

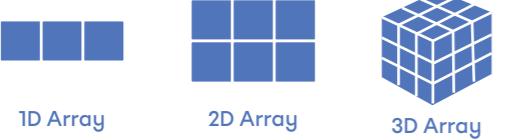
NumPy

The NumPy logo consists of the word "NumPy" in a large, bold, blue sans-serif font. To the right of "NumPy" is a 3D cube icon composed of blue and teal squares.

Cheat Sheet

```
import numpy as np
```

NumPy Arrays



Creating NumPy Arrays

```
# One dimensional array
>>> np.array([1, 2, 3])
array([1, 2, 3])

# Two dimensional array
>>> np.array([(1,2,3),(4,5,6)])
array([[1, 2, 3],
       [4, 5, 6]])

# Array of zeros with 3 elements
>>> np.zeros(3)
array([0., 0., 0.])

# 2x2 array filled with ones
>>> np.ones((2,2))
array([[1., 1.],
       [1., 1.]))

# 3x3 identity matrix
>>> np.eye(3)
array([[1., 0., 0.],
       [0., 1., 0.],
       [0., 0., 1.]))

# 2x3 array filled with 4s
>>> np.full((2,3),4)
array([[4, 4, 4],
       [4, 4, 4]])

# 1D array of length 3 with random integers between 0 and 10
>>> np.random.randint(0, 10, size=3)
array([8, 2, 7])

# 2x3 array of random numbers between 0 and 1
>>> np.random.rand(2, 3)
array([[0.60738971, 0.84209229, 0.64491569],
       [0.50990228, 0.25552633, 0.56484642]]))

# 2x3 array of random numbers from a normal distribution with
and standard deviation 1
>>> np.random.normal(0, 1, (2, 3))
array([[ 0.36527849, -2.48435406,  0.77739812],
       [ 0.07923544, -0.30833118,  0.32393125]]))

# 1D array with values from 0 to 10 with a step of 3
>>> np.arange(0,10,3)
array([0, 3, 6, 9])

# 1D array with 3 evenly spaced values between 0 and 10
>>> np.linspace(0,10,3)
array([ 0.,  5., 10.])
```

Attributes of Arrays

```
>>> a = np.array([(1,2,3),(4,5,6)])  
  
>>> a.ndim  
2  
  
>>> a.shape  
(2, 3)  
  
>>> a.size  
6  
  
>>> a.dtype  
dtype('int32')  
  
>>> a.astype(float)  
array([[1., 2., 3.],  
       [4., 5., 6.]])
```

Data Manipulation

Reshaping

```
>>> a = np.array([(1,2,3),(4,5,6)])
array([[1, 2, 3],
       [4, 5, 6]])

>>> a.flatten()
array([1, 2, 3, 4, 5, 6])

>>> a.resize((6,1))
array([[1],
       [2],
       [3],
       [4],
       [5],
       [6]])

>>> a.T
array([[1, 4],
       [2, 5],
       [3, 6]])

>>> b = np.random.randint(1, 10, size=9)
array([4, 7, 1, 8, 2, 3, 3, 7, 8])

>>> b.reshape(3, 3)
array([[4, 7, 1],
       [8, 2, 3],
       [3, 7, 8]])
```

Combining

```
>>> a = np.array([1, 2, 3, 4, 5])
>>> b = np.array([6, 7, 8, 9, 10])

# Concatenating arrays horizontally (along
# columns)
>>> np.hstack((a, b))
array([ 1,  2,  3,  4,  5,  6,  7,  8,  9, 10])

# Concatenating arrays vertically (along rows)
>>> np.vstack((a, b))
array([[ 1,  2,  3,  4,  5,
       6,  7,  8,  9, 10]])

>>> np.concatenate((a, b))
array([ 1,  2,  3,  4,  5,  6,  7,  8,  9,
       10])

>>> np.stack((a, b))
array([[ 1,  2,  3,  4,  5],
       [ 6,  7,  8,  9, 10]])
```

Indexing / Fancy Indexing

```
>>> a = np.array([1, 2, 3, 4, 5])
array([1, 2, 3, 4, 5])

# Access the first element of the array
>>> a[0]
1

>>> indices = np.array([0, 2, 4])
# Access elements at indices 0, 2, and 4
>>> a[indices]
array([1, 3, 5])

>>> b = np.array([[1, 2, 3], [4, 5, 6], [7, 8, 9]])
array([[1, 2, 3],
       [4, 5, 6],
       [7, 8, 9]])

# Access the element at row 1, column 2 of the array
>>> b[1, 2]
6
```

Slicing

```
# Slice the array from index 2 to index 4 (exclusive)
>>> a[2:5]
array([3, 4, 5])

# Slice the array from index 1 to the end
>>> a[1:]
array([2, 3, 4, 5])

# Slice the array from the beginning to index 2(exclusive)
>>> a[:3]
array([1, 2, 3])

# Slice the array in reverse order
>>> a[::-1]
array([5, 4, 3, 2, 1])

# Slice the first row
>>> b[0, :]
array([1, 2, 3])

# Slice the second column
>>> b[:, 1]
array([2, 5, 8])
```

Subsetting

```
>>> a[a > 2]
array([3, 4, 5])

>>> a[(a > 1) & (a < 4)]
array([2, 3])

>>> a[a != 3]
array([1, 2, 4, 5])

>>> np.where(a > 2, a, 0)
array([0, 0, 3, 4, 5])
```

Numerical Operations

Vector Operations

```
>>> a = np.array([1, 2, 3, 4, 5])
>>> b = np.array([6, 7, 8, 9, 10])

>>> a+b
array([ 7,  9, 11, 13, 15])

>>> a-b
array([-5, -5, -5, -5, -5])

>>> b/a
array([6. , 3.5 , 2.66666667, 2.25 , 2. ])

>>> b//a
array([6, 3, 2, 2, 2], dtype=int32)

>>> a*b
array([ 6, 14, 24, 36, 50])

>>> a**b
array([ 1, 128, 6561, 262144, 9765625], dtype=int32)

>>> np.exp(a)
array([ 2.71828183, 7.3890561 , 20.08553692, 54.59815003, 148.4131591 ])

>>> np.log1p(b)
array([1.94591015, 2.07944154, 2.19722458, 2.30258509, 2.39789527])

>>> np.sin(a)
array([-0.95892427, 0.84147098],
 [ 0.84147098, 0.14112001])

>>> np.sqrt(b)
array([3.46410162, 3.16227766])

>>> np.ceil(np.array([1.2, 2.3, 3.4, 4.5, 5.6, 7.6]))
array([2., 3., 4., 5., 6., 8.])

>>> np.floor(np.array([1.2, 2.3, 3.4, 4.5, 5.6, 7.6]))
array([1., 2., 3., 4., 5., 7.])

>>> np.round(np.array([1.2, 2.3, 3.4, 4.5, 5.6, 7.6]))
array([1., 2., 3., 4., 6., 8.])

>>> np.abs(np.array([-1, 2, -3, -4, 5, 6, -7]))
array([1, 2, 3, 4, 5, 6, 7])
```

Statistics

```
>>> a = np.array([1, 2, 3, 4, 5])  
  
>>> np.mean(a)  
3.0  
  
>>> a.sum()  
15  
  
>>> a.max()  
5  
  
>>> a.min()  
1  
  
>>> np.var(a)  
2.0  
  
>>> np.std(a)  
1.4142135623730951  
  
>>> np.corrcoef(a)  
1.0
```

Linear Algebra

```
# Solve the linear system


$$\begin{cases} 5x_0 + x_1 = 12 \\ x_0 + 3x_1 = 10 \end{cases}$$


>>> a = np.array([[5, 1], [1, 3]])
>>> b = np.array([12, 10])

>>> np.linalg.solve(a, b)
array([1.85714286, 2.71428571])

# Solve the linear system


$$\begin{cases} x_0 + 2x_1 + 3x_2 = 10 \\ 4x_0 + 5x_1 + 6x_2 = 11 \\ 7x_0 + 8x_1 + 9x_2 = 12 \end{cases}$$


>>> a = np.array([[1, 2, 3], [4, 5, 6], [7, 8, 9]])
>>> b = np.array([10, 11, 12])

>>> np.linalg.solve(a, b)
array([ 25.33333333, 41.66666667, -16.])
```

PANDAS



Cheat Sheet

```
import pandas as pd
```

Data Structures

Series

```
# Create a series
>>> s = pd.Series([1, 2, 3], index=['A', 'B', 'C'])

Index   Data
A        1
B        2
C        3

>>> s.index
Index(['A', 'B', 'C'], dtype='object')

>>> s.dtype
dtype('int64')

>>> s.size
3

>>> s.ndim
1

>>> s.values
array([1, 2, 3], dtype=int64)
```

DataFrame

```
# Create a DataFrame
>>> df = pd.DataFrame({'ID': [1, 2, 3],
                      'Name': ['Alex', 'Brian', 'David'],
                      'Profession': ['DA', 'DE', 'DS'],
                      index = ['1 (One)', '2 (Two)', '3 (Three)'])

          ID    Name  Profession
1 (One)  1    Alex       DA
2 (Two)  2    Brian      DE
3 (Three) 3    David      DS
```

```
# Get the shape of the DataFrame
>>> df.shape
(3, 3)
```

```
# Get the column names of the DataFrame
>>> df.columns
Index(['ID', 'Name', 'Profession'], dtype='object')
```

```
# Get the index information of the DataFrame
>>> df.index
Index(['1 (One)', '2 (Two)', '3 (Three)'], dtype='object')
```

```
# Get the data types of the DataFrame columns
>>> df.dtypes
```

```
ID      int64
Name     object
Profession  object
dtype: object
```

```
# Check if there are any missing values in the DataFrame
>>> df.isnull().values.any()
False
```

```
# Get the count of missing values in each column
>>> df.isnull().sum()
```

```
ID      0
Name     0
Profession  0
dtype: int64
```

```
# Get the count of non-null values in each column
>>> df.count()
```

```
ID      3
Name     3
Profession  3
dtype: int64
```

```
# Generate descriptive statistics of the DataFrame (transposed)
>>> df.describe().T
```

Importing & Exporting

```
# Read a CSV file and create a DataFrame
>>> pd.read_csv('filename.csv')

# Save the DataFrame to a CSV file
df.to_csv('filename.csv')

# Read an Excel file and create a DataFrame
>>> pd.read_excel('filename.xlsx')

# Save the DataFrame to an Excel file
>>> df.to_excel('filename.xlsx')
```

Data Manipulation

Selection

Selecting Rows

```
# Select a specific row by its label
>>> df.loc['1 (One)']:
```

```
# Select a specific row by its index
>>> df.loc[1]
```

Selecting Columns

```
# Select a single column by its name
>>> df['ID']
```

```
# Select multiple columns by their name
>>> df[['ID', 'Profession']]
```

```
# Select a single column using label-based indexing
>>> df.loc[:, 'Name']
```

```
# Select a single column using integer-based indexing
>>> df.iloc[:, 1]
```

Selecting Rows and Columns

```
# Select a specific cell by its label-based row and column indices
>>> df.loc['2 (Two)', 'Name']
```

```
# Select a specific cell by its integer-based row and column indices
>>> df.iloc[1, 1]
```

```
# Select specific rows and columns using label-based indexing
>>> df.loc['1 (One)', ['ID', 'Profession']]
```

```
# Select specific rows and columns using integer-based indexing
>>> df.iloc[1, ['ID', 'Profession']]
```

Conditional Selection

```
# Select rows based on a condition
>>> df[df['ID'] > 1]
```

```
# Select rows based on multiple conditions using logical operators
>>> df[(df['ID'] > 2) & (df['Profession'] == 'DS')]
```

```
# Select rows where a column value is present in a given list
>>> df.loc[df['Name'].isin(['Alex', 'David'])]
```

Deleting & Adding

```
# Drop rows with any missing values
>>> df.dropna()
```

```
# Drop columns with any missing values
>>> df.dropna(axis=1)
```

```
# Drop columns with fewer than n non-null values
>>> df.dropna(axis=1, thresh=n)
```

```
# Fill missing values with a specified value
df.fillna(value)
```

Filtering & Sorting

```
# Filter columns using regular expression matching
```

```
>>> df.filter(regex=regex)
```

```
# Sort DataFrame by values in a specific column in descending order
```

```
>>> df.sort_values('ID', ascending=False)
```

Editing

```
# Perform quantile-based discretization of values in 'col' into 3 bins with custom labels
```

```
>>> pd.qcut(df['col'], 3, labels=qcut_labels)
```

```
# Perform value-based discretization of values in 'col' into custom bins with labels
```

```
>>> pd.cut(df['col'], bins=cut_bins, labels=cut_labels)
```

```
# Create a pivot table using 'col1' as index, 'col2' as columns, and 'col3' as values
```

```
>>> df.pivot(index='col1', columns='col2', values='col3')
```

```
# Reset the index of the DataFrame
```

```
>>> df.reset_index()
```

```
# Rename specific columns of the DataFrame
```

```
>>> df.rename(columns = rename_list)
```

Grouping & Aggregation

```
# Perform grouping operation on 'col' and obtain a GroupBy object
```

```
>>> df.groupby('col')
```

```
# Perform grouping operation on multiple columns ('col1' and 'col2') and calculate the mean of 'col3'
```

```
>>> df.groupby(['col1', 'col2']).agg({"col3": "mean"})
```

```
# Perform grouping operation on multiple columns ('col1' and 'col2') and calculate the mean of 'col3' and count of 'col4'
```

```
>>> df.groupby(['col1', 'col2']).agg({'col3': 'mean', 'col4': 'count'})
```

```
# Grouping the DataFrame by the first level of the index
```

```
>>> df.groupby(level=0)
```

```
# An example
```

```
>>> df = pd.DataFrame({
    'Name': ['John', 'Alice', 'John', 'Alice', 'Bob'],
    'City': ['New York', 'Paris', 'London', 'Paris', 'London'],
    'Age': [30, 25, 35, 28, 40],
    'Salary': [50000, 60000, 55000, 45000, 70000]})
```

```
>>> df.groupby('City').agg({'Salary': 'mean', 'Age': 'max'})
```

Name	City	Age	Salary		Salary	Age
				City		
				London		62500.0
				New York		50000.0
				Paris		52500.0
						28

Joining

```
# Merge operation
```

```
>>> df1 = pd.DataFrame({
    'key': ['A', 'B', 'C', 'D'],
    'value1': [1, 2, 3, 4]})
```

```
>>> df2 = pd.DataFrame({
    'key': ['B', 'D', 'E', 'F'],
    'value2': [5, 6, 7, 8]})
```

```
>>> pd.merge(df1, df2, on='key', how='inner')
```

Key	Value1
0	A
1	B
2	C
3	D

Key	Value2
0	B
1	D
2	E
3	F

Key	Value1	Value2
0	A	1
1	B	2
2	C	3
3	D	4

```
>>> df4 = pd.DataFrame({
    'key': ['B', 'D', 'E', 'F'],
    'value4': ['orange', 'grape', 'kiwi', 'lemon']})
```

```
>>> df3.set_index('key').join(df4.set_index('key'), how='inner')
```

Key	Value3
0	A
1	B
2	C
3	D

Key	Value2
0	B
1	D
2	E
3	F

Value1	Value2
B	banana
D	date
E	kiwi
F	lemon

```
# Example for concatenate operation
```

```
>>> df5 = pd.DataFrame({
    'A': ['A0', 'A1', 'A2'],
    'B': ['B0', 'B1', 'B2']})
```

```
>>> df6 = pd.DataFrame({
    'A': ['A3', 'A4', 'A5'],
    'B': ['B3', 'B4', 'B5']})
```

```
>>> pd.concat([df5, df6], axis=0)
```

A	B
0	A0
1	A1
2	A2
0	A3
1	A4
2	A5
0	B0
1	B1
2	B2
0	B

SCIKIT-LEARN



Cheat Sheet



Preprocessing

Splitting data into train and test sets

```
from sklearn.model_selection import train_test_split  
X = df[["Independent Variables"]]  
y = df[["Target Variable"]]  
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2)
```

Handling Missing Values

```
from sklearn.impute import SimpleImputer, KNNImputer

# Dropping Rows or Columns
df.dropna(axis=0)
df.dropna(axis=1)

# Imputation
imputer = SimpleImputer(strategy='mean')
X_train_imputed = imputer.fit_transform(X_train)
X_test_imputed = imputer.fit_transform(X_test)

# K-Nearest Neighbors (KNN) Imputation
knn_imputer = KNNImputer()
X_train_knn_imputed = knn_imputer.fit_transform(X_train)
X_test_knn_imputed = knn_imputer.fit_transform(X_test)
```

Handling Outliers

```

def suppress_outliers_iqr(df, col_name, q1=0.25, q3=0.75, multiplier=1.5):
    quartile1 = df[col_name].quantile(q1)
    quartile3 = df[col_name].quantile(q3)
    interquartile_range = quartile3 - quartile1
    lower_bound = quartile1 - multiplier * interquartile_range
    upper_bound = quartile3 + multiplier * interquartile_range

    # Suppress outliers by replacing them with the lower/upper bounds
    df.loc[df[col_name] < lower_bound, col_name] = lower_bound
    df.loc[df[col_name] > upper_bound, col_name] = upper_bound

# Example usage
suppress_outliers_iqr(dataframe, 'A')

```

Feature Scaling

```
from sklearn.preprocessing import StandardScaler, MinMaxScaler

# Numeric feature scaling (StandardScaler)
scaler = StandardScaler()
X_train_scaled = scaler.fit_transform(X_train)
X_test_scaled = scaler.transform(X_test)

# Numeric feature scaling (MinMaxScaler)
scaler = MinMaxScaler()
X_train_scaled = scaler.fit_transform(X_train)
X_test_scaled = scaler.transform(X_test)
```

Feature Encoding

```
from sklearn.preprocessing import OneHotEncoder

# Categorical feature encoding (One-Hot Encoder)
encoder = OneHotEncoder()
X_train_encoded = encoder.fit_transform(X_train)
X_test_encoded = encoder.transform(X_test)
```

Supervised Learning Algorithms

```
# Linear Regression
from sklearn.linear_model import LinearRegression
model = LinearRegression()
model.fit(X_train, y_train)

# Logistic Regression
from sklearn.linear_model import LogisticRegression
model = LogisticRegression()
model.fit(X_train, y_train)

# K-Nearest Neighbors (KNN)
from sklearn.neighbors import KNeighborsClassifier
model = KNeighborsClassifier()
model.fit(X_train, y_train)

# CART
from sklearn.tree import DecisionTreeClassifier
model = DecisionTreeClassifier()
model.fit(X_train, y_train)

# Random Forests
from sklearn.ensemble import RandomForestClassifier
model = RandomForestClassifier()
model.fit(X_train, y_train)

# GBM
from sklearn.ensemble import GradientBoostingClassifier
model = GradientBoostingClassifier()
model.fit(X_train, y_train)

# XGBoost
!pip install xgboost
from xgboost import XGBClassifier
model = XGBClassifier()
model.fit(X_train, y_train)

# LightGBM
!pip install lightgbm
from lightgbm import LGBMClassifier
model = LGBMClassifier()
model.fit(X_train, y_train)

# Catboost
!pip install catboost
from catboost import CatBoostClassifier
model = CatBoostClassifier()
model.fit(X_train, y_train)
```

Unsupervised Learning Algorithms

```
# K-Means
from sklearn.cluster import KMeans
model = KMeans(n_clusters=3)
model.fit(df)

# Hierarchical Clustering
from sklearn.cluster import AgglomerativeClustering
model = AgglomerativeClustering(n_clusters=5),
model.fit(df)

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

Model Evaluation

Classification metrics

```
from sklearn.metrics import accuracy_score, precision_score,  
recall_score, f1_score, roc_auc_score, classification_report,  
confusion_matrix  
  
y_pred = model.predict(X_test)  
  
# Confusion Matrix  
confusion_matrix(y_test, y_pred)  
  


|        |       | Predicted |      |
|--------|-------|-----------|------|
|        |       | False     | True |
| Actual | True  | TN        | FN   |
|        | False | FN        | TP   |

  
# Accuracy = ( TP + TN ) / ( TP + TN + FP + FN )  
accuracy_score(y_test, y_pred)  
  
# Precision = TP / ( TP + FP )  
precision_score(y_test, y_pred)  
  
# Recall = TP / ( TP + FN )  
recall_score(y_test, y_pred)  
  
# F1 Score = TN / ( TN + FP )  
f1_score(y_test, y_pred)  
  
# AUC  
roc_auc_score(y_test, y_pred)  
  
# Classification Report  
classification_report(y_test, y_pred)
```

Regression metrics

```
from sklearn.metrics import mean_absolute_error,  
mean_squared_error, r2_score  
from sklearn.model_selection import cross_val_score  
  
y_pred = model.predict(X_test)  
  
# MAE  
mean_absolute_error(y_test, y_pred)  
  
# MSE  
mean_squared_error(y_test, y_pred)  
  
# RMSE  
np.sqrt(mean_squared_error(y, y_pred))
```

Clustering metrics

```
from sklearn.metrics import adjusted_rand_score,  
homogeneity_score, v_measure_score  
  
# Adjusted Rand Index  
adjusted_rand_score(y_true, y_pred)  
  
# Homogeneity  
homogeneity_score(y_true, y_pred)  
  
# V-measure  
metrics.v_measure_score(y_true, y_pred))
```

Model Tuning

```

from sklearn.model_selection import GridSearchCV, RandomizedSearchCV

# Logistic Regression parameter optimization
lr_param_grid = {'C': [0.1, 1, 10], 'penalty': ['l1', 'l2']}
lr_grid_search = GridSearchCV(LogisticRegression(), lr_param_grid, cv=5)
lr_grid_search.fit(X_train, y_train)
lr_grid_search.best_params_ # Best parameters for Logistic Regression
lr_grid_search.best_score_ # Best score for Logistic Regression

# Decision Tree parameter optimization
dt_param_grid = {'max_depth': [None, 5, 10], 'min_samples_split': [2, 5, 10]}
dt_grid_search = GridSearchCV(DecisionTreeClassifier(), dt_param_grid, cv=5)
dt_grid_search.fit(X_train, y_train)
dt_grid_search.best_params_ # Best parameters for Decision Tree
dt_grid_search.best_score_ # Best score for Decision Tree

# Random Forest parameter optimization
rf_param_grid = {'n_estimators': [100, 200, 300], 'max_depth': [None, 5, 10]}
rf_random_search = RandomizedSearchCV(RandomForestClassifier(), rf_param_grid,
n_iter=5, cv=5)
rf_random_search.fit(X_train, y_train)
rf_random_search.best_params_ # Best parameters for Random Forest
rf_random_search.best_score_ # Best score for Random Forest

# K-Nearest Neighbors parameter optimization
knn_param_grid = {'n_neighbors': [3, 5, 7], 'weights': ['uniform','distance']}
knn_random_search = RandomizedSearchCV(KNeighborsClassifier(), knn_param_grid,
n_iter=5, cv=5)
knn_random_search.fit(X_train, y_train)
knn_random_search.best_params_ # Best parameters for K-Nearest Neighbors
knn_random_search.best_score_ # Best score for K-Nearest Neighbors

# GBM parameter optimization
gbm_param_grid = {'n_estimators': [100, 200, 300], 'max_depth': [None, 5, 10]}
gbm_grid_search = GridSearchCV(GradientBoostingClassifier(), gbm_param_grid, cv=5)
gbm_grid_search.fit(X_train, y_train)
gbm_grid_search.best_params_ # Best parameters for GBM
gbm_grid_search.best_score_ # Best score for GBM

# LightGBM parameter optimization
lgb_param_grid = {'n_estimators': [100, 200, 300], 'max_depth': [None, 5, 10]}
lgb_grid_search = GridSearchCV(lgb.LGBMClassifier(), lgb_param_grid, cv=5)
lgb_grid_search.fit(X_train, y_train)
lgb_grid_search.best_params_ # Best parameters for LightGBM
lgb_grid_search.best_score_ # Best score for LightGBM

# XGBoost parameter optimization
xgb_param_grid = {'n_estimators': [100, 200, 300], 'max_depth': [None, 5, 10]}
xgb_grid_search = GridSearchCV(xgb.XGBClassifier(), xgb_param_grid, cv=5)
xgb_grid_search.fit(X_train, y_train)
xgb_grid_search.best_params_ # Best parameters for XGBoost
xgb_grid_search.best_score_ # Best score for XGBoost

# CatBoost parameter optimization
cat_param_grid = {'iterations': [100, 200, 300], 'depth': [4, 6, 8]}
cat_random_search = RandomizedSearchCV(CatBoostClassifier(), cat_param_grid,
n_iter=5, cv=5)
cat_random_search.fit(X_train, y_train)
cat_random_search.best_params_ # Best parameters for CatBoost
cat_random_search.best_score_ # Best score for CatBoost

# K-means parameter optimization
kmeans_param_grid = {'n_clusters': [3, 5, 7], 'init': ['k-means++', 'random']}
kmeans_grid_search = GridSearchCV(KMeans(), kmeans_param_grid, cv=5)
kmeans_grid_search.fit(X_train)
kmeans_grid_search.best_params_ # Best parameters for K-Means
kmeans_grid_search.best_score_ # Best score for K-Means

# PCA parameter optimization
pca_param_grid = {'n_components': [2, 5, 10]}
pca_grid_search = GridSearchCV(PCA(), pca_param_grid, cv=5)
pca_grid_search.fit(X_train)
pca_grid_search.best_params_ # Best parameters for PCA
pca_grid_search.best_score_ # Best score for PCA

```

Data Visualization

Cheat Sheet



Importing Libraries

```
import matplotlib.pyplot as plt  
import seaborn as sns
```

Basic Line Plot

```
plt.plot(x, y, linestyle = '-', color = 'b')  
plt.xlabel('X-axis label')  
plt.ylabel('Y-axis label')  
plt.title('Title')  
plt.show()
```

Example

```
# Sample dataset  
months = ['Jan', 'Feb', 'Mar', 'Apr', 'May']  
revenue = [10000, 15000, 12000, 18000, 20000]
```

```
# Create the line plot  
plt.plot(months, revenue, marker='o', linestyle='-', color='b')  
plt.xlabel('Months')  
plt.ylabel('Revenue ($)')  
plt.title('Monthly Revenue')  
plt.grid(True)  
plt.show()
```



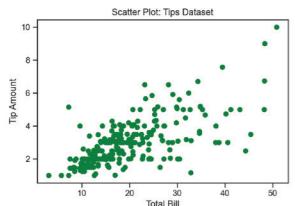
Scatter Plot

```
plt.scatter(x, y, c='color')  
plt.xlabel('X-axis label')  
plt.ylabel('Y-axis label')  
plt.title('Title')  
plt.show()
```

Example

```
# Load the tips dataset  
df = sns.load_dataset("tips")  
  
# Extract the total bill and tip amounts  
total_bill = df["total_bill"].values  
tip_amount = df["tip"].values
```

```
# Create a scatter plot  
plt.scatter(total_bill, tip_amount, c='green')  
plt.xlabel('Total Bill')  
plt.ylabel('Tip Amount')  
plt.title('Scatter Plot: Tips Dataset')  
plt.show()
```



Bar Plot

```
plt.bar(x, y, color = 'b', width = 0.5, height = 0.3)  
plt.xlabel('X-axis label')  
plt.ylabel('Y-axis label')  
plt.title('Title')  
plt.show()
```

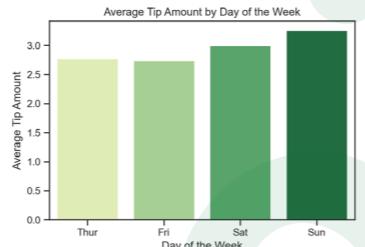
```
# Horizontal bar  
plt.bart(x, y)  
plt.show()
```

Example

```
# Load the 'tip' dataset from seaborn  
tips = sns.load_dataset('tips')
```

```
# Calculate the average tip amount for each day of the week  
avg_tip_by_day = tips.groupby('day')['tip'].mean()
```

```
# Create the bar plot  
sns.barplot(x=avg_tip_by_day.index, y=avg_tip_by_day.values)  
plt.xlabel('Day of the Week')  
plt.ylabel('Average Tip Amount')  
plt.title('Average Tip Amount by Day of the Week')  
plt.show()
```



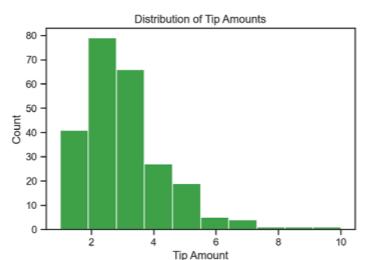
Histogram

```
plt.hist(data, bins=10)  
plt.xlabel('X-axis label')  
plt.ylabel('Y-axis label')  
plt.title('Title')  
plt.show()
```

Example

```
# Load the 'tip' dataset from seaborn  
tips = sns.load_dataset('tips')
```

```
# Create the histogram  
sns.histplot(data=tips, x='tip', bins=10, color='green')  
plt.xlabel('Tip Amount')  
plt.ylabel('Count')  
plt.title('Distribution of Tip Amounts')  
plt.show()
```



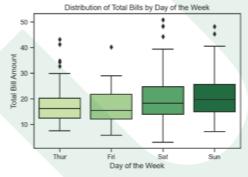
Box Plot

```
sns.boxplot(data, x, y, color)  
plt.xlabel('X-axis label')  
plt.ylabel('Y-axis label')  
plt.title('Title')  
plt.show()
```

Example

```
# Load the 'tip' dataset from seaborn  
tips = sns.load_dataset('tips')
```

```
# Create the box plot with a different color  
sns.boxplot(data=tips, x='day', y='total_bill',  
            color='green')  
plt.xlabel('Day of the Week')  
plt.ylabel('Total Bill Amount')  
plt.title('Distribution of Total Bills by Day of the Week')  
plt.show()
```



Heatmap

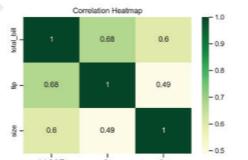
```
sns.heatmap(data, annot=True)  
plt.xlabel('X-axis label')  
plt.ylabel('Y-axis label')  
plt.title('Title')  
plt.show()
```

Example

```
# Load the 'tip' dataset from seaborn  
tips = sns.load_dataset('tips')
```

```
# Calculate the correlation matrix  
correlation_matrix = tips.corr()
```

```
# Create the heatmap  
sns.heatmap(correlation_matrix, annot=True, cmap='YlGn')  
plt.title('Correlation Heatmap')  
plt.show()
```



Pie Chart

```
plt.pie(y, labels=labels, colors=[colors])  
plt.title('Title')  
plt.show()
```

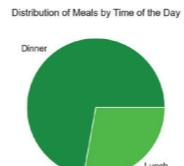
Example

```
# Load the 'tip' dataset from seaborn  
tips = sns.load_dataset('tips')
```

```
# Calculate the count of meals for each time of the day  
meal_counts = tips['time'].value_counts()
```

```
# Specify custom labels and colors  
labels = meal_counts.index  
colors = ['darkgreen', 'aquamarine']
```

```
# Create the pie chart  
plt.pie(meal_counts, labels=labels, colors=colors)  
plt.title('Distribution of Meals by Time of the Day')  
plt.show()
```



Multiple Subplots

```
fig, axs = plt.subplots(2, 2)  
axs[0, 0].plot(x1, y1)  
axs[0, 1].scatter(x2, y2)  
axs[1, 0].bar(x3, height3)  
axs[1, 1].hist(data4, bins=10)  
plt.show()
```

Example

```
# Load the 'tips' dataset from seaborn  
tips = sns.load_dataset('tips')
```

```
# Prepare the data for plotting  
x1 = tips['total_bill']  
y1 = tips['tip']  
x2 = tips['size']  
y2 = tips['tip']  
x3 = tips['day']  
height3 = tips['total_bill']  
data4 = tips['total_bill']
```

```
# Create the subplots  
fig, axs = plt.subplots(2, 2)
```

```
# Plot 1: Line plot  
axs[0, 0].plot(x1, y1)  
axs[0, 0].set_xlabel('Total Bill')  
axs[0, 0].set_ylabel('Tip')  
axs[0, 0].set_title('Tip Amount vs Total Bill')
```

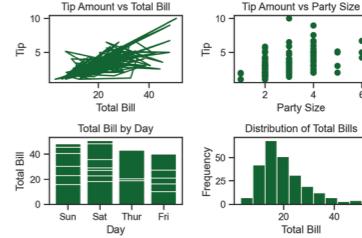
```
# Plot 2: Scatter plot  
axs[0, 1].scatter(x2, y2)  
axs[0, 1].set_xlabel('Party Size')  
axs[0, 1].set_ylabel('Tip')  
axs[0, 1].set_title('Tip Amount vs Party Size')
```

```
# Plot 3: Bar plot  
axs[1, 0].bar(x3, height3)  
axs[1, 0].set_xlabel('Day')  
axs[1, 0].set_ylabel('Total Bill')  
axs[1, 0].set_title('Total Bill by Day')
```

```
# Plot 4: Histogram  
axs[1, 1].hist(data4, bins=10)  
axs[1, 1].set_xlabel('Total Bill')  
axs[1, 1].set_ylabel('Frequency')  
axs[1, 1].set_title('Distribution of Total Bills')
```

```
# Adjust the layout and spacing  
fig.tight_layout()
```

```
# Display the plot  
plt.show()
```



Customizing Plots

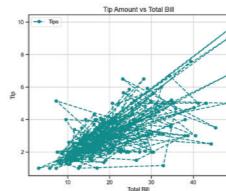
```
plt.figure(figsize=(8, 6))  
plt.plot(x, y, color='red', linestyle='--', linewidth=2,  
         marker='o', markersize=6)  
plt.xlabel('X-axis label', fontsize=12)  
plt.ylabel('Y-axis label', fontsize=12)  
plt.title('Title', fontsize=14)  
plt.legend(['Legend'])  
plt.grid(True)  
plt.show()
```

Example

```
# Load the 'tip' dataset from seaborn  
tips = sns.load_dataset('tips')
```

```
# Prepare the data for plotting  
x = tips['total_bill']  
y = tips['tip']
```

```
# Create the plot  
plt.figure(figsize=(8, 6))  
plt.plot(x, y, color='red', linestyle='--', linewidth=2,  
         marker='o', markersize=6)  
plt.xlabel('Total Bill', fontsize=12)  
plt.ylabel('Tip', fontsize=12)  
plt.title('Tip Amount vs Total Bill', fontsize=14)  
plt.legend(['tips'])  
plt.grid(True)  
plt.show()
```



Probability Cheat Sheet

Probability Basics

Experiment : A process that results in an outcome.

Sample Space (S) : The set of all possible outcomes of an experiment.

Event (E) : A subset of the sample space.

Probability (P) : A measure of the likelihood of an event occurring.

Probability Axioms

Non-Negativity : For all events A , $P(A) \geq 0$.

Additivity : For mutually exclusive events A and B ,

$$P(A \cup B) = P(A) + P(B).$$

Normalization : The probability of the entire sample space S is 1:

$$P(S) = 1.$$

Discrete Uniform Law (Classical Probability) : For equally likely outcomes,

$$P(E) = \frac{\text{Number of favorable outcomes}}{\text{Total number of outcomes}}$$

Relative Frequency Probability :

$$P(E) = \frac{\text{Frequency of } E \text{ occurring}}{\text{Total number of trials}}$$

Conditional Probability $P(E|F)$: Probability of E given that F has occurred. It is calculated as:

$$P(E|F) = \frac{P(E \cap F)}{P(F)} \quad \text{where } P(F) \neq 0.$$

Bayes' Theorem

Bayes' Theorem helps us update our initial beliefs (prior probabilities) based on new information (likelihood) to arrive at a more accurate or informed estimate of the probability of an event occurring (posterior probability). Mathematically, Bayes' Theorem is stated as follows:

$$P(F|E) = \frac{P(F|E) \cdot P(E)}{P(F)}$$

Random Variables

Random Variable : A random variable assigns a numerical value to each outcome in a sample space. It can be discrete or continuous.

Expected Value : The expected value of a random variable represents the average or mean value it takes over all possible outcomes. It provides a measure of the central tendency of the distribution.

Variance : The variance of a random variable measures the extent to which the values of the variable deviate from its expected value. It quantifies the spread or dispersion of the distribution.

Discrete Random Variables

For a discrete random variable X , which takes on distinct values from a finite or countable set, we have the following key concepts:

Probability Mass Function (PMF): The PMF $P(X=x)$ of a discrete random variable X gives the probability that X takes on the value x . It describes the distribution of probabilities across all possible values of X .

Expected value of a discrete random variable X is calculated as the sum of each value x weighted by its probability $P(X=x)$:

$$\bullet \quad E(X) = \sum_x x \cdot P(X=x)$$

Variance of a discrete random variable X is calculated as the average of the squared differences between each value x and the expected value $E(X)$:

$$\bullet \quad \text{Var}(X) = E(X^2) - (E(X))^2$$

Continuous Random Variables

For a continuous random variable X , which can take on any value within a certain range, we have the following key concepts:

Probability Density Function (PDF): The PDF $f(x)$ of a continuous random variable X provides the relative likelihood that X falls within a specific interval. Unlike discrete random variables, the probability that X takes on a specific value is generally zero. Instead, the probability is associated with intervals.

Expected value of a continuous random variable X is the integral of x multiplied by the $f(x)$ over the entire range of X :

$$\bullet \quad E(X) = \int_{-\infty}^{\infty} x \cdot f(x) dx$$

Variance of a continuous random variable X is calculated similarly:

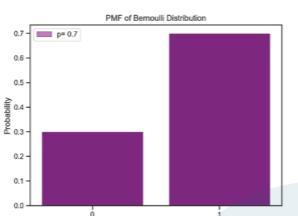
$$\bullet \quad \text{Var}(X) = E(X^2) - (E(X))^2$$

Common Probability Distributions

Discrete Distributions

Bernoulli Distribution

The Bernoulli distribution models a random experiment with two possible outcomes, usually labeled as "success" and "failure".

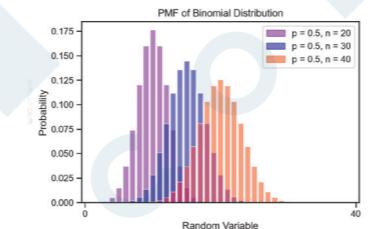


- PMF : $P(X=x) = p^x \cdot (1-p)^{1-x}$
- Parameters : p (probability of success)
- Expected Value : $E(X) = p$
- Variance : $\text{Var}(X) = p \cdot (1-p)$

Binomial Distribution

The binomial distribution describes the number of successes (S) in a fixed number of independent Bernoulli trials. Let us say you conduct " n " trials and you're interested in how many times you will get a success. The binomial distribution provides the probabilities of getting 0, 1, 2, ..., n successes.

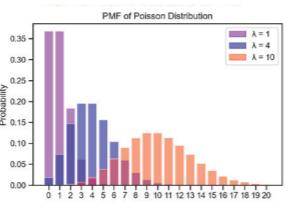
- PMF : $P(X=k) = C(n,k) \cdot p^k \cdot (1-p)^{n-k}$
- Parameters : n (number of trials), p (probability of success)
- Expected Value : $E(X) = n \cdot p$
- Variance : $\text{Var}(X) = n \cdot p \cdot (1-p)$



Poisson Distribution

The Poisson distribution is a probability distribution that describes the number of events that occur in a fixed interval of time or space, given a certain average rate of occurrence.

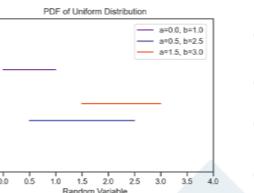
- PMF : $P(X=k) = \frac{\lambda^k \cdot e^{-\lambda}}{k!}$
- Parameters : λ (average rate of events)
- Expected Value : $E(X) = \lambda$
- Variance : $\text{Var}(X) = \lambda$



Continuous Distributions

Uniform Distribution

The uniform distribution is a probability distribution that describes a situation where all outcomes in a given range are equally likely. In other words, each possible value within the range has the same probability of occurring.

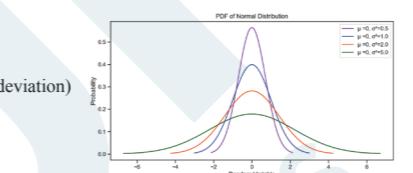


- PMF : $f(x) = \frac{1}{b-a}$ for $a \leq x \leq b$
- Parameters : a (lower bound), b (upper bound)
- Expected Value : $E(X) = a + b$
- Variance : $\text{Var}(X) = \frac{(b-a)^2}{12}$

Normal (Gaussian) Distribution

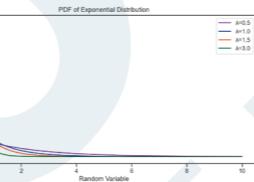
The normal distribution, also known as the Gaussian distribution or bell curve, is a fundamental concept in statistics and probability theory. It describes a symmetrical probability distribution that is characterized by its bell-shaped curve.

- PDF : $f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$
- Parameters : μ (mean), σ (standard deviation)
- Expected Value : $E(X) = \mu$
- Variance : $\text{Var}(X) = \sigma^2$



Exponential Distribution

The exponential distribution models the time between successive events in a process where events occur randomly and independently at a constant average rate.

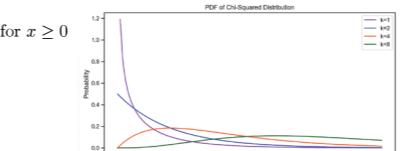


- PDF : $f(x) = \lambda e^{-\lambda x}$ for $x \geq 0$
- Parameters : λ (rate parameter)
- Expected Value : $E(X) = \frac{1}{\lambda}$
- Variance : $\text{Var}(X) = \frac{1}{\lambda^2}$

Chi-Squared Distribution

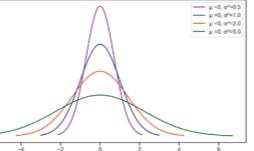
The chi-squared distribution emerges from the sum of the squares of independent standard normal random variables.

- PDF : $f(x) = \frac{1}{2^{k/2}\Gamma(k/2)} x^{k/2-1} e^{-x/2}$ for $x \geq 0$
- Parameters : k (degrees of freedom)
- Expected Value : $E(X) = k$
- Variance : $\text{Var}(X) = 2k$



Student's t-Distribution

The t-distribution models the variability of sample means when drawing small samples from a population. It is characterized by its bell-shaped curve, similar to the normal distribution.

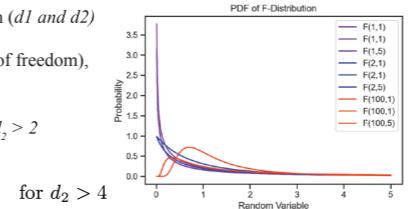


- PDF : Depends on degrees of freedom (v)
- Parameters : v (degrees of freedom)
- Expected Value : $E(X) = 0$ for $v > 1$
- Variance : $\text{Var}(X) = \frac{v}{v-2}$

F-Distribution

The F-distribution models the ratio of two independent chi-squared distributions divided by their respective degrees of freedom.

- PDF : Depends on degrees of freedom (d_1 and d_2)
- Parameters : d_1 (numerator degrees of freedom), d_2 (denominator degrees of freedom)
- Expected Value : $E(X) = \frac{d_2}{d_2-2}$ for $d_2 > 2$
- Variance : $\text{Var}(X) = \frac{2d_2^2(d_1+d_2-2)}{d_1(d_2-2)^2(d_2-4)}$ for $d_2 > 4$



Fundamental Results

The Law of Large Numbers

Weak Law of Large Numbers (WLLN) : The sample mean of a large number of independent and identically distributed random variables approaches the expected value.

Strong Law of Large Numbers (SLLN) : The sample mean converges almost surely to the expected value.

```
# Import necessary libraries
import random
import matplotlib.pyplot as plt

# Function to simulate rolling a fair six-sided die
def roll_die():
    return random.randint(1, 6)

# Number of trials to perform for each sample size
num_trials = 10000

# List of sample sizes to investigate
sample_sizes = [10, 50, 100, 500, 1000, 5000]

# Iterate through each sample size
for sample_size in sample_sizes:
    trial_means = []

    # Perform num_trials trials for the current sample size
    for _ in range(num_trials):
        # Simulate rolling the die 'sample_size' times and calculate the mean
        sample = [roll_die() for _ in range(sample_size)]
        mean = sum(sample) / sample_size
        trial_means.append(mean)

    # Create a histogram of trial means with 20 bins and add it to the plot
    plt.hist(trial_means, bins=20, alpha=0.5)
    label=f'Sample size: {sample_size}'

    # Label the x and y axes and set the title of the plot
    plt.xlabel('Sample Mean')
    plt.ylabel('Frequency')
    plt.title('Law of Large Numbers')

    plt.legend()

plt.show()
```

The Central Limit Theorem

The distribution of the sum (or average) of a large number of independent, identically distributed random variables approaches a normal distribution.

```
# Import necessary libraries
import numpy as np
import matplotlib.pyplot as plt

# Define the size of the population and generate a random population with an exponential distribution
population_size = 100000
population = np.random.exponential(scale=2, size=population_size)

# Define a list of sample sizes that we want to explore
sample_sizes = [10, 50, 100, 500]

# Iterate through each sample size
for sample_size in sample_sizes:
    # Generate 1000 random samples of the specified size from the population
    sample_means = [np.mean(np.random.choice(population, size=sample_size)) for _ in range(1000)]

    # Create a histogram of the sample means with 30 bins and add it to the plot
    plt.hist(sample_means, bins=30, alpha=0.5)
    label=f'Sample size: {sample_size}'

    # Add labels and a title to the plot
    plt.xlabel('Sample Mean')
    plt.ylabel('Frequency')
    plt.title('Central Limit Theorem')

    plt.legend()
```