

## Worksheet 5

### Machine learning assignment 5

#### Answers :

1. R-squared is a commonly used measure of the goodness of fit of a regression model. It is a proportion (between 0 and 1) that expresses how much of the variation in the dependent variable is explained by the independent variables in the model. RSS (residual sum of squares) is the sum of the squared differences between the predicted values and the actual values. It is used to measure the deviation of the model from the actual data.

In general, R-squared is a better measure of goodness of fit because it is easily interpretable and it ranges between 0 and 1, making it easy to compare models. However, it does have some limitations. For example, it can be artificially inflated if you add more variables to the model, even if they are not actually related to the response variable. RSS can also be used to compare models but it can be difficult to interpret because it is a squared value and it does not have a fixed range of possible values. In conclusion, R-squared is a more commonly used measure of goodness of fit in regression

2. In a regression analysis, TSS (Total Sum of Squares) is a measure of the total variation in the dependent variable (y). It is calculated as the sum of the squared differences between the mean of the dependent variable and each individual value of the dependent variable.

ESS (Explained Sum of Squares) is a measure of the variation in the dependent variable that is explained by the independent variables in the model. It is calculated as the sum of the squared differences between the predicted values (based on the model) and the mean of the dependent variable.

RSS (Residual Sum of Squares) is a measure of the variation in the dependent variable that is not explained by the independent variables in the model. It is calculated as the sum of the squared differences between the actual values of the dependent variable and the predicted values (based on the model).

The equation relating these three metrics is:

$$TSS = ESS + RSS$$

This equation states that the total variation in the dependent variable (TSS) is equal to the variation explained by the independent variables (ESS) plus the variation not explained by the independent variables (RSS).

In summary, TSS represents the total deviation of the y variable, ESS represents the deviation that is explained by the model while RSS represents the deviation that is not explained by the model.

3. Regularization is a technique used in machine learning to prevent overfitting. Overfitting occurs when a model is overly complex and is able to fit the noise or random fluctuations in the training data, rather than just the underlying pattern. This can result in a model that performs well on the training data but poorly on new, unseen data.

Regularization addresses this issue by adding a penalty term to the objective function that the model is trying to optimize. This penalty term discourages the model from assigning too much importance to any one feature or parameter, which helps to prevent overfitting.

There are two main types of regularization used in machine learning:

L1 regularization, also known as Lasso regularization, adds a penalty term to the objective function that is proportional to the absolute value of the coefficients. This results in some coefficients becoming exactly zero, effectively eliminating them from the model.

L2 regularization, also known as Ridge regularization, adds a penalty term to the objective function that is proportional to the square of the coefficients. This results in all coefficients becoming smaller, but none of them become exactly zero.

Regularization is particularly useful when dealing with high-dimensional or sparse data, where the number of features or variables is much larger than the number of observations. It's also useful when we have a lot of correlated features in our data, regularization can help us in selecting the most important features. Regularization can also help in improve the model's generalization ability and make model more robust to noise.

4. The Gini impurity index is a measure of the impurity or disorder of a set of data. It is commonly used in decision tree-based algorithms such as the CART (Classification and Regression Trees) algorithm to determine the best split point for a given dataset.

The Gini impurity index is calculated by taking the probability of a random sample from the dataset being labeled incorrectly if it were randomly labeled according to the class distribution of the dataset. The formula for the Gini impurity index is:

$$\text{Gini} = 1 - (p_1^2 + p_2^2 + \dots + p_n^2)$$

Where  $p_i$  is the proportion of the  $i$ th class in the dataset.

For example, if a dataset has two classes, A and B, with proportions of 0.6 and 0.4 respectively, the Gini impurity index would be:

$$\text{Gini} = 1 - (0.6^2 + 0.4^2) = 0.48$$

The Gini impurity ranges between 0 and 1, and a low Gini impurity value indicates a high degree of purity or order, while a high Gini impurity value indicates a high degree of disorder or impurity. In decision tree algorithms, the Gini impurity index is used to determine the best split point at each node of the tree by evaluating all possible split points

and choosing the one that results in the greatest decrease in Gini impurity. The goal is to split the data into subsets with the highest possible purity, which will lead to the most accurate predictions

5. Yes, unregularized decision trees are prone to overfitting. Decision trees are a type of model that recursively splits the data into subsets based on the feature that results in the greatest decrease in impurity (such as Gini impurity index). As the tree grows deeper and deeper, it can continue to split the data into increasingly specific subsets, fitting the noise and random fluctuations in the training data, rather than just the underlying pattern.

This leads to a model that has a low training error, but performs poorly on new, unseen data because it is too complex and sensitive to the noise in the training data. The model is memorizing the training data rather than generalizing to new unseen data.

Regularization techniques such as pruning, early stopping, or adding a penalty term to the objective function can help to prevent overfitting in decision trees by limiting the complexity of the model. Pruning is the process of removing branches of a decision tree that do not contribute to the prediction accuracy, this can help in reducing the overfitting.

Another regularization technique is early stopping, it is used to stop growing the tree when the accuracy on validation data stops improving. Another method is adding a penalty term to the objective function, such as L1 or L2 regularization, this can help to reduce the magnitude of the coefficients and prevent the model from assigning too much importance to any one feature.

6. An ensemble technique in machine learning is a method that combines multiple models in order to improve the overall performance and robustness of the final model. Ensemble techniques are based on the idea that multiple models working together can often achieve better results than a single model working alone. There are several ensemble techniques that can be used in machine learning.

**Bagging (Bootstrap Aggregating):** Bagging is an ensemble technique that involves training multiple models independently on different subsets of the data, and then combining their predictions through a majority vote or average. Bagging is particularly useful for reducing the variance of high-variance models such as decision trees.

**Boosting:** Boosting is an ensemble technique that involves training multiple models sequentially, where each model tries to correct the mistakes of the previous model. The final prediction is made by combining the predictions of all the models. Boosting is particularly useful for reducing the bias of low-bias models such as linear regression.

**Stacking:** Stacking is an ensemble technique that involves training multiple models independently and then using their predictions as inputs to a final model (called a meta-model) that makes the final prediction. Stacking is particularly useful for combining models with complementary strengths and weaknesses.

**Random Forest:** Random Forest is a specific type of ensemble technique that combines the predictions of multiple decision trees. The decision trees are trained independently on different subsets of the data, and the final prediction is made by averaging the predictions of all the trees. Random Forest is widely used for both classification and regression problems and known to work well for high-dimensional data.

Ensemble techniques are generally considered to improve performance compared to a single model, but they also increase the complexity and computational cost of the model. Therefore, it is important to consider the trade-off between performance improvement and computational cost when deciding whether to use an ensemble technique.

7. Bagging and Boosting are both ensemble techniques in machine learning that combine multiple models to improve the overall performance of the final model. However, the key difference between the two is in the way they train and combine the models.

Bagging (Bootstrap Aggregating) is an ensemble technique that involves training multiple models independently on different subsets of the data. The subsets are created by randomly sampling the data with replacement (bootstrapping). The final prediction is made by combining the predictions of all the models through a majority vote or average. Bagging is particularly useful for reducing the variance of high-variance models such as decision trees.

Boosting, on the other hand, is an ensemble technique that involves training multiple models sequentially. Each model is trained to correct the mistakes of the previous model. The final prediction is made by combining the predictions of all the models. Boosting is particularly useful for reducing the bias of low-bias models such as linear regression.

In summary, Bagging uses multiple models independently and combines their predictions through averaging or voting, while Boosting uses multiple models sequentially and combines their predictions by assigning a weight to each model's prediction. This sequential training process aims to correct the mistakes of the previous model in each iteration and to improve the performance.

8. In Random Forests, each tree is trained on a different subset of the data, known as a bootstrap sample. The bootstrap sample is created by randomly sampling the data with replacement. Since some observations may not be included in the bootstrap sample used to train a particular tree, those observations are referred to as out-of-bag (OOB) observations. The out-of-bag error is the error rate calculated on the OOB observations for each tree.

The OOB error is an estimate of the error rate of the entire random forest without using cross-validation or a separate test set. It is also a way to estimate the generalization error of the model, by using an unbiased estimate of the error rate.

It is calculated by predicting the OOB observations for each tree in the forest, and then averaging the error rates across all the trees. This gives an estimate of the error rate of the random forest without using a separate test set.

In general, the OOB error provides a good estimate of the generalization error for the random forest, but it may not always be the best estimate. It can be useful in determining the optimal number of

trees in the forest, as the OOB error will usually decrease as the number of trees in the forest increases, and then eventually stabilize or decrease slightly.

It is important to note that the OOB error should not be used as a replacement of cross validation, it is a way to have an idea of the generalization error of the model, but it can be not always accurate, especially when the sample size is small or the data distribution is not homogeneous.

9. K-fold cross-validation is a technique used to evaluate the performance of a machine learning model. It is used to estimate the generalization error of a model, which is the error rate of the model when applied to new, unseen data.

The basic idea behind K-fold cross-validation is to partition the data into K equally-sized "folds" (or subsets), then train the model K times, each time using a different fold as the test set and the remaining K-1 folds as the training set. The performance of the model is then averaged across all K iterations. This process is also called rotation estimation.

Here are the steps of a typical K-fold cross-validation procedure:

1. Split the data into K equal-sized "folds"
2. For each fold, use it as a test set and the remaining K-1 folds as the training set
3. Train the model on the training set and evaluate its performance on the test set
4. Repeat steps 2 and 3 K times, with a different fold as the test set each time
5. Average the performance of the model across all K iterations

The value of K should be chosen based on the size of the data set, with a larger value of K providing a more robust estimate of the generalization error, but at the cost of a longer computational time. A common choice for K is 10, which is called 10-fold cross-validation.

K-fold cross-validation is a widely used technique for evaluating the performance of a machine learning model, as it provides a robust estimate of the generalization error and is relatively simple to implement.

10. Hyperparameter tuning, also known as hyperparameter optimization or hyperparameter selection, is the process of selecting the best set of hyperparameters for a machine learning model. Hyperparameters are parameters that are not learned from the data during training, but are set prior to training. Examples of hyperparameters include the learning rate in gradient descent, the number of trees in a random forest, or the regularization term in a linear regression model.

The process of hyperparameter tuning is done to improve the performance of the model on unseen data. Hyperparameters control the behaviour and the capacity of the model, so finding the best set of hyperparameters can lead to better generalization performance.

There are several techniques for hyperparameter tuning, including:

Grid search: Grid search is a simple and widely used technique for hyperparameter tuning. It involves specifying a set of possible values for each hyperparameter and training the model for all possible combinations of the hyperparameters. The best set of hyperparameters is then selected based on the performance of the model on a validation set.

Random search: Random search is similar to grid search, but instead of trying all possible combinations of the hyperparameters, a random set of combinations is tried. Random search has been shown to be more efficient than grid search, especially when the number of hyperparameters is large.

Bayesian optimization: Bayesian optimization is a more sophisticated technique for hyperparameter tuning. It uses a probabilistic model to model the relationship between the hyperparameters and the performance of the model, and then chooses the next set of hyperparameters to evaluate based on the current estimates of the model.

11. A large learning rate in Gradient Descent can cause several issues that can negatively impact the performance of the model.

1. Divergence: A large learning rate can cause the cost function to diverge, meaning that the cost increases instead of decreasing. When the learning rate is too high, the algorithm will take large steps and may overshoot the minimum of the cost function, causing the cost to increase instead of decrease.
2. Oscillation: A large learning rate can cause the cost function to oscillate, meaning that the cost fluctuates instead of decreasing monotonically. This can happen when the learning rate is too high, causing the algorithm to take steps that are too large and cause the cost to oscillate around the minimum.
3. Local Minima: A large learning rate can cause the algorithm to converge to a local minimum instead of the global minimum. When the learning rate is too high, the algorithm may overshoot the global minimum and converge to a local minimum that is not the optimal solution.
4. Slow Convergence: A large learning rate can cause the algorithm to converge slowly, meaning that it takes a long time for the cost to reach a stable value. When the learning rate is too high, the algorithm takes large steps that may not be efficient, causing the algorithm to converge slowly.
5. Noise Sensitivity: A large learning rate can cause the algorithm to be sensitive to noise in the data, meaning that small variations in the data can cause large variations in the cost. When the learning rate is too high, the algorithm takes large steps that may be sensitive to noise in the data, causing the algorithm to converge to a suboptimal solution.

In summary, a large learning rate in Gradient Descent can cause the algorithm to diverge, oscillate, converge to a local minimum, converge slowly and be sensitive to noise in the data, which can negatively impact the performance of the model. To avoid these issues, it is usually recommended to use a small learning rate, and gradually increase the learning rate as the algorithm converges. This is known as annealing the learning rate. Additionally, adaptive learning rate methods like Adam, Adagrad, and Adadelat can automatically adapt the learning rate during training and can help to prevent these issues.

12. Logistic regression is a linear model, so it is not well-suited for classification of non-linear data.

The logistic regression model uses a linear combination of the input features to make a prediction, which is then transformed into a probability using the sigmoid function. The decision

boundary of the logistic regression model is a linear hyperplane, which can only separate data that is linearly separable. Non-linear data, however, cannot be separated by a linear hyperplane and thus logistic regression is not able to classify non-linear data well.

There are several ways to overcome this limitation:

1. **Feature Engineering:** By transforming the input features using non-linear functions, the data can be transformed into a higher dimensional space where it becomes linearly separable.
2. **Non-Linear Model:** Non-linear models like Decision Trees, Random Forest, Neural Networks and Support Vector Machines (SVMs) are able to capture the non-linearity in the data and are able to classify non-linear data well.
3. **Non-Linear Logistic Regression :** Non-Linear logistic regression is a variant of logistic regression, where the input features are transformed by non-linear functions before being passed to the model. This allows the model to learn non-linear decision boundaries, and thus classify non-linear data.

In summary, Logistic Regression is a linear model, so it is not well-suited for classification of non-linear data. However, by using feature engineering, or by using non-linear models and non-linear logistic regression, it is possible to classify non-linear data.

13. Adaboost and Gradient Boosting are both ensemble techniques used to improve the performance of weak learners by combining multiple models to form a stronger model. However, they differ in the way they combine the weak learners.

Adaboost (Adaptive Boosting) is an iterative algorithm that adjusts the weights of the training examples at each iteration, so that the next weak learner focuses more on the misclassified examples from the previous iteration. After each iteration, the weights are re-adjusted, and the process is repeated until a stopping criterion is met. The final model is a weighted combination of all the weak learners, where the weights are determined by the performance of each weak learner.

Gradient Boosting, on the other hand, is a forward-stage-wise algorithm. The idea behind gradient boosting is to fit a new model to the negative gradient of the loss function with respect to the current model's predictions. Gradient Boosting uses decision trees as the base learners, and at each iteration, it fits a new decision tree to the negative gradient of the loss function. The final model is a sum of all the decision trees.

In summary, Adaboost is an iterative algorithm that adjusts the weights of the training examples at each iteration to improve the performance of the weak learners. Gradient Boosting is a forward-stage-wise algorithm that fits a new model to the negative gradient of the loss function using decision tree as base learners. The final model is a sum of all decision trees.

14. The bias-variance trade-off is a fundamental concept in machine learning that refers to the trade-off between a model's ability to fit the training data well (low bias) and its ability to generalize well to new data (low variance).

Bias refers to the difference between the predictions of a model and the true values of the target variable. A model with high bias is said to be underfitting, which means that it is not able to capture the underlying patterns in the data. This results in high error on both the training and test data.

Variance, on the other hand, refers to the variability of a model's predictions for a given input. A model with high variance is said to be overfitting, which means that it is too complex and is fitting to the noise in the data. This results in low error on the training data but high error on the test data.

The goal in machine learning is to find a balance between bias and variance to achieve a good generalization performance on unseen data. A model with low bias and low variance is said to have good generalization performance. A model with high bias and high variance is said to have poor generalization performance.

To balance the trade-off we can use regularization, early stopping, increasing or decreasing complexity of the model, cross-validation, or ensemble methods.

In summary, the bias-variance trade-off refers to the trade-off between a model's ability to fit the training data well (low bias) and its ability to generalize well to new data (low variance). A good machine learning model should have a balance of bias and variance to achieve good generalization performance on unseen data.

15. Linear kernel: The linear kernel is the simplest kernel function in Support Vector Machine (SVM). It is used when the data is linearly separable. The linear kernel computes the dot product of the input features. In other words, it considers the input features as they are, without applying any transformation.

RBF (Radial Basis Function) kernel: The RBF kernel is a non-linear kernel function that is used when the data is not linearly separable. The RBF kernel applies a non-linear transformation to the input features by mapping them into a higher-dimensional space where a linear boundary can be found. The RBF kernel is defined as the exponential of the negative Euclidean distance between the input features.

Polynomial kernel: The polynomial kernel is also a non-linear kernel function that is used when the data is not linearly separable. It maps the input features into a higher-dimensional space and finds a polynomial boundary to separate the data. The polynomial kernel is defined as the dot product of the input features raised to a certain power.

In summary, the linear kernel is used when the data is linearly separable, the RBF kernel applies a non-linear transformation to the input features by mapping them into a higher-dimensional space and finds a linear boundary, and the polynomial kernel maps the input features into a higher-dimensional space and finds a polynomial boundary to separate the data.



## Statistics Worksheet 5

### **Answers:**

1. d) Expected
2. c) Frequencies
3. c) 6
4. b) Chisquared distribution
5. c) F Distribution
6. b) Hypothesis
7. a) Null Hypothesis
8. a) Two tailed
9. b) Research Hypothesis
10. a) np