fence we take 1.5 times the IQR and then subtract this value from Q1 and add this value to Q3. This gives us the minimum and maximum fence posts that we compare each observation to. Any observations that are more than 1.5 IQR below Q1 or more than 1.5 IQR above Q3 are considered outliers. This is the method that Minitab uses to identify outliers by default.
12. In Bagging the result is obtained by averaging the responses of the N learners (or majority vote). However, Boosting assigns a second set of weights, this time for the N classifiers, in order to take a weighted average of their estimates. While they are built independently for Bagging, Boosting

tries to add new models that do well where previous models fail. Only Boosting determines weights for the data to tip the scales in favor of the most difficult cases. It is an equally weighted average for Bagging and a weighted average for Boosting, more weight to those with better performance on training data. Only Boosting tries to reduce bias. On the other hand, Bagging may solve the over-fitting problem, while Boosting can increase it.

13. 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. Every time you add an independent variable to a model, the R-squared increases, even if the independent variable is insignificant. It never declines. Whereas Adjusted R-squared increases only when independent variable is significant and affects dependent variable.

$$\text{Adjusted } R^2 = \left\{ 1 - \left[\frac{(1 - R^2)(n - 1)}{(n - k - 1)} \right] \right\}$$

14.

Normalisation	Standardisation
Scaling is done by the highest and the lowest values.	Scaling is done by mean and standard deviation.
It is applied when the features are of separate scales.	It is applied when we verify zero mean and unit standard deviation.
Scales range from 0 to 1	Not bounded
Affected by outliers	Less affected by outliers
It is applied when we are not sure about the data distribution	It is used when the data is Gaussian or normally distributed
It is also known as Scaling Normalization	It is also known as Z-Score

15. Cross-validation is a technique in which we train our model using the subset of the data-set and then evaluate using the complementary subset of the

data-set. Advantages of cross-validation: More accurate estimate of out-of-sample accuracy. More “efficient” use of data as every observation is used for both training and testing. There is a disadvantage because the cross validation process can become a lengthy one. It depends on the number of observations in the original sample and your chosen value of ‘p.’