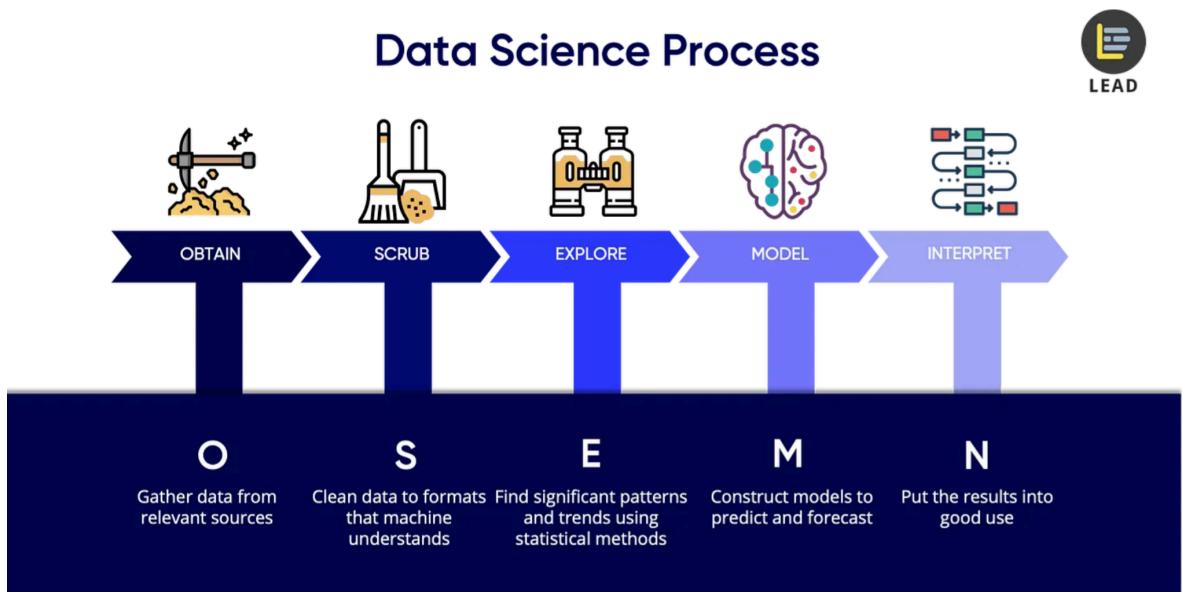


# Data Science Project Stages



Let's dive deeper into each step of the data science process, including when to use different machine learning models:

## 1. Collecting Data:

- Gather data from diverse sources such as databases, APIs, files, or web scraping.
- Choose data that aligns with project objectives and ensure it's of high quality and relevance.

## 2. Cleaning Data:

- Handle missing values, outliers, duplicates, and inconsistencies in the dataset.
- Impute missing values using techniques like mean, median, mode, or advanced methods like MICE.
- Detect and correct data entry errors or inconsistencies.
- Standardize or normalize numerical features to a common scale to prevent biases.

## 3. Exploratory Data Analysis (EDA):

- Explore data distribution, correlations, and relationships between variables.
- Visualize data using histograms, scatter plots, box plots, and correlation matrices.
- Identify patterns, trends, and anomalies that inform further analysis.

# DATA CLEANING CHECKLIST

## Up-to-date data



Data should be up-to-date in order to obtain maximum value from the data analysis.

## Missing values



Count missing values and analyze where in the data they are missing. Missing values can disrupt some analyses and skew the results.

## Duplicates



Duplicate IDs indicate multiple records for one person, e.g. someone holds multiple functions at the same time.

## Numerical outliers



Numerical outliers are fairly easy to detect and remove. Define minimum and maximum to spot outliers easily.

## Check IDs



Check data labels of all the fields to see whether some categorical values are mislabeled.

## Define valid output



Define valid data labels for categorical data. Define data ranges for numerical variables. Non-matching data is presumably wrong.

Data encoding is an essential step in data preprocessing, especially when dealing with categorical variables. Here's the detailed information about data encoding:

## Data Encoding:

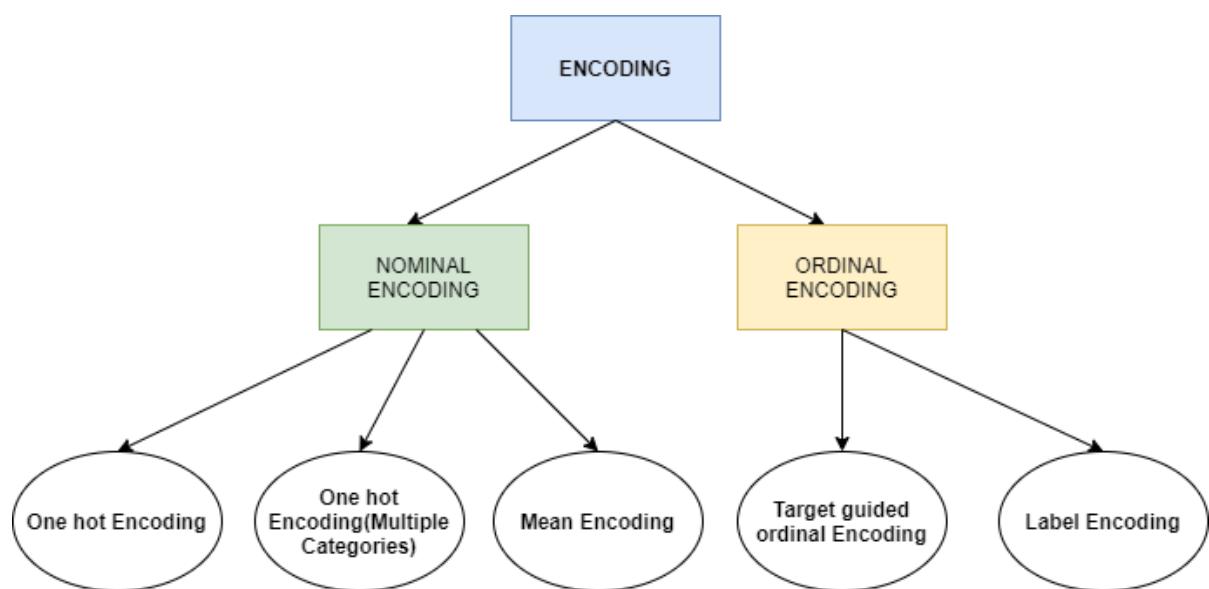
- **Categorical Variables:** Categorical variables are non-numeric variables that represent categories or groups.
- **Machine Learning Models:** Most machine learning algorithms require numeric input data, so categorical variables need to be converted into a numerical format before feeding them into models.

**Original categorical column**

Origin
USA
Japan
Europe
USA
Europe

**One-Hot encoded columns**

Origin_USA	Origin_Japan	Origin_Europe
1	0	0
0	1	0
0	0	1
1	0	0
0	0	1

**Techniques for Encoding Categorical Variables:****1. One-Hot Encoding:**

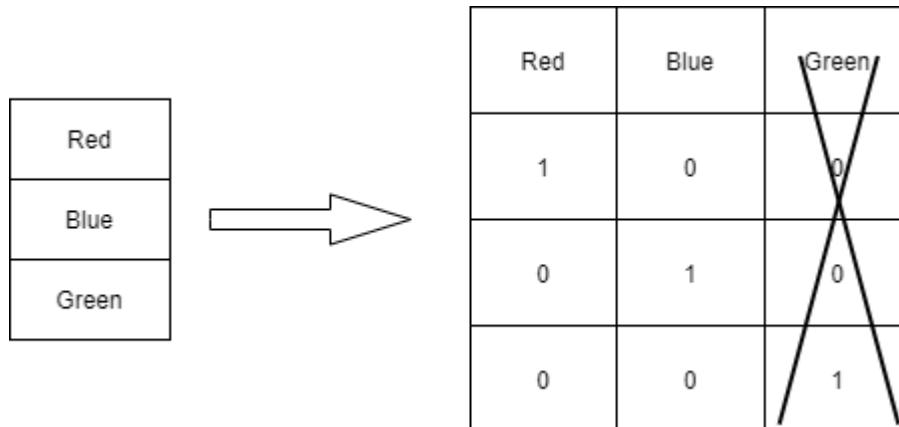
- Converts categorical variables into a binary format where each category becomes a new binary feature.
- Creates a new binary column for each category, with a value of 1 indicating the presence of the category and 0 otherwise.
- Suitable for nominal categorical variables (categories with no inherent order).
- Avoids assigning ordinal relationships between categories.
- Example: Encoding "Gender" with categories "Male", "Female", "Other" would create three binary columns.

**2. Label Encoding:**

- Assigns a unique integer to each category in the variable.
- Encodes ordinal categorical variables (categories with a natural ordering).
- Not suitable for non-ordinal categorical variables as it may introduce unintended ordinal relationships.
- Example: Encoding "Size" with categories "Small", "Medium", "Large" would assign integers 0, 1, and 2 respectively.

**3. Ordinal Encoding:**

- Similar to label encoding but assigns integers based on the order of categories.
- Preserves the ordinal relationship between categories.
- Suitable for ordinal categorical variables.
- Example: Encoding "Education Level" with categories "High School", "Bachelor's Degree", "Master's Degree" would assign integers based on the educational attainment level.



### When to Apply Data Encoding:

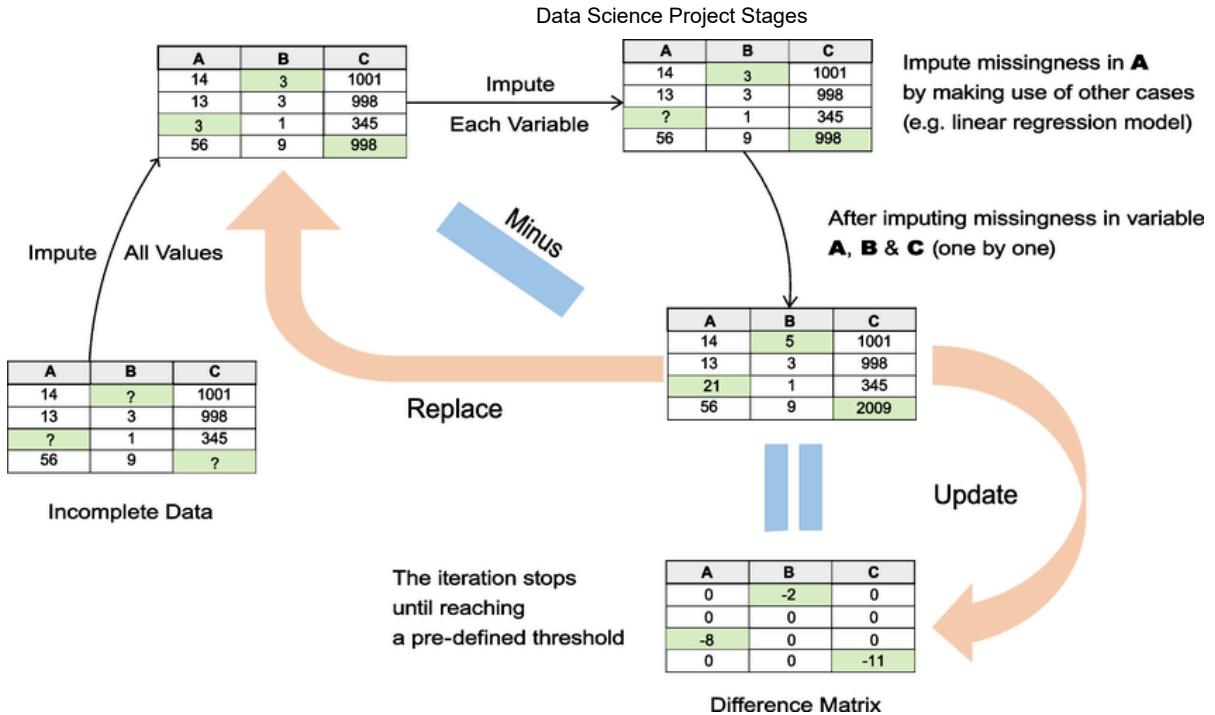
- **One-Hot Encoding:** Use one-hot encoding for nominal categorical variables (categories with no inherent order) to avoid introducing unintended ordinal relationships.
- **Label Encoding:** Apply label encoding for ordinal categorical variables (categories with a natural ordering) when the order of categories matters.
- **Ordinal Encoding:** Use ordinal encoding for ordinal categorical variables to preserve the ordinal relationships between categories.

### When Not to Apply Data Encoding:

- **One-Hot Encoding:** Avoid using one-hot encoding for ordinal categorical variables as it doesn't capture the ordinal relationships between categories.
- **Label Encoding:** Be cautious when applying label encoding to non-ordinal categorical variables as it may introduce unintended ordinal relationships between categories.
- **Ordinal Encoding:** Avoid using ordinal encoding for nominal categorical variables where there is no natural order among categories.

By appropriately encoding categorical variables, you ensure that the machine learning models can effectively interpret and utilize categorical information in the dataset, leading to more accurate and reliable predictions.

## MICE & SMOTE



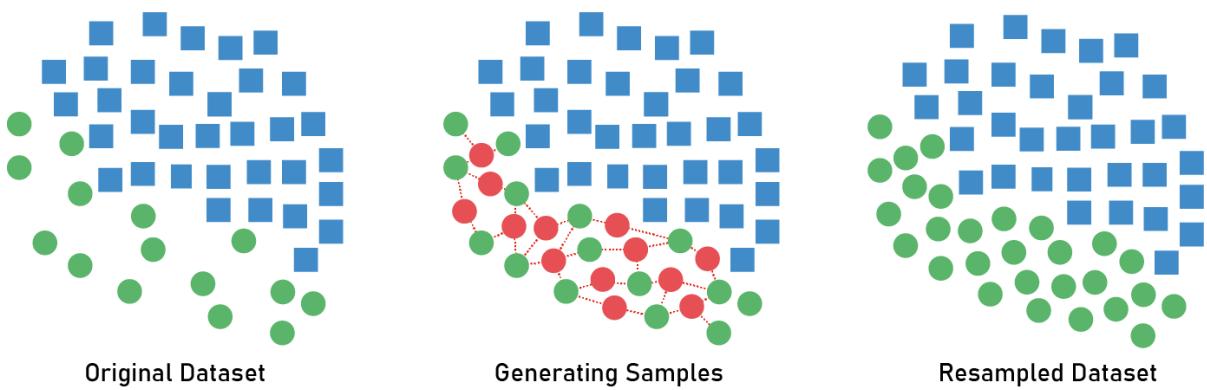
### 1. Imputing Missing Values (MICE):

- Apply MICE (Multiple Imputation by Chained Equations) to handle missing values, especially when missingness is not completely random.
- Use MICE for datasets with complex missing data patterns or when preserving relationships between variables is crucial.

### 2. Handling Imbalanced Data (SMOTE):

- Apply SMOTE (Synthetic Minority Over-sampling Technique) to address class imbalance in classification tasks.
- Use SMOTE when the minority class is underrepresented and improving classification performance is essential.

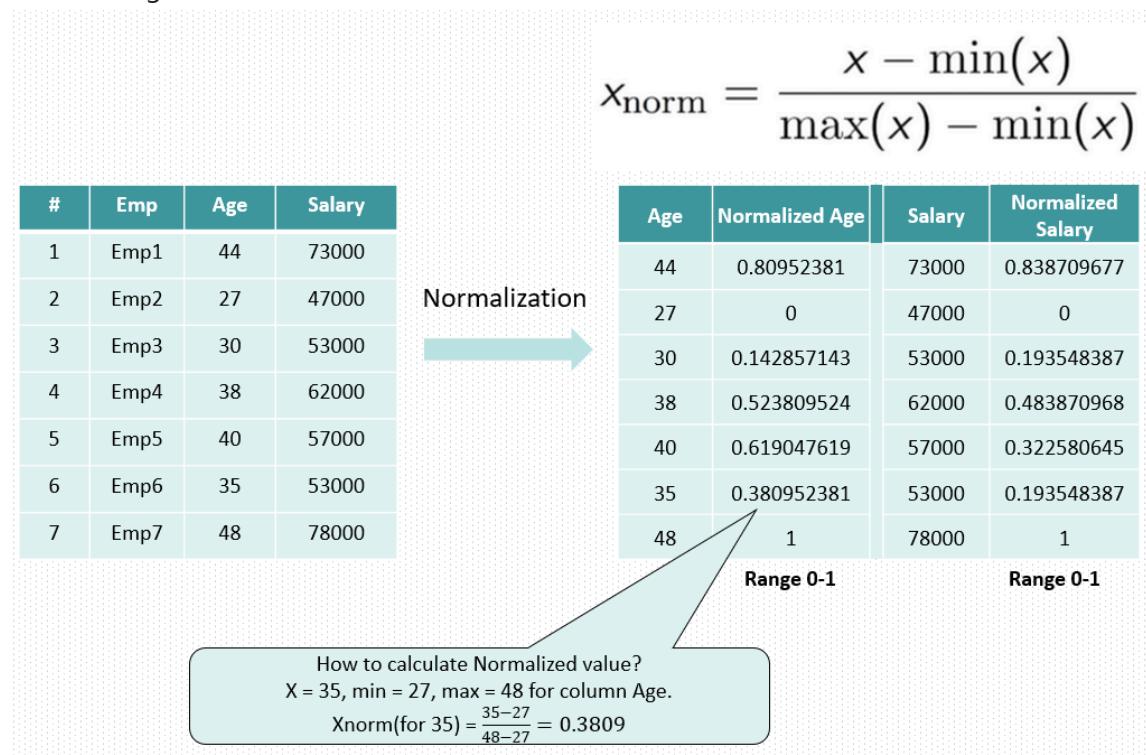
## Synthetic Minority Oversampling Technique



### 1. Scaling Data:

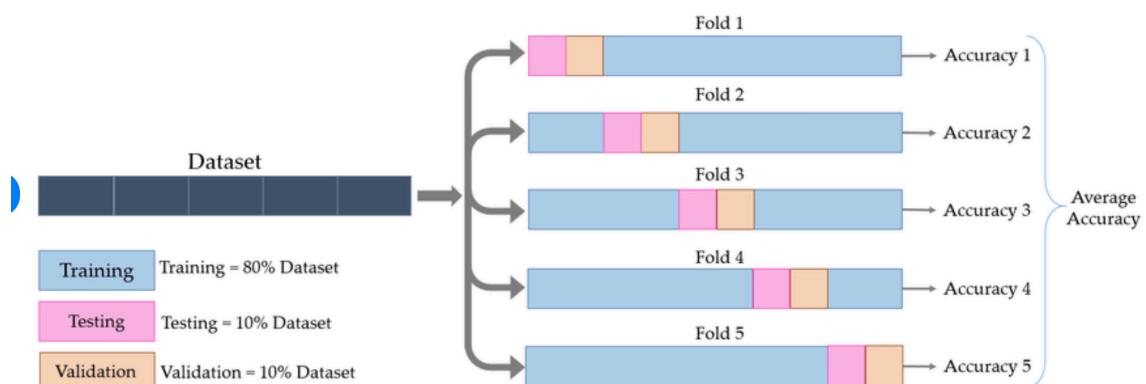
- Scale numerical features to a common range using techniques like standardization (subtract mean, divide by standard deviation) or normalization (scale features to a range between 0 and 1).

- Scaling helps prevent features with larger scales from dominating the model during training.



## 1. Train-Test Split:

- Split the dataset into training and testing sets using a predefined ratio (e.g., 70-30, 80-20).
- Ensure that the split maintains the distribution of target classes in classification tasks.
- Avoid data leakage by ensuring that the testing set is not used during model training.



Visual representation of the training, test, and validation split using cross-validation.

## Train-Test Split Logic:

### 1. Logic:

- The primary goal of machine learning models is to generalize well to unseen data. Therefore, we need to evaluate the model's performance on data it hasn't seen during training.
- The train-test split allows us to simulate real-world scenarios where the model encounters new, unseen data.

## 2. Train Set:

- The training set is used to train the model's parameters. The model learns patterns, relationships, and features present in the training data.
- The larger the training set, the better the model can learn from the available data.

## 3. Test Set:

- The test set is held out from the model during training and is used to evaluate its performance after training.
- By evaluating the model on unseen data, we get an estimate of its generalization ability and how well it performs on new, unseen examples.

## 4. Validation Set:

- The validation set is an additional subset of the data used during the model training process, primarily for hyperparameter tuning and model selection.
- It helps prevent overfitting by providing an independent dataset for model evaluation and parameter optimization.
- The validation set's performance metrics guide the selection of the best-performing model architecture and hyperparameters.

# Real-Life Data Science Example:

**Scenario:** Predicting House Prices

## Explanation:

- **Data Collection:** Suppose we have a dataset containing information about houses, including features like size, number of bedrooms, location, and price.
- **Train-Test Split:** We divide the dataset into a training set (e.g., 70% of the data) and a test set (e.g., 30% of the data).
- **Training:** We use the training set to train our machine learning model to predict house prices based on the features provided.
- **Validation (Optional):** We might use a validation set for hyperparameter tuning and model selection, iterating over different model architectures, and hyperparameters to find the optimal combination.
- **Testing:** Once the model is trained, we evaluate its performance on the test set, which contains unseen data. We calculate metrics such as mean squared error (MSE) or root mean squared error (RMSE) to assess how well the model generalizes to new, unseen houses.
- **Generalization:** If the model performs well on the test set, we can have confidence that it will likely perform well on new houses that it hasn't seen during training.

## Conclusion:

The train-test split is a critical technique in machine learning and data science for evaluating model performance and ensuring that models generalize well to new, unseen data. By using separate training, validation, and testing sets, we can develop and validate robust models that accurately predict outcomes in real-world scenarios.

## Cross Validation:

Cross-validation is a resampling technique used to assess the performance of a predictive model. It helps estimate how the model will generalize to an independent dataset by iteratively splitting the available data into training and validation sets.

## Types of Cross Validation:

### 1. K-Fold Cross Validation:

- **Procedure:** The dataset is divided into k folds of approximately equal size. The model is trained k times, each time using k-1 folds for training and the remaining fold for validation.
- **Advantages:**
  - Provides a more robust estimate of model performance compared to a single train-test split.
  - Utilizes the entire dataset for training and validation, reducing bias.
- **Disadvantages:**
  - Computationally expensive, especially for large datasets or complex models.
  - May result in high variance depending on the choice of k.

### 2. Stratified K-Fold Cross Validation:

- **Procedure:** Similar to k-fold cross-validation, but ensures that each fold contains approximately the same proportion of target classes as the original dataset.
- **Advantages:**
  - Suitable for imbalanced datasets where the distribution of target classes is uneven.
  - Provides more reliable estimates of model performance for classification tasks.
- **Disadvantages:**
  - May be less effective for regression tasks or when the target variable is continuous.

### 3. Leave-One-Out Cross Validation (LOOCV):

- **Procedure:** Each observation in the dataset is used as the validation set once, while the rest of the data is used for training.
- **Advantages:**
  - Provides a nearly unbiased estimate of model performance.
  - Suitable for small datasets or when computational resources are limited.
- **Disadvantages:**
  - Can be computationally expensive, especially for large datasets.
  - Variance of the estimate may be high, leading to instability in results.

### 4. Leave-P-Out Cross Validation:

- **Procedure:** Similar to LOOCV but leaves p observations out for validation, rather than just one.
- **Advantages:**
  - Less computationally intensive compared to LOOCV, while still providing a reliable estimate of model performance.
  - Allows for fine-tuning the trade-off between computational cost and variance in the estimate.
- **Disadvantages:**
  - Still requires considerable computational resources, especially for large p values.
  - May not be suitable for highly imbalanced datasets or when p is chosen arbitrarily.

## Advantages and Disadvantages of Cross Validation:

- **Advantages:**
  - Provides a more accurate estimate of model performance compared to a single train-test split.
  - Reduces the risk of overfitting by evaluating the model on multiple validation sets.
  - Utilizes the entire dataset for training and validation, maximizing data efficiency.
- **Disadvantages:**
  - Can be computationally expensive, especially for large datasets or complex models.
  - May result in high variance, particularly with small sample sizes or when the dataset is highly imbalanced.
  - Requires careful consideration of the appropriate cross-validation strategy based on the characteristics of the dataset and the modeling task.

In summary, cross-validation is a valuable technique for assessing the performance of predictive models, but it's essential to choose the appropriate type of cross-validation based on the specific requirements and constraints of the problem at hand. Each type of cross-validation has its advantages and disadvantages, and understanding these can help researchers and practitioners make informed decisions when evaluating model performance.

## Hyperparameter Tuning:

Hyperparameter tuning is the process of selecting the optimal hyperparameters for a machine learning model to maximize its performance on unseen data. Hyperparameters are parameters that are set before the learning process begins and control aspects of the learning process itself.

### Importance of Hyperparameter Tuning:

#### 1. Model Performance:

- Well-tuned hyperparameters can significantly improve a model's performance, leading to better accuracy, precision, and recall on unseen data.
- Hyperparameter tuning helps in achieving the best possible generalization performance of the model.

#### 2. Generalization:

- Properly tuned hyperparameters prevent overfitting or underfitting by finding the right balance between model complexity and generalization.
- They ensure that the model captures the underlying patterns in the data without memorizing noise.

#### 3. Efficiency:

- Hyperparameter tuning optimizes the learning process, reducing the time and computational resources required to train the model.
- It helps in identifying the most promising hyperparameter configurations early in the experimentation process.

## Techniques for Hyperparameter Tuning:

## 1. Manual Search:

- Manually adjusting hyperparameters based on domain knowledge and intuition.
- Simple but time-consuming and may not always lead to optimal results.

## 2. Grid Search:

- Exhaustively searches through a specified subset of hyperparameter combinations.
- Evaluates model performance for each combination using cross-validation and selects the best-performing one.
- Computationally expensive, especially for a large search space.

## 3. Random Search:

- Randomly samples hyperparameter combinations from a predefined search space.
- Allows for a more efficient exploration of the hyperparameter space compared to grid search.
- Often yields comparable or better results with fewer computational resources.

## 4. Bayesian Optimization:

- Uses probabilistic models to learn the relationship between hyperparameters and model performance.
- Iteratively explores the hyperparameter space by selecting promising points based on past evaluations.
- Efficiently balances exploration and exploitation to converge to the optimal hyperparameters.

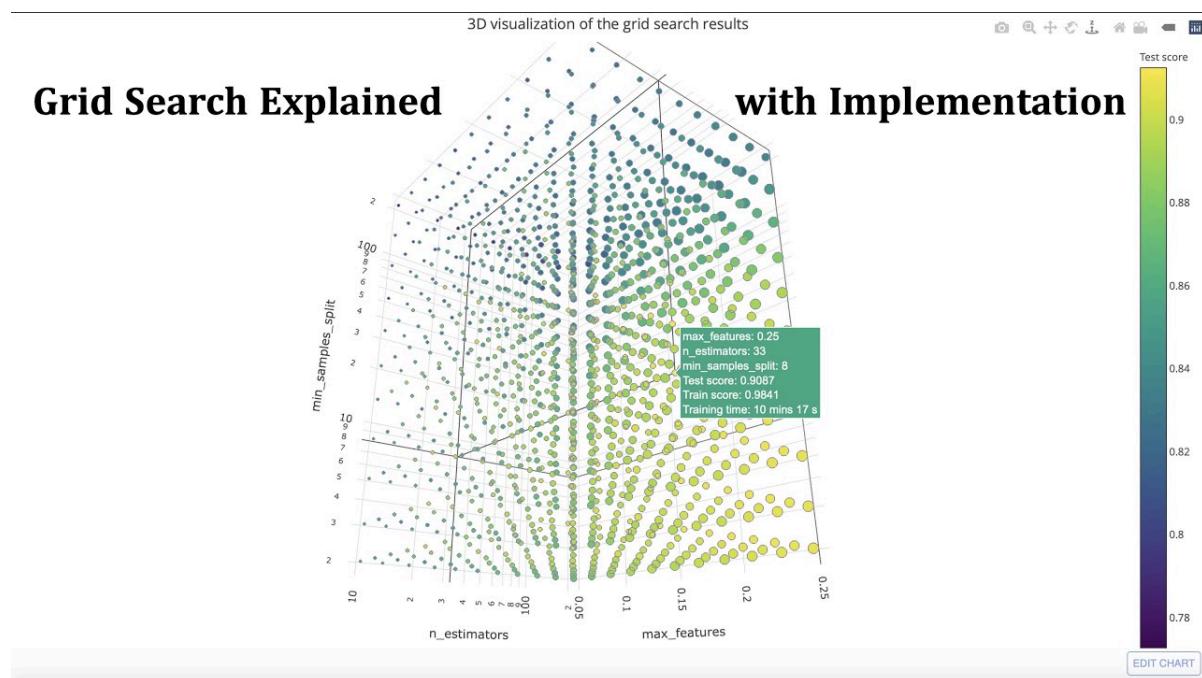
## 5. Automated Hyperparameter Tuning Libraries:

- Libraries like scikit-learn's `GridSearchCV`, `RandomizedSearchCV`, and `Optuna` provide built-in functions for hyperparameter tuning.
- These libraries streamline the process of hyperparameter optimization and facilitate experimentation with different algorithms and hyperparameters.

## Best Practices for Hyperparameter Tuning:

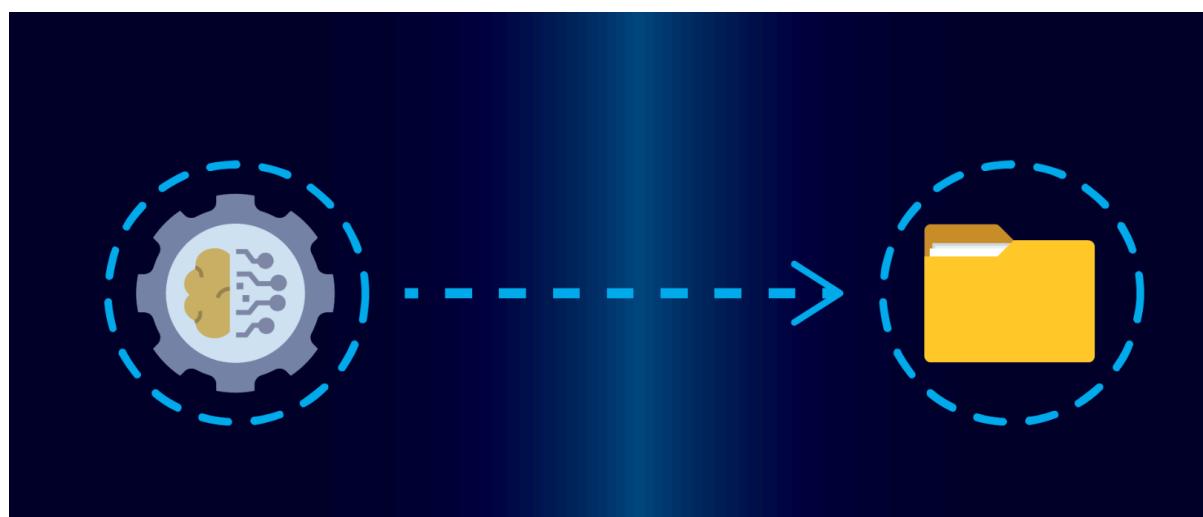
1. **Define Evaluation Metrics:** Choose appropriate evaluation metrics (e.g., accuracy, precision, recall, F1 score) based on the problem domain and objectives.
2. **Use Cross-Validation:** Employ cross-validation to robustly estimate the model's performance and prevent overfitting during hyperparameter tuning.
3. **Start with Coarse Search:** Begin with a coarse search over a wide range of hyperparameters to identify promising regions of the search space.
4. **Refine with Fine Search:** Narrow down the search space and perform a finer search around the promising hyperparameter combinations identified during the coarse search.
5. **Consider Computational Resources:** Take into account the available computational resources and time constraints when choosing the hyperparameter tuning strategy.
6. **Experiment and Iterate:** Hyperparameter tuning is an iterative process. Experiment with different hyperparameter configurations and iterate based on the performance feedback.

Hyperparameter tuning is a crucial step in the model development pipeline, allowing machine learning practitioners to optimize model performance and achieve better generalization on unseen data. By leveraging appropriate tuning techniques and best practices, practitioners can unlock the full potential of their machine learning models.



## 1. Model Building:

- Choose appropriate machine learning algorithms based on the nature of the problem and data characteristics:
  - Linear Models: Logistic Regression (classification), Linear Regression (regression)
  - Tree-Based Models: Decision Trees, Random Forests, Gradient Boosting Machines
  - Support Vector Machines (SVM)
  - Neural Networks: Multi-layer Perceptron (MLP)
  - K-Nearest Neighbors (KNN)
  - Naive Bayes
- Consider the interpretability, scalability, and complexity of the model when selecting algorithms.



## 1. Model Evaluation:

- Evaluate model performance using appropriate metrics:

- Classification: Accuracy, Precision, Recall, F1 Score, ROC Curve, AUC Score
- Regression: Mean Absolute Error (MAE), Mean Squared Error (MSE), R-squared
- Choose the best-performing model based on evaluation results and consider trade-offs between different metrics.

### 1. Pickling:

- Serialize the trained model using pickle or joblib to save it as a binary file.
- Pickling allows you to save the model's state, including parameters and trained weights, for later use or deployment.

### 2. Deployment:

- Deploy the trained model in a production environment, such as a web application or API, for real-time predictions.
- Monitor and update the deployed model regularly to maintain performance and accuracy.

When to Use Each ML Model:

- **Linear Models:** Suitable for problems with linear relationships between features and target variables. Often used for regression and binary classification tasks.
- **Tree-Based Models:** Effective for handling non-linear relationships and capturing complex patterns in the data. Useful for classification and regression tasks.
- **Support Vector Machines (SVM):** Ideal for binary classification problems with complex decision boundaries. SVMs work well with high-dimensional data and can handle non-linear relationships using kernel tricks.
- **Neural Networks (MLP):** Powerful for modeling complex and non-linear relationships in large datasets. Particularly effective for image recognition, natural language processing, and sequence prediction tasks.
- **K-Nearest Neighbors (KNN):** Simple and intuitive for classification and regression tasks. KNN is effective when there's sufficient labeled data and the decision boundary is not complex.
- **Naive Bayes:** Suitable for text classification and spam filtering tasks. Naive Bayes assumes independence between features, making it fast and efficient for large datasets with many features.

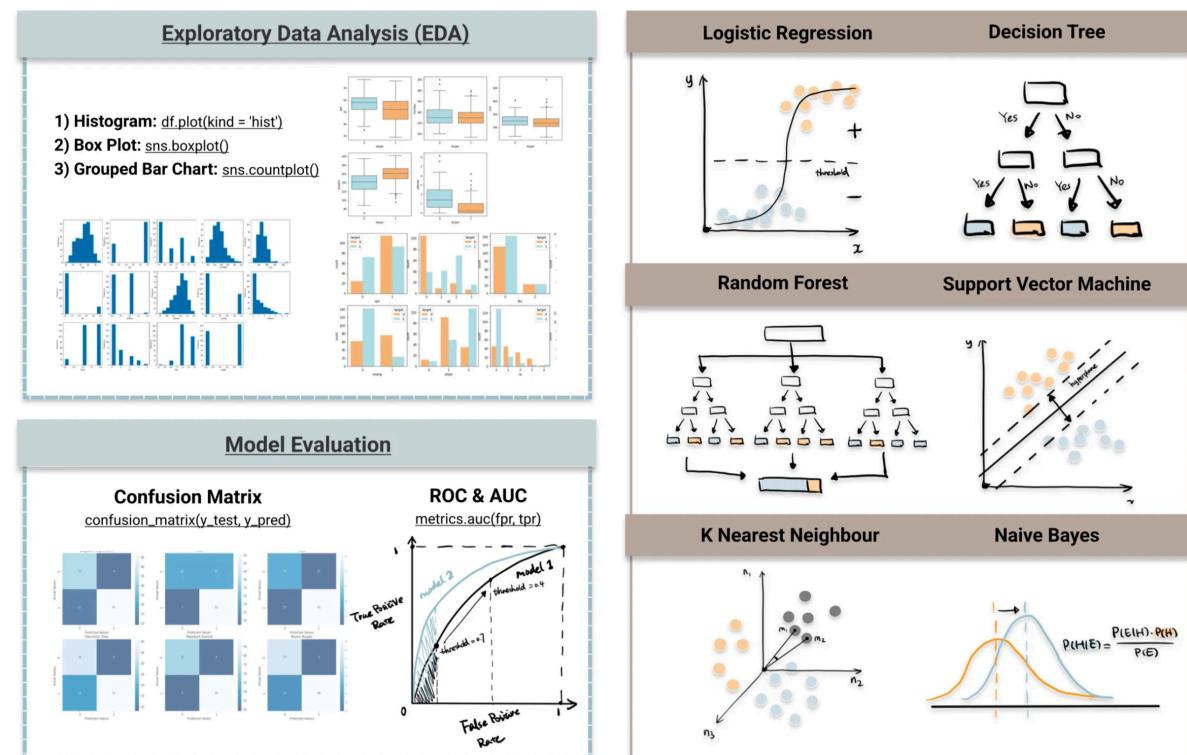
By considering the characteristics of each machine learning model and the requirements of the problem at hand, you can select the most appropriate algorithms to achieve optimal performance in your data science projects.

Model Type	Use Case	Characteristics	When to Choose
Linear Regression	Predicting continuous target variables	Assumes linear relationship between features and target	- When the relationship between features and target is approximately linear - Interpretability is important
Logistic Regression	Binary classification	Estimates probabilities of class membership	- Binary classification tasks - Linear decision boundaries
Decision Trees	Classification and regression tasks	Non-parametric, can capture complex relationships	- Non-linear relationships between features and target

Model Type	Use Case	Characteristics	When to Choose
Random Forests	Classification and regression tasks	Ensemble of decision trees, reduces overfitting	- Interpretability is not a priority - High-dimensional datasets - Avoid overfitting
Support Vector Machines (SVM)	Classification and regression tasks	Find hyperplanes with maximum margin between classes	- Non-linear relationships - Small to medium-sized datasets
K-Means Clustering	Clustering tasks	Divides data into k clusters based on similarity	- Unlabeled data - Exploration of data patterns
Principal Component Analysis (PCA)	Dimensionality reduction	Identifies orthogonal axes capturing maximum variance	- High-dimensional datasets - Visualization and feature selection
Neural Networks	Classification, regression, and complex pattern recognition	Can learn complex relationships in data	- Large datasets with complex patterns - Deep learning tasks

The choice of the model depends on several factors such as the nature of the data, the problem at hand, the interpretability of the model, the size of the dataset, and the computational resources available. Understanding the characteristics and capabilities of each model helps in selecting the most appropriate one for a given task.

## Machine Learning Algorithms - Classification



## 1. Supervised Learning Models:

### 1. Linear Regression:

- Used for predicting continuous target variables based on one or more input features.
- Assumes a linear relationship between input and output variables.

### 2. Logistic Regression:

- Used for binary classification tasks, where the target variable has two possible outcomes.
- Estimates the probability that an observation belongs to a particular class.

### 3. Decision Trees:

- Non-parametric supervised learning models used for classification and regression tasks.
- Divides the input space into regions and predicts the target variable based on the majority class or average value in each region.

### 4. Random Forests:

- Ensemble learning method that constructs multiple decision trees during training and outputs the mode or mean prediction of the individual trees.
- Provides improved performance and reduced overfitting compared to single decision trees.

### 5. Support Vector Machines (SVM):

- Supervised learning models used for classification and regression tasks.
- Find the hyperplane that best separates classes in the feature space, maximizing the margin between classes.

## 2. Unsupervised Learning Models:

### 1. K-Means Clustering:

- Unsupervised learning algorithm used for clustering tasks.
- Divides the dataset into k clusters based on similarity between data points.

### 2. Hierarchical Clustering:

- Clustering algorithm that builds a hierarchy of clusters by recursively merging or splitting clusters.
- Does not require the number of clusters to be specified in advance.

### 3. Principal Component Analysis (PCA):

- Dimensionality reduction technique used to transform high-dimensional data into a lower-dimensional space while preserving most of the variance.
- Finds orthogonal axes (principal components) that capture the maximum variance in the data.

## 3. Semi-Supervised Learning Models:

### 1. Self-Training:

- Iterative algorithm that starts with a small labeled dataset and progressively incorporates unlabeled data by training a model on the labeled data and using it to label unlabeled instances with high confidence.

### 2. Co-Training:

- Ensemble-based approach where multiple models are trained on different subsets of features or labeled data, and they exchange information to improve performance.

## 4. Reinforcement Learning Models:

### 1. Q-Learning:

- Model-free reinforcement learning algorithm used to learn the optimal action-selection policy for an agent in a Markov decision process.
- Learns the quality (Q-value) of taking a particular action in a particular state.

### 2. Deep Q-Networks (DQN):

- Deep learning model that combines deep neural networks with Q-learning to handle high-dimensional input spaces and complex environments.

## 5. Neural Network Models:

### 1. Feedforward Neural Networks (FNN):

- Simplest form of neural networks where information flows in one direction, from input to output layers.
- Used for various tasks like classification, regression, and pattern recognition.

### 2. Convolutional Neural Networks (CNN):

- Specifically designed for processing grid-like data, such as images.
- Consist of convolutional layers that apply filters to extract spatial hierarchies of features.

### 3. Recurrent Neural Networks (RNN):

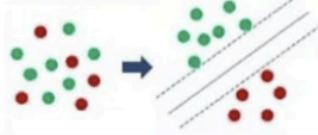
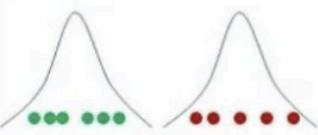
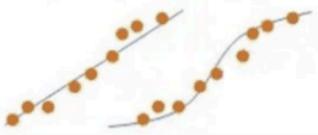
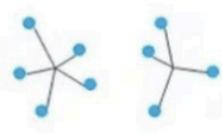
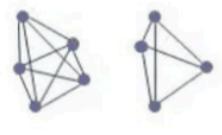
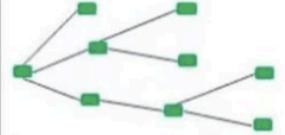
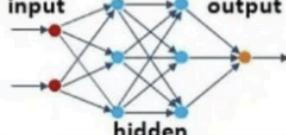
- Designed to handle sequential data with temporal dependencies.
- Contains recurrent connections that allow information to persist over time.

### 4. Long Short-Term Memory (LSTM) Networks:

- A variant of RNNs with memory cells that can maintain information over long sequences.
- Suitable for tasks like speech recognition, language modeling, and time series prediction.

Each type of machine learning model has its strengths, weaknesses, and specific use cases.

Understanding the characteristics and capabilities of different models is essential for selecting the most appropriate model for a given task.

Algorithm	Keyword	Diagram
<b>Support Vector Machines (SVM)</b>	<b>Vector on Points</b>	
<b>Naïve Bayes</b>	<b>Probability Distribution</b>	
<b>Linear Regression Logistic Regression</b>	<b>Straight Line Logarithmic Line</b>	
<b>K-Means</b>	<b>Kernel (<i>central</i>) Mean</b>	
<b>K-Nearest Neighbour</b>	<b>Neighbouring Points</b>	
<b>Decision Trees</b>	<b>Tree Branches</b>	
<b>Neural Networks</b>	<b>Network with Layers of elements</b>	

## Neural Networks:

Neural networks are a class of machine learning models inspired by the structure and function of the human brain. They consist of interconnected layers of artificial neurons (nodes) that process input data to produce output predictions. Neural networks have gained popularity due to their ability to learn complex patterns and relationships in data, making them suitable for various tasks, including classification, regression, and pattern recognition.

### Key Components of Neural Networks:

#### 1. Neurons (Nodes):

- Neurons are the basic computational units of a neural network.
- Each neuron receives input signals, performs computations using activation functions, and produces an output signal.

#### 2. Layers:

- Neural networks consist of multiple layers of interconnected neurons organized into input, hidden, and output layers.

- The input layer receives the initial input data, while the output layer produces the final predictions.
- Hidden layers perform intermediate computations and extract features from the input data.

### 3. Weights and Biases:

- Each connection between neurons is associated with a weight parameter that represents the strength of the connection.
- Biases are additional parameters added to neurons, influencing the neuron's activation threshold and introducing flexibility into the model.

### 4. Activation Functions:

- Activation functions introduce non-linearities into the neural network, allowing it to learn complex relationships in the data.
- Common activation functions include sigmoid, tanh, ReLU (Rectified Linear Unit), and softmax.

### 5. Forward Propagation:

- During forward propagation, input data is fed into the neural network, and computations are performed layer by layer to produce output predictions.
- Each layer's output serves as the input to the subsequent layer until the final output is generated.

### 6. Backward Propagation (Training):

- Backward propagation is the process of updating the weights and biases of the neural network to minimize the difference between predicted and actual outputs (i.e., the loss function).
- It involves calculating gradients of the loss function with respect to the model parameters and adjusting the parameters using optimization algorithms like gradient descent.

## Types of Neural Networks:

### 1. Feedforward Neural Networks (FNN):

- The simplest form of neural networks where information flows in one direction, from input to output layers.
- Widely used for tasks like classification and regression.

### 2. Convolutional Neural Networks (CNN):

- Specifically designed for processing grid-like data, such as images.
- Consist of convolutional layers that apply filters to extract spatial hierarchies of features.

### 3. Recurrent Neural Networks (RNN):

- Designed to handle sequential data with temporal dependencies.
- Contains recurrent connections that allow information to persist over time.

### 4. Long Short-Term Memory (LSTM) Networks:

- A variant of RNNs with memory cells that can maintain information over long sequences.

- Suitable for tasks like speech recognition, language modeling, and time series prediction.

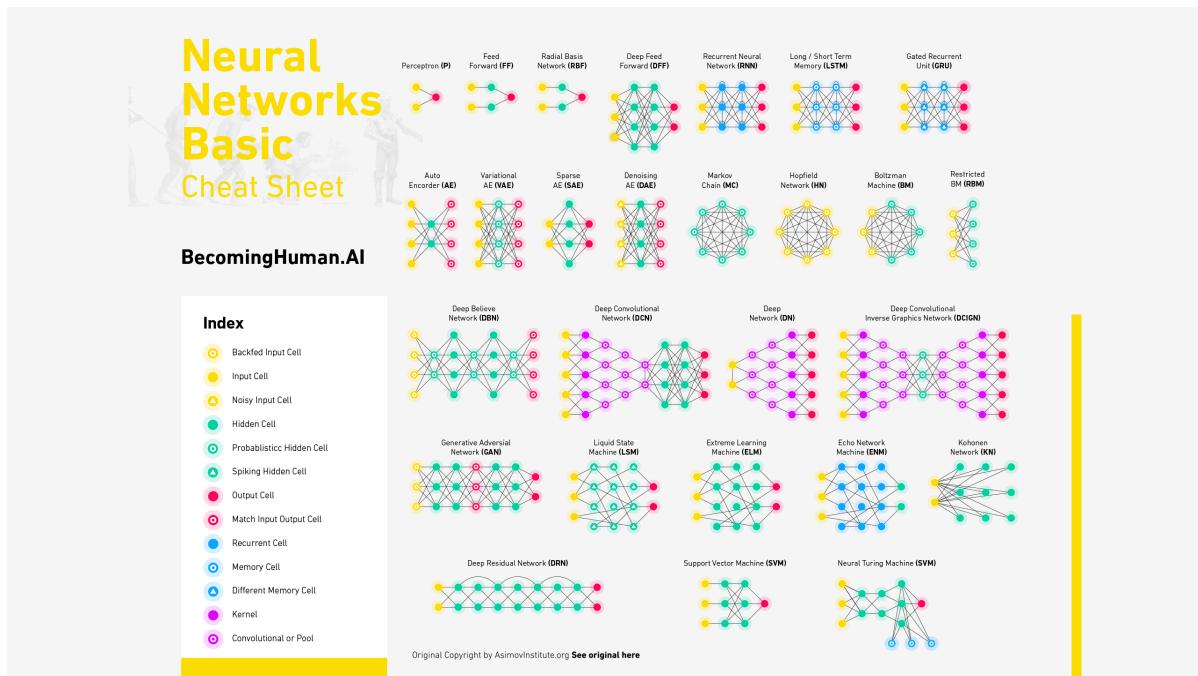
## 5. Generative Adversarial Networks (GAN):

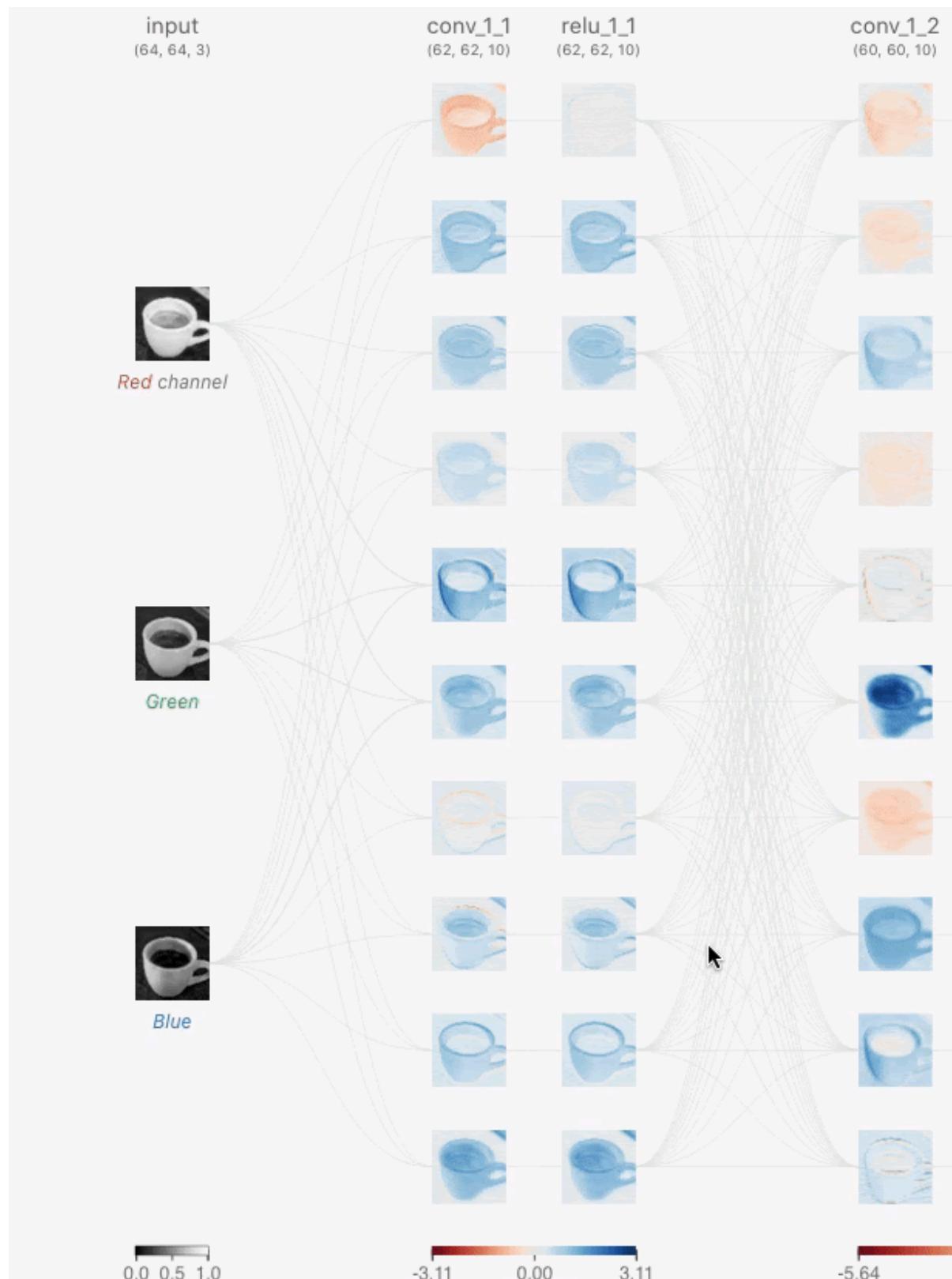
- Comprise two neural networks, a generator and a discriminator, trained simultaneously in a competitive setting.
- Used for generating realistic synthetic data, image-to-image translation, and unsupervised learning.

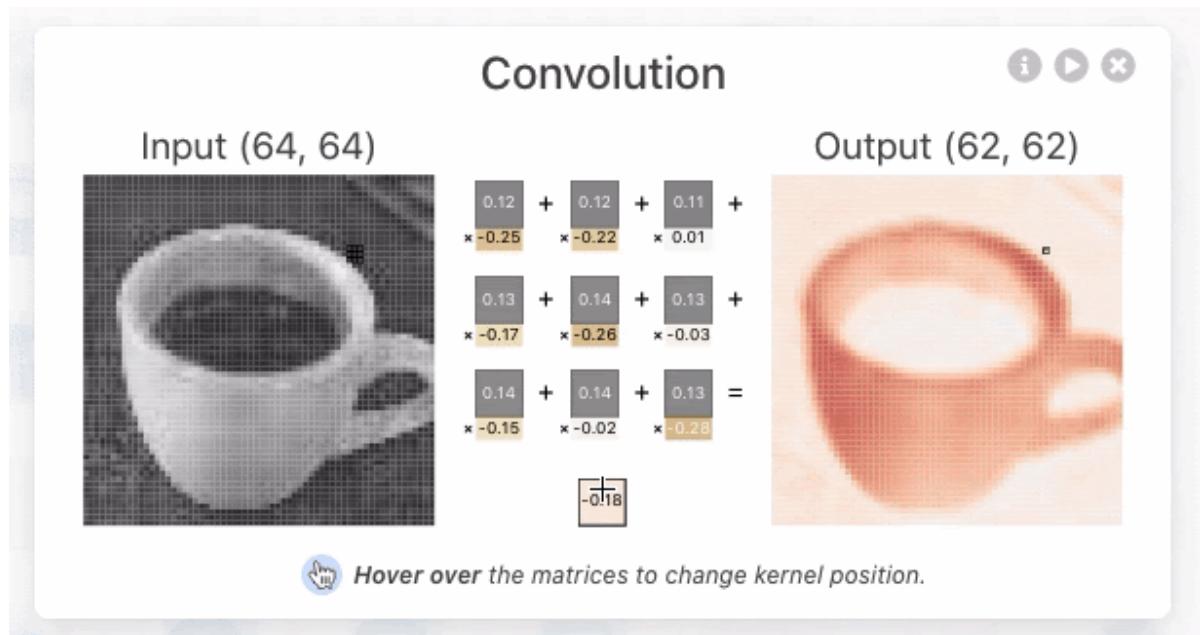
## Applications of Neural Networks:

- **Image Recognition and Computer Vision:** CNNs are widely used for tasks like object detection, image classification, and facial recognition.
- **Natural Language Processing (NLP):** RNNs and LSTM networks are used for tasks like sentiment analysis, language translation, and speech recognition.
- **Healthcare:** Neural networks are employed for disease diagnosis, medical image analysis, drug discovery, and personalized treatment recommendations.
- **Finance:** Used for fraud detection, stock market prediction, algorithmic trading, and credit risk assessment.
- **Autonomous Vehicles:** Neural networks play a crucial role in self-driving cars for tasks like object detection, path planning, and decision-making.

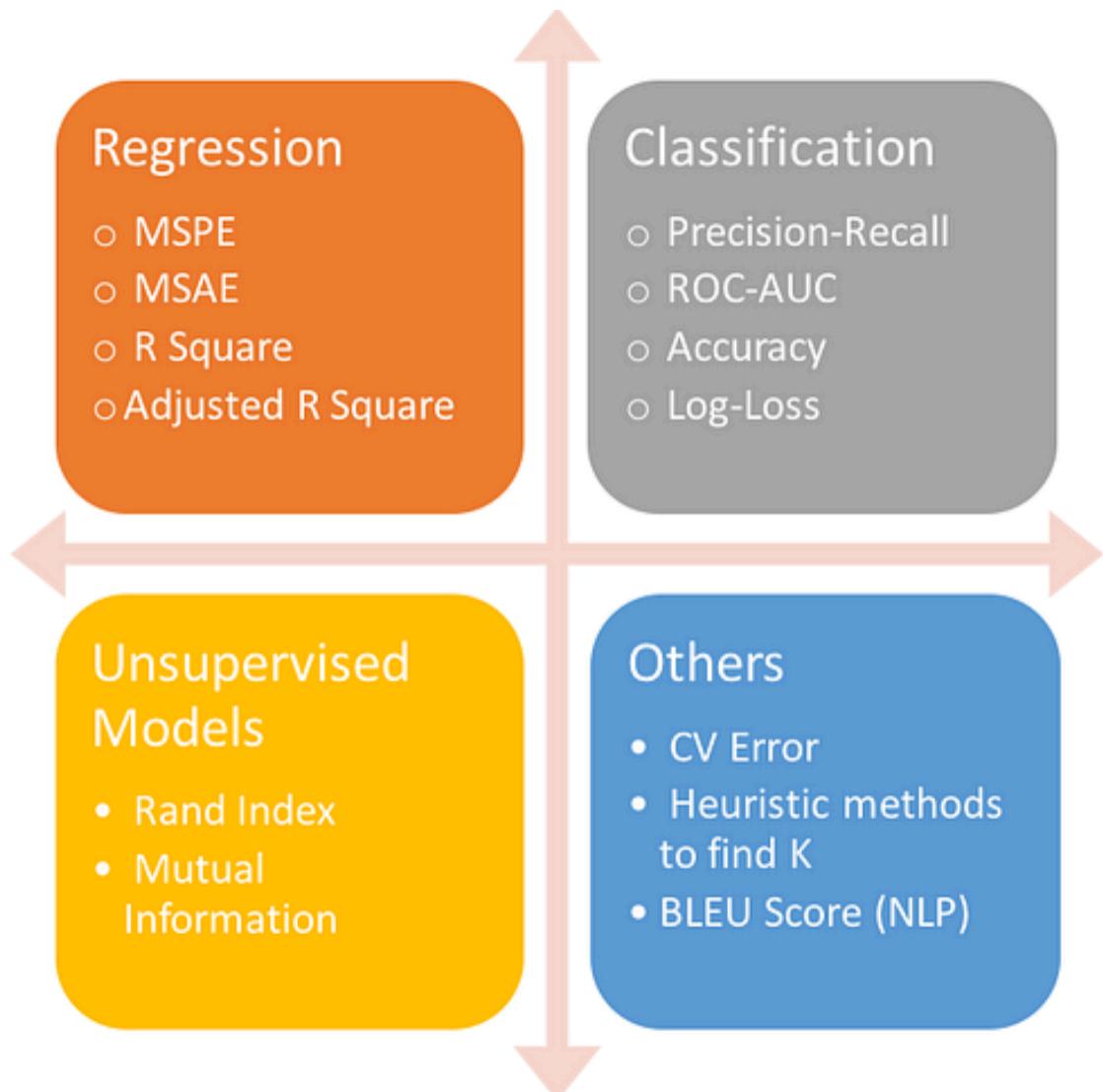
In summary, neural networks are powerful machine learning models capable of learning complex patterns and relationships in data across various domains. With advancements in deep learning techniques and computing resources, neural networks continue to drive innovation and make significant contributions to artificial intelligence and data science.

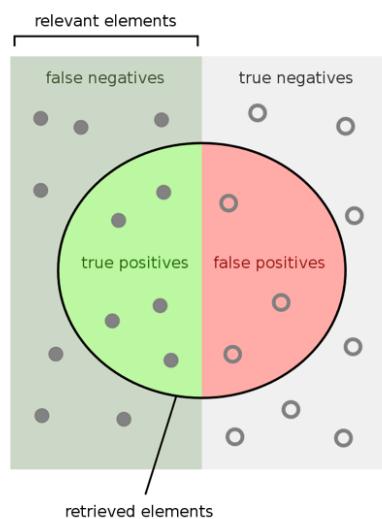






<https://poloclub.github.io/cnn-explainer/>





How many retrieved items are relevant?

$$\text{Precision} = \frac{\text{true positives}}{\text{true positives} + \text{false positives}}$$



How many relevant items are retrieved?

$$\text{Recall} = \frac{\text{true positives}}{\text{true positives} + \text{false negatives}}$$



The predicted value is positive and its positive

ACTUAL VALUES

		Positive	Negative
PREDICTED VALUES	Positive	TP	FP
	Negative	FN	TN

Type I error :  
The predicted value is positive but it False

Positive

TP

FP

Negative

FN

TN

Type II error :  
The predicted value is negative but its positive

The predicted value is Negative and its Negative

## 1 Problem Definition

High level type of problem to be solved like: optimize operations, reduce risk, or provide insights...



EDA

Exploratory Data Analysis.  
Obtaining a realistic overview of the data landscape.



Modeling

Iterative modeling approach to find the best fitting model with the highest accuracy, stability, performance, and (sometimes) transparency



Production

Hosting or providing a model to provide insights on existing data or predict on new data.

## 2 Review EDA Readout

### Data

Getting access to the Data.  
Ensuring permissions and compliance. Making data computationally ready



Feature Engineering

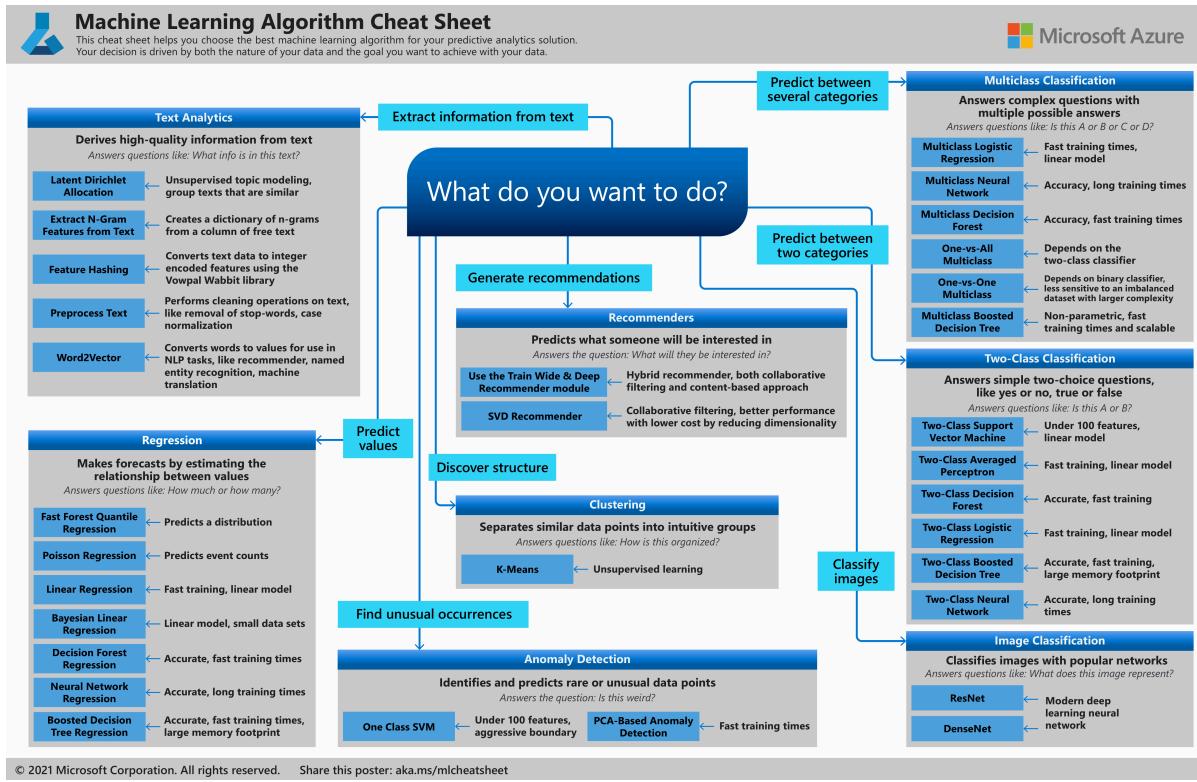
Converting all the data points to something useful to the model.  
Also Feature Selection and Reduction of noise.

## 3 Model Review Analysis

### Error Analysis

Analyzing the model performance in perspective of the business goals. Set thresholds, Gain business approval.

## 4 Final Readout



© 2021 Microsoft Corporation. All rights reserved. Share this poster: aka.ms/mlcheatsheet



Visit KDnuggets.com for more cheatsheets and additional learning resources.

**Scikit-learn CheatSheet**

Scikit-learn is an open-source Python library for all kinds of predictive data analysis. You can perform classification, regression, clustering, dimensionality reduction, model tuning, and data preprocessing tasks.

**Loading the Data**

```
Classification
from sklearn import datasets
X, y = datasets.load_wine(return_X_y=True)

Regression
diabetes = datasets.load_diabetes()
X, y = diabetes.data, diabetes.target
```

**Training And Test Data**

```
from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(
    X, y, random_state=0)
```

**Preprocessing the Data****Standardization**

```
from sklearn.preprocessing import StandardScaler
scaler = StandardScaler()
scaled_X_train = scaler.fit_transform(X_train)
scaled_X_test = scaler.transform(X_test)
```

**Normalization**

```
from sklearn.preprocessing import Normalizer
norm = Normalizer()
norm_X_train = norm.fit_transform(X_train)
norm_X_test = norm.transform(X_test)
```

**Binaryization**

```
from sklearn.preprocessing import Binarizer
binary = Binarizer(threshold=0.0)
binary_X = binary.fit_transform(X)
```

**Encoding Categorical Features**

```
from sklearn.preprocessing import LabelEncoder
lab_enc = LabelEncoder()
y = lab_enc.fit_transform(y)
```

**Imputer**

```
from sklearn.impute import SimpleImputer
imp_mean = SimpleImputer(missing_values=np.nan,
    strategy='mean')
imp_mean.fit_transform(X_train)
```

**Supervised Learning Model****Linear Regression**

```
from sklearn.linear_model import LinearRegression
lr = LinearRegression()
```

**Support Vector Machines**

```
from sklearn.svm import SVC
svm_svc = SVC(kernel='linear')
```

**Naive Bayes**

```
from sklearn.naive_bayes import GaussianNB
gnb = GaussianNB()
```

**Unsupervised Learning Model****Principal Component Analysis**

```
from sklearn.decomposition import PCA
pca = PCA(n_components=2)
```

**K Means**

```
from sklearn.cluster import KMeans
kmeans = KMeans(n_clusters=5, random_state=0)
```

**Model Fitting****Supervised Learning**

```
lr.fit(X_train, y_train)
svm_svc.fit(X_train, y_train)
```

**Unsupervised Learning**

```
model = pca.fit_transform(X_train)
kmeans.fit(X_train)
```

**Prediction****Supervised Learning**

```
y_pred_lr = lr.predict_proba(X_test)
```

```
y_pred_svm = svm_svc.predict(X_test)
```

**Unsupervised Learning**

```
y_pred_kmeans = kmeans.predict(X_test)
```

**Evaluation****Accuracy Score**

```
lr.score(X_test, y_test)
```

```
from sklearn.metrics import accuracy_score
accuracy_score(y_test, y_pred)
```

**Classification Report**

```
from sklearn.metrics import classification_report
print(classification_report(y_test, y_pred))
```

**Mean Squared Error**

```
from sklearn.metrics import mean_squared_error
mean_squared_error(y_test, y_pred)
```

**R2 Score**

```
from sklearn.metrics import r2_score
r2_score(y_test, y_pred)
```

**Adjusted Rand Index**

```
from sklearn.metrics import adjusted_rand_score
adjusted_rand_score(y_test, y_pred)
```

**Cross-Validation**

```
from sklearn.model_selection import cross_val_score
cross_val_score(lr, X, y, cv=5, scoring='f1_macro')
```

**Model Tuning**

```
from sklearn.model_selection import GridSearchCV
parameters = {'kernel':('rbf','linear'), 'C':[1, 10]}
model = GridSearchCV(svm_svm, parameters)
model.fit(X_train, y_train)
print(model.best_score_)
print(model.best_estimator_)
```

Subscribe to KDnuggets News

# Machine Learning Algorithms Cheat Sheet

Machine Learning Model	Hyperparameter Tuning Techniques	Parameters to Reduce Overfitting	Use Case	Evaluation Metrics
Linear Regression	Grid Search, Random Search	Regularization (L1, L2), Feature Selection, Cross-Validation	Regression Analysis	Mean Squared Error, R-squared
Logistic Regression	Grid Search, Random Search	Regularization (L1, L2), Feature Selection, Cross-Validation	Binary Classification	Accuracy, Precision, Recall, F1-score
Decision Trees	Grid Search, Random Search	Pruning, Maximum Depth, Minimum Samples Split, Minimum Samples Leaf	Classification, Regression	Accuracy, Gini Impurity, Entropy,

Machine Learning Model	Hyperparameter Tuning Techniques	Parameters to Reduce Overfitting	Use Case	Evaluation Metrics
				Mean Squared Error
Random Forest	Grid Search, Random Search	Number of Trees, Maximum Features, Maximum Depth, Minimum Samples Split	Classification, Regression	Accuracy, Mean Squared Error, Out-of-Bag Error
Gradient Boosting	Grid Search, Random Search	Learning Rate, Number of Trees, Maximum Depth, Minimum Samples Split	Classification, Regression	Accuracy, Mean Squared Error
Support Vector Machines	Grid Search, Random Search	Kernel Selection, Regularization Parameter (C), Kernel Coefficient (gamma)	Classification, Regression	Accuracy, Precision, Recall, F1-score
k-Nearest Neighbors	Grid Search, Random Search	Number of Neighbors, Distance Metric, Weighting Scheme	Classification, Regression	Accuracy, Mean Squared Error
Neural Networks	Grid Search, Random Search	Number of Layers, Number of Neurons per Layer, Learning Rate, Regularization (Dropout, L2), Batch Size	Classification, Regression	Accuracy, Mean Squared Error
Naive Bayes	Grid Search, Random Search	Smoothing Parameter (Laplace, Lidstone)	Text Classification, Spam Filtering	Accuracy, Precision, Recall, F1-score
Ridge Regression	Grid Search, Random Search	Regularization Strength (Alpha)	Regression Analysis	Mean Squared Error, R-squared
Lasso Regression	Grid Search, Random Search	Regularization Strength (Alpha)	Regression Analysis	Mean Squared Error, R-squared
Elastic Net	Grid Search, Random Search	Regularization Strength (Alpha), L1 Ratio	Regression Analysis	Mean Squared Error, R-squared

- Regularization:** Adding penalty terms to the loss function to discourage complex models.
- Feature Selection:** Choosing relevant features and removing irrelevant or redundant ones.
- Cross-Validation:** Using techniques like k-fold cross-validation to evaluate model performance.
- Pruning:** Trimming branches of decision trees to reduce complexity.
- Maximum Depth:** Limiting the depth of trees to prevent overfitting.
- Minimum Samples Split/Leaf:** Setting thresholds for splitting nodes in decision trees.
- Number of Trees:** In ensemble methods like Random Forest, controlling the number of trees.
- Learning Rate:** In gradient boosting algorithms, adjusting the step size for each iteration.
- Kernel Coefficient (gamma):** In SVMs, controlling the influence of individual training samples.
- Number of Neighbors:** In k-Nearest Neighbors, determining the number of neighbors to consider.
- Distance Metric:** Selecting the appropriate distance metric for k-NN.
- Weighting Scheme:** Assigning weights to neighbors in k-NN based on distance.
- Number of Layers/Neurons:** In neural networks, controlling model complexity.
- Batch Size:** The number of samples propagated through the network at once during training.

15. **Smoothing Parameter:** In Naive Bayes, adjusting the smoothing parameter to account for unseen data.
16. **L1/L2 Ratio:** In Elastic Net, balancing between L1 and L2 regularization.
17. **Alpha:** Regularization strength in Ridge, Lasso, and Elastic Net regression.

These parameters and tuning techniques play a crucial role in mitigating overfitting and optimizing the performance of machine learning models.

## Common Impurity Measures for Decision Trees

Here's a table outlining common impurity measures used in decision trees for classification problems, along with their formulas and explanations:

Impurity Measure	Formula	Description	Suitable for	Considerations
<b>Entropy</b>	$H(S) = - \sum (p_i * \log_2(p_i))$	Measures the level of randomness or uncertainty in a dataset.	Multi-class classification	Sensitive to class imbalance
<b>Gini Impurity</b>	$\text{Gini}(S) = 1 - \sum (p_i^2)$	Measures the probability of a randomly chosen sample being misclassified.	Two-class or multi-class classification	Less sensitive to class imbalance than Entropy
<b>Classification Error</b>	$\text{Error}(S) = (1 - \max(p_i))$	Measures the proportion of samples misclassified by assigning them to the most frequent class.	Two-class classification	Simplest measure, but doesn't consider class distribution

### Explanation of Symbols:

- $H(S)$ : Entropy of the dataset  $S$
- $p_i$ : Proportion of samples belonging to class  $i$
- $\text{Gini}(S)$ : Gini impurity of the dataset  $S$

### Choosing the Right Impurity Measure:

The best choice depends on several factors:

- **Number of classes:** Entropy is generally preferred for multi-class problems, while Gini impurity can be used for both two-class and multi-class problems.
- **Class imbalance:** Gini impurity is less sensitive to imbalanced class distributions compared to Entropy.
- **Interpretability:** Gini impurity might be easier to interpret as it represents the probability of misclassification.

### Additional Considerations:

- Other impurity measures exist, like the Misclassification Cost function, which assigns different costs to misclassifications based on the specific problem.
- The chosen impurity measure influences how the decision tree splits the data at each node.

### Remember:

Understanding these impurity measures is crucial for interpreting and optimizing decision tree models.

## Feature Selection Techniques

While there isn't a single formula encompassing all feature selection techniques, here's a table explaining some common methods along with their approaches:

Technique	Description	Approach (No Formula Involved)	Suitable for	Considerations
<b>Filter Methods</b>	Select features based on statistical properties	- Univariate methods: Rank features based on measures like correlation with target variable (e.g., Pearson correlation, Chi-square) - Multivariate methods: Use feature importance scores from models like Random Forests	High-dimensional datasets	Limited understanding of feature relationships
<b>Wrapper Methods</b>	Evaluate feature subsets using a machine learning model	- Train the model with different feature combinations - Select the subset leading to the best model performance	Any dataset size	Computationally expensive
<b>Embedded Methods</b>	Feature selection integrated into the model training process	- Leverage techniques like L1/L2 regularization in models like Lasso regression - Features with low coefficients are considered less important	High-dimensional datasets	Feature importance interpretation depends on the model

### Notes:

- **Correlation:** Measures the linear relationship between two features. A high correlation (positive or negative) suggests they might be redundant.
- **Chi-square:** Statistical test for independence between categorical features and the target variable. A low p-value indicates dependence, potentially suggesting a good feature.
- **Lasso regression:** A regression model that penalizes the sum of absolute coefficients (L1), driving some coefficients to zero, effectively removing those features.
- **Random Forests:** Ensemble learning method where feature importance is calculated based on the number of times a feature is used in splitting decisions across trees.

### Choosing a Feature Selection Technique:

The best method depends on your data characteristics, model type, and computational resources. Consider:

- **Data size:** Filter methods are generally faster for large datasets.
- **Model interpretability:** Wrapper methods might offer more interpretable feature importance.
- **Computational cost:** Wrapper methods can be expensive for complex models.

Remember, feature selection is an iterative process. Evaluate different techniques and compare their impact on model performance.

Here's a table summarizing common feature selection methods, along with a brief description of each method and the statistical or mathematical intuition behind it:

Feature Selection Method	Description	Statistical/Mathematical Intuition
Filter Methods	Filter methods evaluate the relevance of features based on their statistical properties such as correlation with the target variable, variance, or mutual information. Features are ranked or selected independently of the model.	<ul style="list-style-type: none"> <li>- <b>Correlation:</b> Measures the strength and direction of the linear relationship between each feature and the target variable. High correlation indicates higher importance.</li> <li>- <b>Variance Thresholding:</b> Features with low variance are considered less informative.</li> <li>- <b>Mutual Information:</b> Measures the amount of information obtained about one variable through the other variable. Features with high mutual information with the target variable are considered more informative.</li> </ul>
Wrapper Methods	Wrapper methods select subsets of features based on the performance of a specific machine learning algorithm. Features are selected or eliminated based on their impact on the model's performance during training.	<ul style="list-style-type: none"> <li>- <b>Forward Selection:</b> Iteratively adds one feature at a time to the model and evaluates performance.</li> <li>- <b>Backward Elimination:</b> Starts with all features and removes one feature at a time based on their impact on model performance.</li> <li>- <b>Recursive Feature Elimination (RFE):</b> Ranks features by their importance and recursively removes the least important features.</li> </ul>
Embedded Methods	Embedded methods incorporate feature selection as part of the model training process. Feature importance is derived from the model's coefficients, weights, or other intrinsic properties during training.	<ul style="list-style-type: none"> <li>- <b>Lasso (L1 Regularization):</b> Encourages sparsity in coefficients, effectively eliminating less important features.</li> <li>- <b>Random Forest Feature Importance:</b> Measures the importance of each feature by evaluating the decrease in impurity when the feature is used for splitting in decision trees.</li> <li>- <b>Gradient Boosting Feature Importance:</b> Measures the improvement in model performance when a feature is used for splitting in gradient boosting trees.</li> </ul>
Principal Component Analysis (PCA)	PCA is a dimensionality reduction technique that transforms the original features into a new set of uncorrelated variables called principal components. These components capture the maximum variance in the data and can be used as features.	<ul style="list-style-type: none"> <li>- PCA identifies the directions in which the data varies the most, known as principal components.</li> <li>- By selecting a subset of principal components that capture most of the variance in the data, PCA reduces the dimensionality while retaining as much information as possible.</li> </ul>
t-Distributed Stochastic Neighbor Embedding (t-SNE)	t-SNE is a non-linear dimensionality reduction technique used primarily for visualization. It reduces high-dimensional data to two or three dimensions while preserving the local structure of the data, making it suitable for exploring complex datasets.	<ul style="list-style-type: none"> <li>- t-SNE minimizes the divergence between the pairwise similarities of the high-dimensional data and those of the low-dimensional embedding.</li> <li>- It focuses on preserving the local structure of the data, making it effective for visualizing clusters and patterns in complex datasets.</li> </ul>

These feature selection methods employ various statistical and mathematical techniques to identify and prioritize relevant features while reducing dimensionality and noise in the data. Each method has its strengths and weaknesses and is suited to different types of datasets and modeling tasks.

## Statistical Measures

Here's a table outlining various statistical measures commonly used in data science, along with their formulas, use cases, and brief explanations:

Statistical Measure	Formula	Use Cases	Explanation
<b>Central Tendency</b>			
Mean	$\mu (\text{mu}) = \sum(x_i) / N$	Descriptive statistics	The arithmetic average of a dataset. Sensitive to outliers.
Median	$M = \text{Middle value}$ (odd n) or average of middle two (even n)	Descriptive statistics	The "center" value of a dataset, less sensitive to outliers than mean.
Mode	Value with highest frequency	Descriptive statistics	The most frequently occurring value. May not be unique.
<b>Dispersion</b>			
Variance	$\sigma^2 (\text{sigma squared}) = \sum(x_i - \mu)^2 / N$	Descriptive statistics	Measures spread around the mean. Sensitive to outliers.
Standard Deviation	$\sigma (\text{sigma}) = \sqrt{\sigma^2}$	Descriptive statistics	Square root of variance, easier to interpret than variance.
Range	$R = \text{Max}(x) - \text{Min}(x)$	Descriptive statistics	Difference between maximum and minimum values.
Interquartile Range (IQR)	$IQR = Q_3 - Q_1$	Descriptive statistics	Spread of the middle 50% of the data (between 1st & 3rd quartiles).
<b>Relationship Between Variables</b>			
Pearson Correlation Coefficient	$r = \frac{\sum((x_i - \bar{x})(y_i - \bar{y}))}{\sqrt{(\sum(x_i - \bar{x})^2)(\sum(y_i - \bar{y})^2)}}$	Measures linear relationship between two continuous variables	-1 to 1, strong positive (1), strong negative (-1), or no linear relationship (0). Sensitive to outliers.
Tau (Kendall's Tau Rank Correlation)	$\tau (\text{Tau})$	Measures monotonic relationship between two ranked variables	-1 to 1, strong positive (1), strong negative (-1), or weak relationship. Less sensitive to outliers than Pearson.
Chi-Square Test	$\chi^2 (\text{Chi-Square}) = \sum ((O - E)^2 / E)$	Hypothesis testing	Measures independence between categorical variables.
<b>Hypothesis Testing</b>			
Z-test	$Z = (\bar{x} - \mu_0) / (\sigma / \sqrt{N})$	Compares a sample mean to a hypothesized population mean	Numerical data (assuming normality). Used for significance testing.
t-test		Compares means of two groups	Numerical data (assuming normality).

- One-sample t-test (compares mean to a value).
- Two-sample t-test (independent or paired samples). || ANOVA (Analysis of Variance) | F-statistic = ( MS\_between / MS\_within ) | Compares means of multiple groups | Numerical data (assuming normality). Used for significance testing across multiple groups. |

### Additional Notes:

- $\Sigma$  (sigma) denotes summation.
- $x_i$  (x-subscript-i) represents an individual data point.
- N represents the total number of data points.

- $\mu_0$  (mu-subscript-zero) represents the hypothesized population mean.
- O represents observed frequency.
- E represents expected frequency.
- MS\_between and MS\_within are mean squares from ANOVA calculations.

### **Remember:**

- This table provides a foundation for statistical measures in data science. Many more techniques exist.
- Statistical software can automate calculations, but understanding the concepts is crucial.

### **Variance:**

- The formula  $\sigma^2 = \sum(x_i - \mu)^2 / N$  represents the population variance. In this formula,
- $\sigma^2$ : Is the population variance, which is the square of the population standard deviation
- $\sum(x_i - \mu)^2$ : Is the sum of all the squared deviation scores of the population
- N: Is the number of scores in the population

### **Correlation Coefficient:**

- $x_i$  and  $y_i$  are individual data points, and  $\bar{x}$  and  $\bar{y}$  are the means of the respective variables

**Standard Deviation:** The standard deviation  $\sigma$  is calculated as the square root of the variance  $\sigma^2$ . Mathematically, it is represented as:

$$\sigma = \sqrt{\sigma^2}$$

Where:

- $\sigma$  represents the standard deviation.
- $\sigma^2$  represents the variance.

The variance  $\sigma^2$  is calculated as the average of the squared differences between each data point and the mean of the dataset. The formula for variance is:

$$\sigma^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n}$$

Where:

- $x_i$  represents each individual data point.
- $\bar{x}$  represents the mean of the dataset.
- $n$  represents the total number of data points.

Once the variance is calculated, the standard deviation is obtained by taking the square root of the variance.

I hope this clarifies the standard deviation formula. Thank you for your patience.