# Logistic Regression

Matthew McDonald

# Machine Learning Poll Questions

# Question 1

Which of the following best describes the two main types of predictive modeling used in machine learning to analyze data and make predictions?

a. Parametric and Non-Parametric Learning

b. Supervised and Unsupervised Learning

c. Regression and Classification

d. Batch Learning and Online Learning

# Question 2

In the context of supervised learning, which of the following options correctly identifies the two primary types of tasks that algorithms aim to solve?

a. Clustering and Association

b. Regression and Classification

c. Dimensionality Reduction and Feature Extraction

d. Reinforcement and Batch Learning

# Question 3

In the context of machine learning, which of the following pairs of datasets are crucial for developing, validating, and assessing the performance of predictive models?

a. Primary and Secondary Datasets

b. Quantitative and Qualitative Datasets

c. Training and Testing Datasets

d. Cross-validation and Bootstrap Datasets

# Question 4

In programming, particularly when using languages like R for statistical analysis or simulations, which function is crucial for ensuring the reproducibility of results that involve random number generation?

a.  randomize()

b.  set.seed()

c.  initRandom()

d.  random.seed()

# Question 5

In the preparation of datasets for machine learning, which combination is crucial for effectively training models and ensuring that samples are representative, especially in cases where the target variable classes are imbalanced?

a. Stratified Sampling

b. Systemic Sampling

c. Simple Random Sampling

d. Cluster Sampling

# Question 6

Which of the following metrics is commonly used to assess the performance of regression models by measuring the average magnitude of the errors between the predicted values and the actual values?

a. Accuracy

b. F1 Score

c. Root Mean Squared Error (RMSE)

d. Precision

# Question 7

Which of the following metrics is widely regarded as a robust measure for evaluating the performance of classification models, especially in distinguishing between the model's ability to correctly predict different classes?

a. Confusion Matrix

b. Precision at K

c. Area Under ROC Curve (AUC-ROC)

d. Mean Absolute Error (MAE)

# Question 8

When dealing with machine learning datasets that exhibit high correlation among predictor variables, which feature engineering technique is particularly beneficial to reduce multicollinearity and improve model performance?

a.  Categorical Encoding

b.  Min-Max Scaling

c.  Principal Component Analysis (PCA)

d.  Bucketing/Binning

# Question 9

Which of the following are the most common resampling techniques used in machine learning to estimate model accuracy on a limited data sample?

a. Principal Component Analysis and Linear Discriminant Analysis

b. Gradient Boosting and Random Forests

c. Cross-Validation and Repeated Cross Validation

d. Standardization and Normalization

# Question 10

Why is resampling an important technique during the parameter tuning phase of building a machine learning model?

a. To increase the speed of the model training process

b. To ensure that the model can handle large datasets

c. To avoid overfitting by assessing how well the model generalizes to unseen data

d. To reduce the computational resources required for model training

# Linear and Logistic Regression

# Linear Regression Assumptions

- **Linearity**: The relationship between the dependent and independent variables is linear. This means that the change in the dependent variable due to a one-unit change in any one of the independent variables is constant.

- **Independence**: Observations are independent of each other. This assumption is crucial for the standard errors of the coefficients to be valid.

- **Homoscedasticity**: The variance of the error terms (residuals) is constant across all levels of the independent variables. If the variance of the residuals increases or decreases with the independent variable(s), the model exhibits heteroscedasticity, violating this assumption.

- **Normal Distribution of Errors**: For any fixed value of the independent variable(s), the dependent variable is normally distributed. In practice, this assumption is often interpreted as requiring the residuals of the model to be normally distributed.

# Other Linear Regression Assumptions

- **No or Little Multicollinearity**: Multicollinearity occurs when independent variables in a regression model are highly correlated. This condition can make it difficult to ascertain the individual effect of each independent variable on the dependent variable because it becomes hard to distinguish their individual contributions.

- **No Auto-correlation**: The residuals from the model are not correlated with each other. This assumption is particularly important in time series data, where the residual from one point could be correlated with the residual from a previous point.

- **No Perfect Multicollinearity**: In the regression model, none of the independent variables is a perfect linear function of any other variable(s). This would mean that the matrix of independent variables is of full rank.

# Logistic Regression

- **Definition:** Logistic Regression is a statistical method for analyzing a dataset in which there are one or more independent variables that determine an outcome. The outcome is measured with a dichotomous variable (where there are only two possible outcomes).

- **Purpose:** It is used for classification tasks, specifically binary classification.

- **How it works:**

  - Utilizes the logistic function to model a binary outcome, which can vary between 0 and 1.

  - The logistic function (also called the sigmoid function) ensures that the output can be interpreted as a probability.

# The Logistic Function

**What is the Logistic Function?**

- The logistic function, also known as the sigmoid function, is a mathematical function that outputs a value between 0 and 1.

- Formula: $f(x) = \frac{1}{1+e^{-x}}$

**Key Characteristics:**

- **S-shaped Curve:** The function has an S-shaped curve (sigmoid curve), making it ideal for binary classification.

- **Output Range:** The output ranges from 0 to 1, which can be interpreted as probabilities.

- **Non-linearity:** Introduces non-linearity to models, allowing them to learn complex patterns.

# The Logistic Function



Logistic Function