

Tuesday, November 16, 2021

6:07 PM

GIT INTRO QUESTIONS
GIT ADD,commit,push

LABS 1-7 REVIEW THOSE QUESTIONS

What is Data Science?

Data Science is an interdisciplinary field that involves the use of statistical and computational insights from data.

What is a PIVOT Table?

How do we use pivot tables to help us do
data visualization

Answer business questions

feature engineer?

Aggregate Functions

What are the steps involved in a Data Science project?

Be aware what the code looks like for each of the machine Learning steps

- Define the problem
 - How do we identify a bird using a computer
- Gather data
 - Bird pictures
- Preprocess the data
 - Label our images extract key components into features
-
- Create a model
- Train Model
- Evaluate the model

onal methods to extract

- Evaluate the model
- Deploy the model
- Monitor the performance of the model
- Repeat until desired performance is reached

BE ABLE TO IDENTIFY WHAT STEP IS DOING WHAT IN SCIKITLEARN

What is Machine Learning?

Answer: Machine Learning is a subset of Artificial Intelligence that involves the use of algorithms and models to enable machines to learn from data and make predictions or decisions without being explicitly programmed.

No if else etc...

What are the different types of Machine Learning algorithms?

- Supervised Learning
- Unsupervised Learning
- Reinforcement Learning

What is the difference between supervised and unsupervised learning?

In supervised learning, the algorithm learns from labeled data, whereas in unsupervised learning, the algorithm learns from unlabeled data.

What is classification in Machine Learning?

Answer: Classification is a type of supervised learning algorithm that involves predicting a categorical outcome variable based on one or more input variables.

What is regression in Machine Learning? Answer: Regression is a type of supervised learning algorithm that involves predicting a continuous or numerical outcome variable based on one or more input variables.

What is clustering in Machine Learning? Answer: Clustering is a type of unsupervised learning algorithm that involves grouping similar data points together based on their characteristics.

What is feature engineering? Answer: Feature engineering is the process of selecting and transforming input variables or features used in a machine learning model to improve its performance.

What is overfitting in Machine Learning?

Answer: Overfitting occurs when a machine learning model is too complex and fits the training data too closely, resulting in poor performance on new or unseen data.

gorithms and statistical
t being explicitly

learning, the algorithm

a categorical or discrete

ed learning algorithm
e or more input

arning algorithm that

nd transforming the
ce.

aining data too well, but

What is cross-validation?

Answer: Cross-validation is a technique used to assess the performance of a machine learning model by splitting the data into training and testing sets multiple times and averaging the results.

What is hyperparameter tuning?

Answer: Hyperparameter tuning is the process of selecting the best set of hyperparameters for a machine learning algorithm to optimize its performance.

What is an ROC curve?

Answer: An ROC (Receiver Operating Characteristic) curve is a graphical representation of the true positive rate and the false positive rate of a machine learning model at different thresholds.

What is an AUC score?

1. **Answer: AUC (Area Under the Curve) is a metric used to evaluate the performance of a machine learning model based on the area under the ROC curve.**

2. **What is a confusion matrix?**

3. **Answer: A confusion matrix is a table that shows the number of true positives, false positives, true negatives, and false negatives of a machine learning model.**

What is the difference between a training set and a testing set? Answer: A training set is used to train a machine learning model, whereas a testing set is used to evaluate the performance of the trained model.

What is a decision tree?

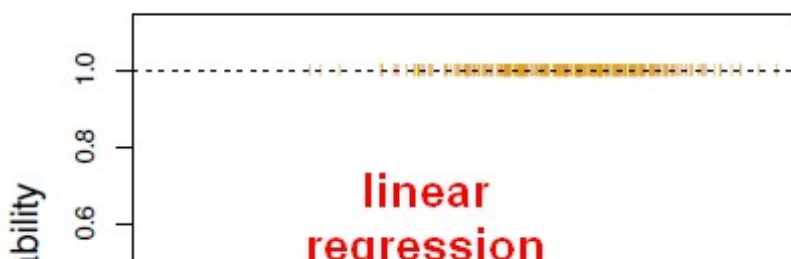
Answer: A decision tree is a tree-like model used in machine learning to make decisions based on a series of rules or conditions.

What is random forest?

Answer: Random forest is an ensemble learning method that uses multiple decision trees to improve the accuracy and reduce the variance of a machine learning model.

What is logistic regression?

Logistic regression is a type of regression analysis used to predict the probability of a binary or categorical outcome based on one or more input variables.



learning model by splitting

ers for a machine learning

**the trade-off between the
olds.**

a machine learning

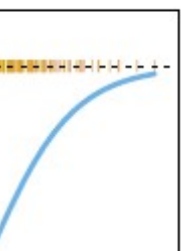
s, true negatives, and

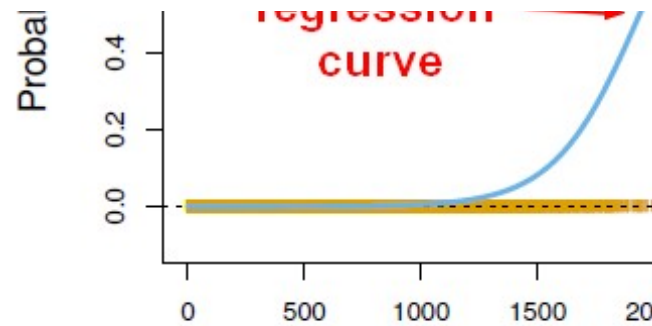
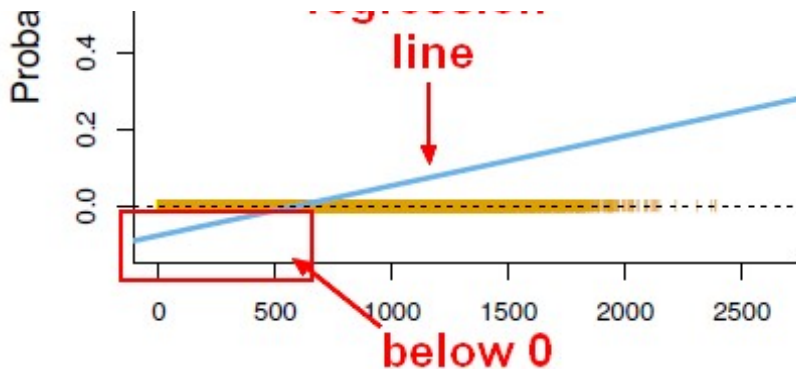
**is used to train a
the model.**

s based on a series of

improve the accuracy and

ome variable based on





What is K-means clustering?

K-means clustering is a type of unsupervised learning algorithm used to partition a set of data points into clusters based on their characteristics.

What is support vector machine (SVM)?

Support vector machine (SVM) is a supervised learning algorithm used to classify data by finding a hyperplane that maximally separates the classes.

What is a neural network?

A neural network is a set of interconnected nodes or neurons that can learn from data and make decisions.

What is NumPy?

What is the difference between Pandas Series and Pandas DataFrame?

Answer: Pandas Series is a one-dimensional labeled array that can hold any data type, whereas Pandas DataFrame is a two-dimensional labeled data structure that can hold multiple Pandas Series.

SERIES

DATAFRAME2D MATRIX

SERIES	SERIES	SERIES



of data points into K

ling the hyperplane that

to make predictions or

Pandas DataFrame is a

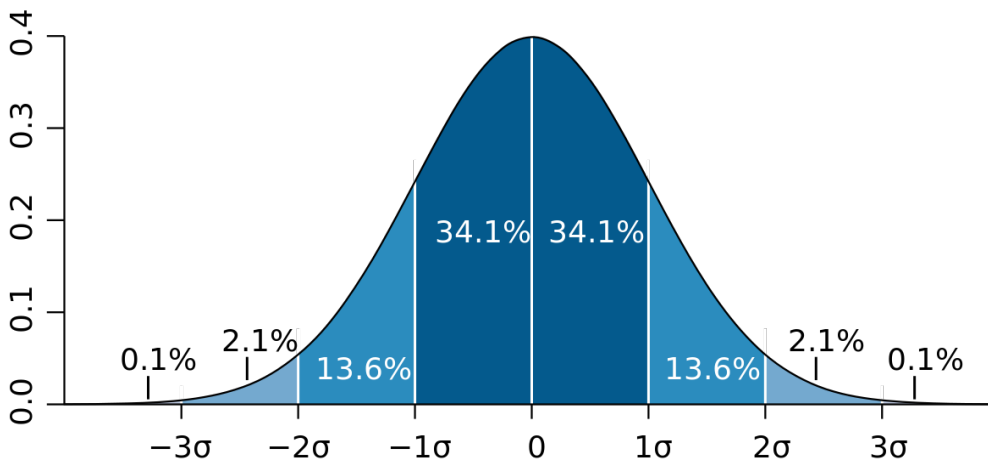
What is the difference between a Python list and a NumPy array?

Answer: NumPy arrays are **more efficient** and provide more **functionality** for numerical computations which are more general-purpose data structures.

What is broadcasting

Broadcasting is a feature in NumPy that allows arrays of different shapes to be used in an operation by automatically aligning their dimensions.

What is the difference between variance and standard deviation? Answer: Variance is a measure of how spread out a set of data is from its mean, calculated by taking the average of the squared differences from the mean, while standard deviation is the square root of the variance and provides a measure of spread in the same units as the original data.



Standard Deviation

$$\sqrt{\frac{\sum (x - \bar{x})^2}{n - 1}}$$

76	84	69	92	58
89	73	97	85	77

$$\bar{x} = \frac{\text{Sum}}{n}$$

...ing than Python lists,

...ithmetic operations by

...a measure of how
...differences from the
...ure of the spread of the

- Supervised Learning: Supervised learning is a machine learning algorithm that learns from labeled data, which means the data has input and output variables. The goal of supervised learning is to learn a model that can predict the output variables from the input variables based on the labeled data.
- Common Algorithms in Supervised Learning: The following algorithms are commonly used in supervised learning:
 - Linear Regression
 - Logistic Regression
 - Decision Trees
 - Random Forest
 - Support Vector Machines (SVM)
 - K-Nearest Neighbors (KNN)
 - Naive Bayes
 - Neural Networks
- Classification vs. Regression: Supervised learning can be divided into two main categories: classification and regression.
 - Classification involves predicting a categorical variable, such as predicting whether an email is spam or not.
 - Regression involves predicting a continuous variable, such as predicting house prices based on features like the number of bedrooms.
- Feature Selection: Feature selection is the process of selecting the most relevant features or variables from the original set of features.

a type of machine learning where the means it is provided with input-output pairs. The goal is to learn a mapping from input variables to the output variable using the training data. Some common algorithms used in

egorical output variable. For example, "spam" or "not spam". Continuous output variable. For example, "price" or "size". Features like size, location, and number of

rocess of selecting a subset of relevant features to use in the model. The objective

What is the purpose of data normalization?

11. Data large In Scikit-learn, what is the purpose of the predict method? a) To train the model using provided training data b) To make predictions on new data c) To evaluate the performance of the model d) To visualize the model's decision boundaries
Answer: b) To make predictions on new data

12. Which Scikit-learn function is commonly used to split a dataset into training and testing sets? a) train_test_split b) split_data c) test_train_split d) dataset_split
Answer: a) train_test_split

13. What is the recommended ratio for splitting a dataset into training and testing sets? a) 70:30 b) 80:20 c) 60:40 d) It depends on the size of the dataset and the problem domain
Answer: d) It depends on the size of the dataset and the problem domain

14. What is the purpose of setting a random seed when splitting a dataset into training and testing sets? a) To ensure reproducibility of the results b) To increase the model's accuracy c) To decrease the model's variance d) To speed up the training process
Answer: a) To ensure reproducibility of the results

15. What is the difference between training error and testing error? a) Training error measures the performance of the model on the training data, while testing error measures its performance on unseen data b) Training error measures the performance of the model on unseen data, while testing error measures its performance on the training data c) There is no difference between the two d) Both terms refer to the same concept
Answer: a) Training error measures the performance of the model on the training data, while testing error measures its performance on unseen data

16. What is overfitting in the context of machine learning models? a) When the model performs well on the training data but poorly on unseen data b) When the model performs poorly on both the training and testing data c) When the model has too few parameters to capture the underlying patterns in the data d) When the model converges too quickly during training
Answer: a) When the model performs well on the training data but poorly on unseen data

17. How can overfitting be mitigated in machine learning models? a) By increasing the complexity of the model b) By decreasing the size of the training dataset c) By regularizing the model using techniques like cross-validation d) By training the model for fewer epochs
Answer: c) By regularizing the model or using techniques like cross-validation

18. What is underfitting in the context of machine learning models? a) When the model performs poorly on both the training and testing data b) When the model performs well on both the training and testing data c) When the model has too many parameters d) When the model converges too slowly during training

ard variables with
ing the
ormance of

sting sets? a)

? a) 70:30 b)

g and testing
y c) To

measures the
formance
n data,
ference;
train the

data, while

performs
oorly on
apture the
aining
en data
complexity
model or

l performs

of feature selection is to improve the model
decreasing computational cost, and enhancing

- Feature Standardization: Feature standardization is the process of scaling or normalizing the range of features by transforming the features so that they have a mean of zero and a standard deviation of one. Feature standardization is needed because features with larger scales do not dominate those with smaller scales during training. It helps the model converge faster and improves the stability of the training process.
- Overfitting: Overfitting occurs when a model is too closely tailored to a specific dataset, capturing noise and random fluctuations in the data rather than the underlying pattern. This results in poor generalization performance on unseen data. Overfitting can be prevented by:
 - Using more training data
 - Cross-validation to evaluate model performance on unseen data
 - Feature selection to reduce model complexity
 - Regularization techniques such as L1 and L2
 - Early stopping during model training
- Underfitting: Underfitting occurs when a model is too simple to capture the underlying structure of the data. This results in poor performance on both training and testing data. Underfitting can be prevented by:
 - Increasing model complexity (e.g., adding more layers in a neural network)
 - Adding more features to the model
 - Reducing regularization constraints
 - Using a more complex model algorithm

's performance by reducing overfitting, improving model interpretability. Feature normalization, also known as feature scaling, is the process of adjusting the values of features in the dataset. It involves transforming the data to have a mean of zero and a standard deviation of one. This is useful because it ensures that features with larger scales do not dominate the training process. It helps algorithms converge faster and improves the performance of the models. Overfitting occurs when a model learns the training data too well, capturing noise rather than the underlying patterns. This leads to poor performance on new data. Overfitting can be prevented by:

- Reducing model complexity
- Using cross-validation
- Applying regularization (e.g., L1 and L2 regularization)

A model is too simple to capture the underlying patterns in the data, leading to poor performance both on the training data and on new data. This is known as underfitting. Underfitting can be prevented by:

- Increasing model complexity (e.g., adding more layers to a neural network)

n