Lecture 3

Classification and Prediction, overview
Review of Logistic Regression

**Classification and prediction** are essential data mining techniques used to make sense of data, discover patterns, and provide insights for various applications. Both techniques involve the analysis of historical data to build models that can be used to classify or predict future or unknown data points. Here's an overview of how classification and prediction are used in data mining:

**Classification**:
Classification is the process of categorizing data into predefined classes or categories based on the input attributes. It's a supervised learning technique, meaning it requires labeled training data, where each data point is associated with a known class or category.

***Steps in Classification***:
- *Data Collection*: Gather historical data that includes both the input attributes (features) and the corresponding class labels.
- *Data Preprocessing*: Clean and preprocess the data, handling missing values, outliers, and other data quality issues.
- *Feature Selection*: Identify and select the most relevant features for the classification task, which can help improve model accuracy and reduce computational complexity.
- *Training Data Split*: Divide the dataset into training and testing sets. The training set is used to build the classification model, while the testing set is used to evaluate the model's performance.
- *Model Building*: Select a classification algorithm (e.g., decision trees, support vector machines, or neural networks) and train the model using the training data.
- *Model Evaluation*: Assess the model's performance using the testing data. Common evaluation metrics include accuracy, precision, recall, F1 score, and the receiver operating characteristic (ROC) curve.
- *Model Deployment*: Once the model performs well on the testing data, deploy it to make predictions on new, unseen data.

***Applications of Classification***:
    *Email Spam Detection*: Classify emails as either spam or not spam.
    *Medical Diagnosis*: Categorize patient data into disease or non-disease groups.
    *Sentiment Analysis*: Determine the sentiment of text data (e.g., positive, negative, or neutral).
    *Credit Scoring*: Classify applicants as high-risk or low-risk for credit approval.

**Prediction**:

Prediction, also known as **regression**, is the process of estimating a continuous numerical value or outcome based on input attributes. It is used when the target variable is numeric and is not categorized into classes. Like classification, prediction is a supervised learning technique that requires labeled training data.

***Steps in Prediction***:
- *Data Collection*: Gather historical data with both the input attributes and the corresponding numerical target values.
- *Data Preprocessing*: Clean and preprocess the data, addressing issues such as missing values, outliers, and data quality.
- *Feature Selection*: Choose the most relevant input features for the prediction task.
- *Training Data Split*: Divide the dataset into training and testing sets for model evaluation.
- *Model Building*: Select a regression algorithm (e.g., linear regression, decision tree regression, or support vector regression) and train the model using the training data.
- *Model Evaluation*: Assess the model's performance using the testing data. Evaluation metrics for prediction tasks include mean squared error (MSE), root mean squared error (RMSE), and R-squared.
- *Model Deployment*: Deploy the model to make numerical predictions on new, unseen data.

***Applications of Prediction***:
> *Stock Price Forecasting*: Predict future stock prices based on historical data.
> *Demand Forecasting*: Estimate future demand for products or services.
> *Sales Revenue Prediction*: Forecast sales revenue for a business.
> *Weather Forecasting*: Predict temperature, precipitation, and other weather-related variables.

Classification and prediction are powerful tools in data mining and machine learning, enabling organizations to automate decision-making processes, improve accuracy in various domains, and gain valuable insights from data. These techniques are fundamental to building intelligent systems and making data-driven predictions and classifications.

We talked about many methods of regression (prediction) in 325. In this course, we will focus more on classification methods.

There are various **classification algorithms** used in data mining, and the choice of algorithm depends on the specific characteristics of the data and the problem you are trying to solve. Here are some of the most common classification algorithms in data mining:

***Logistic Regression*** is a simple and widely used algorithm for binary classification tasks. It models the relationship between the independent variables and the probability of a binary outcome.

***Decision Trees*** are used to partition the data into subsets based on the values of input features. They are interpretable and can handle both classification and regression tasks. Popular decision tree algorithms include CART and C4.5.

***Random Forest*** is an ensemble learning method that combines multiple decision trees to improve classification accuracy. It is robust and handles high-dimensional data well.

***Support Vector Machines (SVM)*** is a powerful classification algorithm that finds the optimal hyperplane to separate data into classes. It is effective in high-dimensional spaces and can handle non-linear data using kernel functions.

**K-Nearest Neighbors (K-NN)** is a simple and intuitive algorithm that classifies data points based on the majority class of their nearest neighbors. It's particularly useful for small to moderately sized datasets.

**Naive Bayes** is a probabilistic classification algorithm based on Bayes' theorem. It's often used for text classification, spam detection, and sentiment analysis.

**Neural Networks**, particularly deep learning models, are powerful for complex classification tasks. Convolutional Neural Networks (CNNs) are used for image classification, and Recurrent Neural Networks (RNNs) are applied to sequence data.

**K-Means Clustering**: While primarily a clustering algorithm, K-Means can be used for simple classification tasks by assigning new data points to the cluster with the closest centroid.

**Gradient Boosting Algorithms** like *XGBoost*, *LightGBM*, and *CatBoost* are widely used for classification tasks. They work well on structured data and are known for their high predictive accuracy.

**Ensemble Methods** combine the predictions of multiple base classifiers to improve accuracy and reduce overfitting. Besides random forests, other ensemble techniques include AdaBoost and Gradient Boosting.

**Nearest Centroid Classifier**: This simple algorithm classifies data points based on their similarity to the centroids of each class. It's suitable for text classification and image classification.

**Gaussian Naive Bayes**: A variation of the Naive Bayes algorithm that assumes a Gaussian distribution for continuous features. It's suitable for datasets with continuous attributes.

**LDA (Linear Discriminant Analysis)** is a dimensionality reduction technique that can also be used for classification tasks. It maximizes the separation between classes by projecting data into a lower-dimensional space.

The choice of classification algorithm depends on factors such as the nature of the data, the size of the dataset, the interpretability of the model, and the specific requirements of the problem at hand. It's common to experiment with multiple algorithms and fine-tune their parameters to achieve the best results for a given classification task.

Let's review Logistic Regression briefly.

**Logistic regression** is a statistical method used for binary classification, where the outcome variable is categorical and has two classes (usually coded as 0 and 1). It models the relationship between one or more independent variables and the probability of a particular outcome occurring.

**How Logistic Regression Works**:
*Sigmoid Function*: Logistic regression uses the logistic function (sigmoid function) to transform the linear combination of input features into a probability between 0 and 1. The sigmoid function is defined as:

$$P(Y = 1) = \frac{1}{1+e^{-(\beta_0+\beta_1 X_1+\beta_2 X_2+\ldots+\beta_n X_n)}}$$

where:

- $P(Y = 1)$ is the probability of the event Y being class 1.
- $e$ is the base of the natural logarithm.
- $\beta_0, \beta_1, \ldots, \beta_n$ are the coefficients to be estimated.
- $X_1, X_2, \ldots, X_n$ are the input features.

This equation is derived from a log-odds model that is a linear function. Solving for the probability results in the equation above.

*Training*: The logistic regression model is trained by adjusting the coefficients to maximize the likelihood of the observed data given the model.

*Decision Boundary*: The model predicts class 1 if the probability $P(Y = 1)$ is greater than or equal to 0.5 and predicts class 0 otherwise. The decision boundary is determined by the threshold (0.5 by default). It can be adjusted if needed.

**Pros**:
- *Interpretability*: Logistic regression coefficients represent the log-odds, making it easy to interpret the impact of each feature on the probability of the outcome.
- *Efficiency*: Logistic regression is computationally efficient and requires less computational resources compared to more complex algorithms.
- *Scalability*: It can handle a large number of features without much risk of overfitting.
- *Probabilistic Output*: Logistic regression provides probabilistic outputs, allowing for more nuanced predictions and decision-making.
- *No Assumptions about Feature Distribution*: Logistic regression does not assume that the input features are normally distributed.

**Cons**:
- *Linear Decision Boundary*: Logistic regression assumes a linear decision boundary, which may not capture complex relationships in the data.
- *Sensitivity to Outliers*: Logistic regression is sensitive to outliers, and their presence can impact the model's performance.
- *Assumption of Independence*: It assumes that the observations are independent of each other, which may not hold in certain situations.
- *Limited to Binary Classification*: Logistic regression is primarily designed for binary classification tasks and may need modifications for multi-class problems.
- *Requires Large Sample Sizes*: It performs better with a large number of samples, and the risk of overfitting increases with a small dataset.
- *No Feature Importance Ranking*: Unlike some tree-based models, logistic regression does not provide a natural feature importance ranking.

In summary, logistic regression is a versatile and interpretable classification algorithm suitable for various applications. However, its simplicity and linearity might be limiting in situations where the underlying

relationships are highly complex or nonlinear. It's crucial to consider the characteristics of the data and the problem at hand when choosing a classification algorithm.

In the next lecture, we'll look at two specific classification algorithms, Support Vector Machines (SVM) and K-Nearest Neighbors (KNN).

Resources:

1. https://www.geeksforgeeks.org/basic-concept-classification-data-mining/
2. https://www.upgrad.com/blog/classification-in-data-mining/
3. https://www.ibm.com/docs/en/db2/10.1.0?topic=algorithms-classification
4. https://stats.oarc.ucla.edu/r/dae/logit-regression/
5. https://www.datacamp.com/tutorial/logistic-regression-R
6. https://www.geeksforgeeks.org/logistic-regression-in-r-programming/
7. https://www.r-bloggers.com/2015/09/how-to-perform-a-logistic-regression-in-r/
8. https://uc-r.github.io/logistic_regression
9. https://towardsdatascience.com/how-to-do-logistic-regression-in-r-456e9cfec7cd
10. https://www.tutorialspoint.com/r/r_logistic_regression.htm
11. https://www.mastersindatascience.org/learning/machine-learning-algorithms/logistic-regression/
12. https://www.statology.org/types-of-logistic-regression/
13. https://www.spiceworks.com/tech/artificial-intelligence/articles/what-is-logistic-regression/
14. https://www.analyticsvidhya.com/blog/2021/08/conceptual-understanding-of-logistic-regression-for-data-science-beginners/