

– Task

To predict which passengers are transported to an alternate dimension, you'll need to follow a series of steps involving data preparation, model building, training, and evaluation. Here's a step-by-step guide on how to do this, assuming you have a dataset:

1. Gather Data

You need a dataset with features related to the passengers and a target variable indicating whether they were transported. If you don't have a dataset, you can create a synthetic one for practice.

2. Data Preparation

1. **Load Data:** Load dataset into a pandas DataFrame.
2. **Explore Data:** Examine the data to understand its structure, distributions, and any missing values.
3. **Clean Data:** Handle missing values, remove duplicates, and correct any inconsistencies.
4. **Feature Engineering:** Create or modify features to better capture the patterns. This might include normalizing numerical features, encoding categorical variables, etc.

3. Split Data

Split data into training and test sets. A common split is 80% for training and 20% for testing.

4. Model Building

1. **Choose a Model:** Select a machine learning model. For classification problems, common choices include Logistic Regression, Decision Trees and Random Forests and many more.
2. **Train the Model:** Train your model on the training data.
3. **Hyperparameter Tuning:** Use techniques like cross-validation to tune hyperparameters for better performance.

5. Evaluation

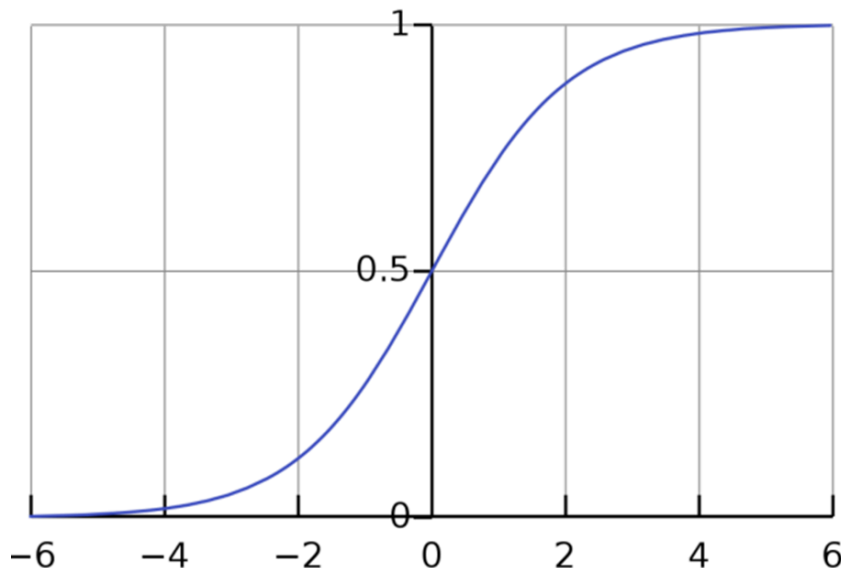
Evaluate the model's performance using the test set. Common metrics include accuracy, precision, recall, F1 score and confusion metrics

6. Prediction

Use the trained model to predict which passengers are transported.

Logistic Regression

Logistic Regression is a statistical method used for binary classification problems, where the outcome is a categorical variable with two possible values, often coded as 0 and 1. It estimates the probability that a given input point belongs to a particular class.



* we use "log loss" as cost function.

$$J(\theta_0, \theta_1) = \begin{cases} -\log(h_0(x)) & \text{if } y=1 \\ -\log(1-h_0(x)) & \text{if } y=0 \end{cases}$$
$$J(\theta_0, \theta_1) = -y \log(h_0(x)) - (1-y) \log(1-h_0(x))$$

here $h_0(x) = \frac{1}{1+e^{-z}} = \frac{1}{1+e^{-(\theta_0 + \theta_1 x_1)}}$

Key Points

- **Purpose:** To predict the probability of a binary outcome based on one or more predictor variables.
- **Algorithm Type:** Supervised learning, linear model.
- **Output:** Probability values between 0 and 1.
- **Decision Boundary:** Logistic regression uses a sigmoid function to model the relationship between the input variables and the probability of a specific outcome.

Applications

- **Medical Diagnosis:** Predicting the presence or absence of a disease.
- **Credit Scoring:** Determining the likelihood of a loan default.
- **Marketing:** Predicting whether a customer will buy a product.

Advantages

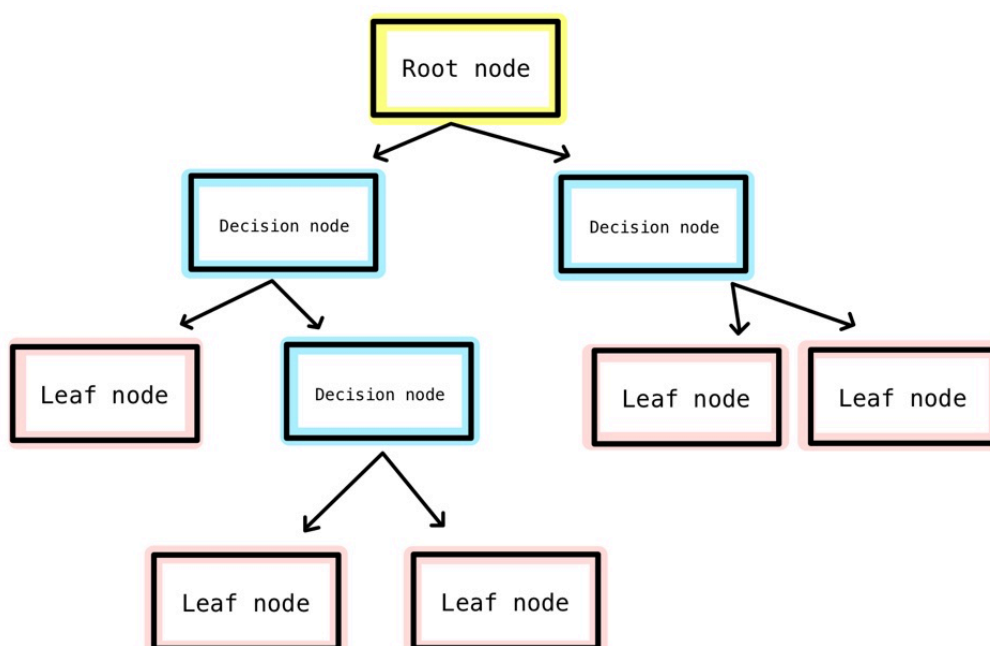
- **Simplicity:** Easy to implement and interpret.
- **Efficiency:** Works well with small to medium-sized datasets.
- **Probability Outputs:** Provides probabilistic interpretations.

Limitations

- **Binary Classification:** Primarily suited for binary outcomes (can be extended to multiclass problems with adjustments).
- **Linear Decision Boundary:** Assumes a linear relationship between the input variables and the log odds of the outcome. Non-linear relationships require transformations or more complex models.

Decision Tree Overview

Decision Trees are a type of supervised learning algorithm used for both classification and regression tasks. They work by recursively splitting the data into subsets based on the value of input features, leading to a tree-like model of decisions.



<p>a) Entropy</p> $H(G) = -P+ \log_2 P+ - P- \log_2 P-$ <p>$P+$ = probability of positive category</p> <p>$P-$ = probability of negative category</p>	<p>b) Gini Impurity</p> $G.I = 1 - \sum_{i=1}^n (p_i)^2$
---	--

Key Points

- **Purpose:** To predict the value of a target variable based on several input features by learning simple decision rules inferred from the data.
- **Algorithm Type:** Supervised learning, non-linear model.
- **Output:** Predicted class (classification) or continuous value (regression).
- **Structure:** Composed of nodes (decision points), branches (decision rules), and leaves (outcomes).

Applications

- **Medical Diagnosis:** Classifying diseases based on symptoms and medical test results.
- **Credit Scoring:** Predicting loan default risk based on borrower characteristics.
- **Customer Segmentation:** Categorizing customers into different groups based on purchasing behavior.

Advantages

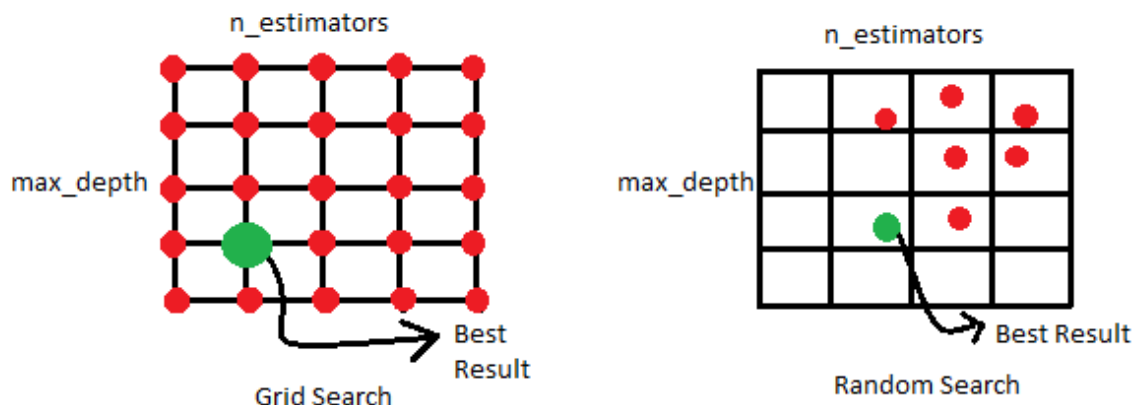
- **Interpretability:** Easy to understand and visualize.
- **No Assumption of Linear Relationship:** Can model complex, non-linear relationships.
- **Feature Importance:** Provides insights into feature importance.

Limitations

- **Overfitting:** Can easily overfit the training data, especially with deep trees.
- **Instability:** Small changes in the data can result in different splits and trees.
- **Bias:** Can be biased towards features with more levels (for categorical variables).

Hyperparameter Tuning with GridSearchCV

Hyperparameter tuning is the process of optimizing a model's performance by finding the best set of hyperparameters. `GridSearchCV` is a popular method in scikit-learn for systematically searching through a specified hyperparameter space to find the best combination for a given model.



Key Points

- **Purpose:** To improve model performance by finding the optimal set of hyperparameters.
- **Algorithm Type:** Model optimization, exhaustive search.
- **Output:** The best combination of hyperparameters and the best model trained with these hyperparameters.

Process

1. **Define Model:** Choose the machine learning model for which you want to tune hyperparameters.
2. **Specify Parameter Grid:** Create a dictionary where keys are the hyperparameter names and values are lists of values to try.
3. **Initialize GridSearchCV:** Create a `GridSearchCV` object with the model, parameter grid, scoring metric, cross-validation strategy, and other options.
4. **Fit the Model:** Fit the `GridSearchCV` object on the training data.
5. **Evaluate:** Access the best hyperparameters and evaluate the model on the test set if needed.

Performance Metrics Overview

Performance metrics are essential for evaluating the effectiveness of a machine learning model. They provide insights into how well a model performs and help in comparing different models or approaches. Here are key performance metrics commonly used for classification tasks:

Key Metrics

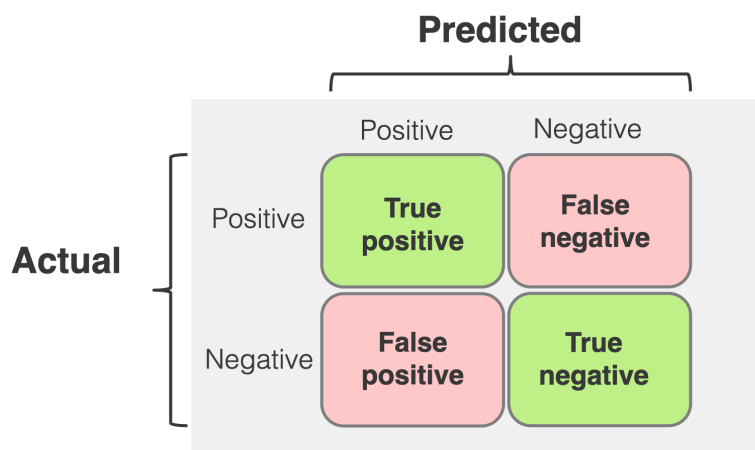
1. **Accuracy Score**
2. **Confusion Matrix**
3. **Precision**
4. **Recall**

Accuracy Score

Definition: The ratio of correctly predicted instances to the total instances.

Confusion Matrix

Definition: A table used to describe the performance of a classification model, showing the actual versus predicted classifications.



Structure:

- **True Positive (TP):** Correct positive predictions.
- **True Negative (TN):** Correct negative predictions.
- **False Positive (FP):** Incorrect positive predictions.
- **False Negative (FN):** Incorrect negative predictions.

Precision

Definition: The ratio of correctly predicted positive observations to the total predicted positives.

Recall

Definition: The ratio of correctly predicted positive observations to all observations in the actual class.