# **Comprehensive Overview of Data Analytics: Tools, Techniques, and Models**

List of subsections within data analytics, including relevant tools, models, techniques, comparisons, and evaluation methods

- 1. Data Collection and Acquisition
  - Techniques: Web Scraping, API Integration, Survey Collection, Sensor Data Acquisition
  - Tools: BeautifulSoup, Scrapy, Requests, Selenium, Postman
  - Models: Data Collection Frameworks, ETL Models

## **Data Analytics: Core Aspects Breakdown**

The core aspects of data analytics, exploratory data analysis (EDA), regressions, machine learning, and deep learning, along with relevant tools, techniques, models, comparisons, and evaluation methods.

### **Data Analytics**

- 1. Data Collection and Acquisition
  - Techniques: Web Scraping, API Integration, Survey Collection, Sensor Data
  - Tools: BeautifulSoup, Scrapy, Requests, Selenium, Postman

# **Data Analyis Models: Table**

Table of models across data analytics, EDA, regressions, machine learning, and deep learning, along with their relationships and differences:

## **Key Differences and Relationships:**

- Summary Statistics vs. Data Visualization: Summary statistics provide basic metrics, while data visualization presents these metrics in graphical form for
- Linear Regression vs. Logistic Regression: Both are regression models, but linear regression is for continuous outcomes, and logistic regression is for binary

Table of models across data analytics, EDA, regressions, machine learning, and deep learning, along with their relationships and differences:		
Model/ Technique	Domain	Related M Technique
Summary Statistics	EDA	Descriptive Statistics
Data Visualization	EDA, Data Analytics	Interactive Dashboard
Linear Regression	Regression	Multiple Li Regression Polynomia Regression

# Most Important Models Notes + Table

Most important models in each area of data analytics, EDA, regressions, machine learning, and deep learning:

### **Data Analytics**

- 1. Descriptive Analytics
  - Summary Statistics Models: Mean, median, mode, standard deviation,
  - Descriptive Visualization Models: Histograms, bar charts, pie charts.
- 2. Predictive Analytics

## **Al Model Comparison Chart**

Model/Tool	Туре	Descript
GPT-40	Language Model	Advanced of GPT-4
DALL+E 3	Image Generation Model	Generater images for description
GPT-3.5	Language Model	Predeces GPT-4
GPT-4	Language Model	Latest ver GPT from
Claude 3 Halku	Language Model	Version or optimized creative to

# **Python Cheat Sheet: Essential Commands and Functions**

# **General Syntax**

Comments

# **Data Types**

- Integers: int
- Floating Point: float
- String: str

# **R Cheat Sheet: Essential Commands and Functions**

## **General Syntax** Comments

# **Data Types**

- Numeric: numeric
- Integer: integer

# **Functions** Basics

**SQL Cheat Sheet: Essential Commands and** 

 Comments **Data Definition Language (DDL)** 

# Create Database

- Drop Database
- Create Table



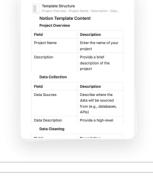
# **Data Science Project Guide Notion Template Template Structure**

Project Overview · Project Name · Description · Objectives · Outcomes · Data Collection · Data Source..

**Notion Template Content** 

**Project Overview Data Collection** 

**Data Cleaning Data Analysis** 



# Python Code Reference Guide for Data Science

**Extras** 

 $\textbf{Table of Contents} \cdot \textbf{Introduction} \cdot \textbf{Python Basics} \cdot \textbf{Data Types} \cdot \textbf{Control Structures} \cdot \textbf{Functions} \cdot \textbf{Librari...}$ What is Natural Language Processing? · Natural Language Processing (NLP) is a sub-field of artificial...

What is Regex? · Regex, short for Regular Expression, is a sequence of characters that forms a search... Basic Stats, Analysis, Data Science Overview

## Comprehensive Overview of Data Analytics: Tools, Techniques, and Models

List of subsections within data analytics, including relevant tools, models, techniques, comparisons, and evaluation methods

## 1. Data Collection and Acquisition

- **Techniques**: Web Scraping, API Integration, Survey Collection, Sensor Data Acquisition
- Tools: BeautifulSoup, Scrapy, Requests, Selenium, Postman
- Models: Data Collection Frameworks, ETL Models
- Comparisons: Batch vs. Real-time Data Collection

# 2. Data Storage and Management

- Techniques: Data Warehousing, Data Lakes, Distributed Storage
- Tools: SQL Databases (MySQL, PostgreSQL), NoSQL Databases (MongoDB, Cassandra), Data Lakes (Amazon S3, Azure Data Lake)
- Models: Relational Models, Document Models, Columnar Models
- Comparisons: SQL vs. NoSQL, On-Premise vs. Cloud Storage

# 3. Data Cleaning and Preparation

- **Techniques**: Data Wrangling, Handling Missing Values, Data Normalization
- Tools: Pandas, OpenRefine, Dask, Alteryx, Trifacta
- Models: Data Cleaning Pipelines, Data Transformation Frameworks
- Comparisons: Manual Cleaning vs. Automated Cleaning

# 4. Descriptive Analytics

- **Techniques**: Descriptive Statistics, Data Visualization
- Tools: Excel, Tableau, Power BI, Matplotlib, Seaborn
- Models: Summary Statistics Models, EDA Frameworks
- Comparisons: Static Reports vs. Interactive Dashboards

# 5. Diagnostic Analytics

- Techniques: Root Cause Analysis, Anomaly Detection, Correlation Analysis
- Tools: Splunk, RapidMiner, SAS, Python (Scikit-learn), R
- Models: Diagnostic Models, Anomaly Detection Algorithms
- Comparisons: Traditional Methods vs. Machine Learning Methods

# 6. Predictive Analytics

- Techniques: Regression Analysis, Time Series Forecasting, Machine Learning
- Tools: Scikit-learn, TensorFlow, Keras, PyTorch, Prophet
- Models: Linear Regression, Decision Trees, Random Forests, ARIMA, LSTM
- Comparisons: Traditional Statistical Models vs. Machine Learning Models

# 7. Prescriptive Analytics

- Techniques: Optimization, Decision Analysis, Simulation
- Tools: Gurobi, CPLEX, AnyLogic, Arena, Simul8
- Models: Linear Programming, Integer Programming, Monte Carlo Simulation
- Comparisons: Deterministic Models vs. Stochastic Models

# 8. Advanced Analytics

- Techniques: Deep Learning, Natural Language Processing (NLP), Computer Vision
- Tools: TensorFlow, PyTorch, spaCy, OpenCV
- Models: Convolutional Neural Networks (CNN), Recurrent Neural Networks (RNN), Transformers
- Comparisons: Rule-based Systems vs. Al-based Systems

# 9. Data Visualization

- Techniques: Graphical Representations, Interactive Dashboards, Geospatial Visualization
   Tools: Tableau, Power BI, QlikView, D3.js, Plotly
- Models: Visualization Frameworks, Geospatial Models
- Comparisons: Static Visualizations vs. Interactive Visualizations
- 10. Statistical Analysis

# • **Techniques**: Inferential Statistics, Hypothesis Testing, Bayesian Statistics

- Tools: R, Python (SciPy, Statsmodels), SAS, SPSS
   Models: T-tests, ANOVA, Bayesian Networks
- Comparisons: Frequentist vs. Bayesian Approaches
- 11. Data EngineeringTechniques: ETL Processes, Data Pipelines, Workflow Orchestration
- Tools: Apache Airflow, Luigi, NiFi, Talend, Informatica
   Models: ETL Models, Data Pipeline Frameworks
- Comparisons: Batch Processing vs. Stream Processing
- 12. Big Data Technologies

# • Techniques: Distributed Computing, Real-Time Data Processing

- Tools: Hadoop, Spark, Flink, Storm, Kafka
   Models: MapReduce, Distributed Data Processing Models
- Comparisons: Hadoop vs. Spark, Batch Processing vs. Stream Processing
- 13. Time Series Analysis
- Techniques: ARIMA Models, Seasonal Decomposition, Forecasting
   Tools: R (forecast package), Python (statsmodels, Prophet)
- Models: ARIMA, SARIMA, Exponential Smoothing, LSTM
   Comparisons: Traditional Time Series Models vs. Neural Network Models
- 14. Optimization Techniques Techniques: Linear and Nonlinear Programming, Stochastic Optimization
- o Tools: Gurobi, CPLEX, OR-Tools
- Comparisons: Exact Methods vs. Heuristic Methods
   Model Evaluation and Validation

Models: Linear Programming, Genetic Algorithms, Simulated Annealing

• **Techniques**: Cross-Validation, Bootstrap Methods, Performance Metrics

Tools: Scikit-learn, MLflow, TensorBoard

Comparisons: Holdout Validation vs. Cross-Validation

Models: Confusion Matrix, ROC Curve, AUC, F1 Score

- 16. Automated AnalyticsTechniques: Automated Machine Learning (AutoML).
- Techniques: Automated Machine Learning (AutoML), Recommendation Engines
   Techniques: AutoML DataBahat, TROT
- Tools: <u>H2O.ai</u>, AutoML, DataRobot, TPOT
   Models: Automated Feature Engineering, Hyperparameter Tuning Models
- Comparisons: Manual Model Building vs. AutoML

# 17. Operationalization and Deployment Output Techniques: Model Serving, CVCD for

- Techniques: Model Serving, CI/CD for ML, A/B Testing
   Tools: Docker, Kubernetes, MLflow, TFX
- Models: Deployment Pipelines, Model Monitoring Frameworks
   Deployment Papel time Papel time
- Comparisons: Batch Deployment vs. Real-time Deployment
   18. Ethics and Data Privacy
- **Techniques**: Data Governance, Privacy-Preserving Data Analysis
- Tools: Differential Privacy Libraries, Data Anonymization Tools
- Models: Fairness Metrics, Bias Detection Models
- Comparisons: GDPR vs. CCPA Compliance, Ethical AI vs. Traditional AI

This list provides a detailed overview of the various subsections of data analytics, including relevant tools, models, techniques, comparisons, and evaluation methods.

# **Data Analytics: Core Aspects Breakdown**

The core aspects of data analytics, exploratory data analysis (EDA), regressions, machine learning, and deep learning, along with relevant tools, techniques, models, comparisons, and evaluation methods.

## **Data Analytics**

# 1. Data Collection and Acquisition

- Techniques: Web Scraping, API Integration, Survey Collection, Sensor Data Acquisition
- Tools: BeautifulSoup, Scrapy, Requests, Selenium, Postman 2. Data Cleaning and Preparation

- Techniques: Data Wrangling, Handling Missing Values, Data Normalization • Tools: Pandas, OpenRefine, Dask, Alteryx, Trifacta
- 3. Descriptive Analytics

- Techniques: Descriptive Statistics, Data Visualization
- Tools: Excel, Tableau, Power BI, Matplotlib, Seaborn Models: Summary Statistics Models, EDA Frameworks
- 4. Diagnostic Analytics

## Techniques: Root Cause Analysis, Anomaly Detection, Correlation Analysis

- Tools: Splunk, RapidMiner, SAS, Python (Scikit-learn), R
- Models: Diagnostic Models, Anomaly Detection Algorithms
- 5. Predictive Analytics

## • Techniques: Regression Analysis, Time Series Forecasting, Machine Learning

- Tools: Scikit-learn, TensorFlow, Keras, PyTorch, Prophet
- Models: Linear Regression, Decision Trees, Random Forests, ARIMA, LSTM
- 6. Prescriptive Analytics

## • Techniques: Optimization, Decision Analysis, Simulation • Tools: Gurobi, CPLEX, AnyLogic, Arena, Simul8

- Models: Linear Programming, Integer Programming, Monte Carlo Simulation

# 1. Techniques

**Exploratory Data Analysis (EDA)** 

# • Data Summarization: Using descriptive statistics to summarize data.

- Data Visualization: Plotting data to identify patterns and outliers.
- Data Profiling: Examining data distributions and relationships.
- Outlier Detection: Identifying anomalies in data.
- 2. Tools • Pandas: Data manipulation and summary statistics.

- Matplotlib: Basic plotting.
- Seaborn: Advanced statistical visualizations.
- Plotly: Interactive visualizations.
- Tableau/Power BI: Business intelligence and interactive dashboards.
- 3. Models • Summary Statistics Models: Mean, median, mode, standard deviation.

- Distribution Analysis: Histograms, box plots, density plots.
- Correlation Analysis: Pearson correlation, Spearman rank correlation.
- 4. Comparisons

# Static Visualizations vs. Interactive Dashboards: Static plots for quick insights vs. interactive

dashboards for detailed exploration. Regressions

# 1. Techniques

# • Simple Linear Regression: Modeling relationship between two variables.

- Multiple Linear Regression: Modeling relationship between one dependent variable and multiple independent variables.
- Polynomial Regression: Modeling non-linear relationships.
- Logistic Regression: Modeling binary outcomes. Ridge and Lasso Regression: Regularization techniques to prevent overfitting.
- 2. Tools
  - Scikit-learn: Comprehensive library for regression models.

# • Statsmodels: Statistical models and hypothesis tests.

- R: Built-in functions and packages like 1m and g1m.
- SAS/SPSS: Commercial tools for regression analysis.
- 3. Models Linear Regression: (y = \beta\_0 + \beta\_1 x + \epsilon)
  - Multiple Linear Regression: (y = \beta\_0 + \beta\_1 x\_1 + \beta\_2 x\_2 + \ldots + \beta\_n x\_n + \epsilon)
- 4. Evaluation R-squared: Proportion of variance explained by the model.

Mean Squared Error (MSE): Average squared difference between observed and predicted

Logistic Regression: (\log\left(\frac{p}{1-p}\right) = \beta\_0 + \beta\_1 x\_1 + \ldots + \beta\_n x\_n)

AIC/BIC: Model selection criteria based on goodness of fit and complexity.

# 5. Comparisons

**Machine Learning** 

1. Techniques

- Linear vs. Logistic Regression: Continuous outcome vs. binary outcome. • Ridge vs. Lasso Regression: L2 regularization vs. L1 regularization.
- Supervised Learning: Training models with labeled data. • Unsupervised Learning: Finding patterns in unlabeled data.

• Semi-supervised Learning: Using a mix of labeled and unlabeled data.

- Reinforcement Learning: Learning through rewards and punishments.
- 2. Tools • Scikit-learn: Comprehensive library for various ML algorithms.

Neighbors (k-NN).

classification.

- TensorFlow/Keras: Deep learning frameworks. • PyTorch: Deep learning library with a dynamic computation graph.
- 3. Models Supervised: Decision Trees, Random Forests, Support Vector Machines (SVM), k-Nearest

XGBoost/LightGBM: Gradient boosting libraries for efficient and accurate predictions.

Unsupervised: K-means Clustering, Hierarchical Clustering, Principal Component Analysis

ROC/AUC: Receiver Operating Characteristic curve and Area Under the Curve for binary

(PCA). • Reinforcement: Q-Learning, Deep Q-Networks (DQN).

4. Evaluation

 Accuracy: Proportion of correctly predicted instances. Precision/Recall/F1 Score: Metrics for classification problems.

• Confusion Matrix: Matrix to evaluate classification performance.

5. Comparisons • Supervised vs. Unsupervised Learning: Labeled vs. unlabeled data.

• Decision Trees vs. Random Forests: Single tree vs. ensemble of trees. • SVM vs. k-NN: Hyperplane separation vs. distance-based classification.

• Cross-Validation: Assessing model performance using different data splits.

**Deep Learning** 

Neural Networks: Feedforward, Convolutional (CNN), Recurrent (RNN).

# Transfer Learning: Using pre-trained models for new tasks. • Generative Adversarial Networks (GANs): Generating new data samples.

1. Techniques

# • Reinforcement Learning: Deep Q-Networks (DQN), Policy Gradients.

- 2. Tools • TensorFlow/Keras: High-level APIs for building and training deep learning models.
  - Theano: Deep learning library (less commonly used now). MXNet: Scalable deep learning framework.

• PyTorch: Deep learning library with dynamic computation graph.

3. Models • Feedforward Neural Networks (FNN): Basic neural networks for structured data.

• Convolutional Neural Networks (CNN): Image and video processing.

- Long Short-Term Memory (LSTM): Improved RNNs for long-term dependencies.
- GANs: Generative models for creating new data samples. 4. Evaluation
  - Metrics: Accuracy, Precision, Recall, F1 Score. • Training Curves: Monitoring loss and accuracy over epochs.
  - Hyperparameter Tuning: Grid search, random search, Bayesian optimization.
- 5. Comparisons • CNN vs. RNN: Spatial data processing vs. sequential data processing.
  - LSTM vs. Vanilla RNN: Handling long-term dependencies vs. basic sequential data.

• GANs vs. Autoencoders: Generative models vs. data compression and reconstruction.

• Recurrent Neural Networks (RNN): Sequential data processing (e.g., time series, text).

• Loss Functions: Cross-entropy for classification, Mean Squared Error for regression.

## **Data Analyis Models: Table**

Table of models across data analytics, EDA, regressions, machine learning, and deep learning, along with

Model/Technique	Domain	Related Models/ Techniques	Key Differences
Summary Statistics	EDA	Descriptive Statistics	Focuses on basic summary metrics like mean, median, and mode
Data Visualization	EDA, Data Analytics	Interactive Dashboards, EDA	Visualization tools like Tableau/Power BI vs. static visualizations in Matplotlib/Seaborn
Linear Regression	Regression	Multiple Linear Regression, Polynomial Regression	Models linear relationships; simpler than multiple or polynomial regression
Logistic Regression	Regression	Linear Regression, Classification Models	Used for binary outcomes vs. continuous outcomes in linear regression
Decision Trees	Machine Learning	Random Forests, Gradient Boosting	Single tree model vs. ensemble methods in Random Forests/ Gradient Boosting
Random Forests	Machine Learning	Decision Trees, Gradient Boosting	Ensemble of multiple decision trees for improved accuracy
Support Vector Machines (SVM)	Machine Learning	Logistic Regression, k- Nearest Neighbors	Hyperplane separatio vs. probability-based logistic regression and distance-based k-NN
k-Nearest Neighbors (k-NN)	Machine Learning	SVM, Decision Trees	Non-parametric method based on distance; differs from SVM's hyperplane and tree's split-based methods
k-Means Clustering	Machine Learning (Unsupervised)	Hierarchical Clustering, DBSCAN	Partitions data into k clusters; differs in methodology from hierarchical and density-based clustering
Principal Component Analysis (PCA)	Machine Learning (Unsupervised)	Linear Discriminant Analysis (LDA)	PCA reduces dimensionality by transforming variable whereas LDA focuses on maximizing class separability
ARIMA	Time Series, Regression	Seasonal Decomposition, LSTM	Traditional time series model vs. advanced neural network methods like LSTM
LSTM (Long Short- Term Memory)	Deep Learning	RNN, GRU	Handles long-term dependencies in sequential data better than vanilla RNN
Feedforward Neural Networks (FNN)	Deep Learning	CNN, RNN	Basic neural networks for structured data vs specialized CNNs for spatial and RNNs for sequential data
Convolutional Neural Networks (CNN)	Deep Learning	RNN, FNN	Best for image and spatial data processin vs. FNNs for structure data and RNNs for sequential data
Generative Adversarial Networks (GANs)	Deep Learning	Autoencoders, Variational Autoencoders (VAE)	GANs generate new data samples, while autoencoders focus of data compression and reconstruction
Transfer Learning	Deep Learning	Fine-Tuning, Feature Extraction	Utilizes pre-trained models for new tasks differs from training models from scratch
Monte Carlo Simulation	Prescriptive Analytics	Optimization Models, Simulation Models	Uses randomness and statistical sampling to model and simulate

constraints.

- **Key Differences and Relationships:** • Summary Statistics vs. Data Visualization: Summary statistics provide basic metrics, while data visualization presents these metrics in graphical form for better understanding.
- **Linear Regression vs. Logistic Regression**: Both are regression models, but linear regression is for continuous outcomes, and logistic regression is for binary outcomes.
- Decision Trees vs. Random Forests: Decision trees are simpler and prone to overfitting, while random forests use multiple trees to improve accuracy and robustness.
- **SVM vs. k-NN**: SVMs use hyperplanes for classification, while k-NN is a distance-based method.
- **k-Means Clustering vs. PCA**: k-Means clusters data, while PCA reduces dimensionality, helping to visualize and understand the structure of high-dimensional data.
- ARIMA vs. LSTM: ARIMA is a traditional time series forecasting method, whereas LSTMs are
- advanced neural networks that handle complex temporal dependencies. • FNN vs. CNN vs. RNN: FNNs are basic neural networks for structured data, CNNs specialize in
- spatial data like images, and RNNs (including LSTM) handle sequential data like time series or text.
- GANs vs. Autoencoders: GANs generate new data similar to the training data, while autoencoders learn efficient representations of data for compression and reconstruction tasks. Monte Carlo Simulation vs. Optimization Models: Monte Carlo simulations use randomness to

understand complex systems, whereas optimization models seek the best solution within

## Most Important Models Notes + Table

Most important models in each area of data analytics, EDA, regressions, machine learning, and deep learning:

## **Data Analytics**

- 1. Descriptive Analytics
  - Summary Statistics Models: Mean, median, mode, standard deviation, variance.
  - Descriptive Visualization Models: Histograms, bar charts, pie charts.

## 2. Predictive Analytics

- Regression Models: Linear Regression, Logistic Regression, Polynomial Regression.
- Time Series Models: ARIMA, Exponential Smoothing.
- Classification Models: Decision Trees, Random Forests, Gradient Boosting Machines (GBM).

## 3. Prescriptive Analytics

- Optimization Models: Linear Programming, Integer Programming.
- Simulation Models: Monte Carlo Simulation.

## **Exploratory Data Analysis (EDA)**

## 1. Data Summarization Models

- **Descriptive Statistics**: Mean, median, mode, range, quartiles.
- **Distribution Analysis**: Histograms, box plots, density plots.

## 2. Data Visualization Models

- Correlation Analysis: Scatter plots, correlation matrices.
- Pattern Detection: Line plots, pair plots.

## 3. Outlier Detection Models

- Box Plots: Visualization for detecting outliers.
- **Z-score Analysis**: Statistical method for detecting outliers.

## Regressions

- 1. Linear Regression
  - Simple Linear Regression: Modeling relationship between two variables.
  - Multiple Linear Regression: Modeling relationship between one dependent variable and multiple independent variables.

## 2. Logistic Regression

- Binary Logistic Regression: Modeling binary outcomes.
- Multinomial Logistic Regression: Modeling outcomes with more than two categories.

## 3. Regularization Techniques

- **Ridge Regression**: L2 regularization to prevent overfitting.
- Lasso Regression: L1 regularization for feature selection and preventing overfitting.

## Machine Learning

## 1. Supervised Learning

- Classification Models: Decision Trees, Random Forests, Support Vector Machines (SVM), k-Nearest Neighbors (k-NN).
- Regression Models: Linear Regression, Ridge Regression, Lasso Regression.

## 2. Unsupervised Learning

- Clustering Models: k-Means Clustering, Hierarchical Clustering, DBSCAN.
- Dimensionality Reduction Models: Principal Component Analysis (PCA), t-Distributed Stochastic Neighbor Embedding (t-SNE).

# 3. Ensemble Learning

- Bagging Models: Random Forests.
- Boosting Models: Gradient Boosting Machines (GBM), AdaBoost, XGBoost.

## 4. Reinforcement Learning • Q-Learning: Model-free reinforcement learning algorithm.

- Deep Q-Networks (DQN): Neural network-based reinforcement learning.
- **Deep Learning**

# 1. Neural Networks

- - Feedforward Neural Networks (FNN): Basic neural networks for structured data. • Deep Neural Networks (DNN): Deeper architectures for more complex patterns.
- 2. Convolutional Neural Networks (CNN)

• Long Short-Term Memory (LSTM): Improved RNNs for long-term dependencies.

# • Basic CNNs: Used for image and video data processing.

- Advanced CNN Architectures: ResNet, VGG, Inception. 3. Recurrent Neural Networks (RNN)
  - Basic RNNs: For sequential data processing.
- Gated Recurrent Units (GRU): Simplified version of LSTM. 4. Generative Models

# • Generative Adversarial Networks (GANs): For generating new data samples.

 Variational Autoencoders (VAE): For generating new data and data compression. 5. Transfer Learning

Pre-trained Models: Using models like BERT, GPT, and ResNet for new tasks.

**Most Important Models** 

**Key Characteristics** 

**Key Comparisons** 

• **Fine-Tuning**: Adapting pre-trained models to specific tasks.

# **Domain**

Descriptive Analytics	Summary Statistics, Histograms	Basic metrics and visualizations for understanding data distributions
Predictive Analytics	Linear Regression, Decision Trees	Models for forecasting and classification
Prescriptive Analytics	Linear Programming, Monte Carlo	Models for optimization and simulating complex systems
EDA	Descriptive Statistics, Box Plots	Initial exploration and understanding of data
Regression	Linear, Logistic, Ridge, Lasso	Modeling relationships between variables, feature selection
Machine Learning	Decision Trees, SVM, k-Means	Supervised and unsupervised learning, clustering, and classification
Deep Learning	CNN, LSTM, GANs	Advanced neural networks for image, sequential data, and generative tasks

application of data analytics, EDA, regressions, machine learning, and deep learning.

These models are essential for their respective domains and are foundational to the practice and

Model/Tool	Туре	Description	Notable Features	Use Case
GPT-4o	Language Model	Advanced version of GPT-4	Enhanced language understanding and generation	Complex text generation, summarization, and analysis
DALL•E 3	Image Generation Model	Generates images from text descriptions	Improved image quality and complexity	Creative design, marketing, and visual storytelling
GPT-3.5	Language Model	Predecessor to GPT-4	Strong language processing capabilities	Customer service, content creation, and translation
GPT-4	Language Model	Latest version of GPT from OpenAl	State-of-the-art language understanding and generation	Advanced research, detailed content generation
Claude 3 Haiku	Language Model	Version of Claude optimized for creative tasks	Specialized in generating haikus	Poetry and creative writing tasks
Claude 3 Opus	Language Model	Version of Claude optimized for creative tasks	Specialized in generating long-form content	Writing novels, articles, and essays
Claude 3 Sonnet	Language Model	Version of Claude optimized for creative tasks	Specialized in generating sonnets	Writing poetry and sonnets
Claude Instant	Language Model	Fast and efficient version of Claude	Quick response times	Real-time customer suppor and quick queries
Claude 2	Language Model	Predecessor to Claude 3	Versatile language processing	General text generation and summarization
Mistral-7b	Language Model	Smaller, efficient language model	Good performance with lower computational requirements	Lightweight applications and mobile devices
_lama-2-70b	Language Model	Large language model from Meta	High capacity for understanding and generating text	Large-scale data analysis and detailed reports
Codellama-34b	Language Model	Code-focused language model	Optimized for programming tasks	Code generation, debugging, and code review
Gemini 1.5 Pro	Language Model	Advanced version of Gemini	Enhanced capabilities	Professional content creation and data analysis
Llama 3 70B	Language Model	Newer version of Llama with 70 billion parameters	High performance and accuracy	Advanced research and high-precision tasks
Llama 3 8B	Language Model	Smaller version of Llama 3	Balanced performance and efficiency	Everyday tasks and small-scale applications
Gemini Pro	Language Model	Professional version of Gemini	Advanced features for professional use	Professional writing and complex data tasks
StepFun	Al Tool	Tool for creating step-by-step instructions	User-friendly interface	Educational content and instructional design
DeepSeek-V2	Al Tool	Tool for deep search and information retrieval	Improved search algorithms	Research and information retrieval
ChatRWKV	Language Model	Chat-focused version of RWKV	Enhanced conversational abilities	Customer interaction and conversational Al
Mixtral 8×7B	Language Model	Mixtral's language model with multiple smaller models	Modular and flexible approach	Modular AI applications and flexible solutions
Gemini 1.0	Language Model	Initial version of Gemini	Foundational model	Basic language tasks and foundational applications
GPT-3.5 Turbo	Language Model	Faster version of GPT-3.5	Quick and efficient processing	Fast content generation and quick queries
GPT-40 pro	Language Model	Professional version of GPT-4	Enhanced for professional applications	Detailed reports and professional documentation
Mixtral 8x22B	Language Model	Mixtral's larger model	Higher capacity and performance	Large-scale analysis and complex problem-solving
Mistral Large	Language Model	Large version of Mistral	Greater capabilities with more resources	Advanced applications and high-demand tasks
GPT-4 Turbo	Language Model	Faster and more efficient version of GPT-4	Optimized for quick tasks	High-speed content generation and real-time analysis

Feature / Model	GPT-4 (OpenAI)	BERT (Google)	T5 (Google)	RoBERTa (Facebook)	GPT-3 (OpenAI)	LLaMA (Meta)
Туре	Transformer- based	Transformer- based	Transformer- based	Transformer- based	Transformer- based	Transformer- based
Release Year	2023	2018	2019	2019	2020	2023
Training Data Size	570GB	3.3 Billion tokens	750GB	160GB	570GB	1.4 Trillion tokens
Parameters	175 Billion	340 Million	11 Billion	355 Million	175 Billion	65 Billion
Architecture	Decoder-only	Encoder-only	Encoder-Decoder	Encoder-only	Decoder-only	Decoder-only
Use Cases	Text generation, NLP	NLP, Question Answering	NLP, Text Summarization	NLP, Text Classification	Text generation, NLP	NLP, Text Generation
Fine-Tuning Capability	Yes	Yes	Yes	Yes	Yes	Yes
Open Source	No	Yes	Yes	Yes	No	Yes
Primary Languages	Multilingual	Multilingual	Multilingual	Multilingual	Multilingual	Multilingual
Strengths	Human-like responses, contextual understanding	Language understanding, low resource use	Flexibility in NLP tasks, language generation	Robustness, accuracy	Creativity, context understanding	Scalability, resource efficiency
Weaknesses	High computational resources required	Limited generation capabilities	High computational resources required	Requires large datasets	High computational resources required	Newer, less tested
Applications	Chatbots, content creation,	Search engines, virtual assistants,	Translation, summarization,	Content moderation,	Chatbots, content creation,	Research, advanced NLP

dialogue systems

sentiment

analysis

virtual assistants

sentiment

analysis

virtual assistants

tasks

# **Python Cheat Sheet: Essential Commands and Functions**

# **General Syntax**

Comments

```
# This is a single-line comment
This is a
multi-line comment
```

# **Data Types**

- Integers: int
- Floating Point: float
- String: str
- List: list
- Tuple: tuple
- Set: set

• Dictionary: dict

# **Variables**

```
x = 5
y = "Hello, World!"
```

# **Basic Operations**

• Arithmetic

```
addition = 5 + 3
subtraction = 5 - 3
multiplication = 5 * 3
division = 5 / 3
floor_division = \frac{5}{3}
modulus = 5 \% 3
exponentiation = 5 ** 3
```

# **Strings**

Concatenation

```
str1 = "Hello"
str2 = "World"
result = str1 + " " + str2
```

Format Strings

```
name = "Alice"
greeting = f"Hello, {name}!"
```

# Lists

Initialization

```
fruits = ["apple", "banana", "cherry"]

    Add Elements
```

fruits.append("orange")

```
    Access Elements
```

first\_fruit = fruits

```
Dictionaries
```

# student = {"name": "John", "age": 25}

Initialization

```
    Access Values

   name = student["name"]
```

**Control Structures** 

if x > 0:

print("Positive") elif x < 0: print("Negative")

Conditional Statements

```
print("Zero")
Loops
 For Loop
    for fruit in fruits:
      print(fruit)
```

count = 5while count > 0:

While Loop

print(count) count -= 1

```
Functions

    Definition

    def greet(name):
       return f"Hello, {name}"
    greeting = greet("Alice")
```

# Libraries • Importing

import numpy as np import pandas as pd

```
Numpy
• Array Initialization
```

• Basic Operations

**Pandas** 

arr = np.array()

```
sum_arr = np.sum(arr)
mean_arr = np.mean(arr)
```

# data = {'Name': ['John', 'Anna', 'Peter', 'Linda'],

• DataFrame Initialization

'Age': } df = pd.DataFrame(data)

```
    Read CSV File

    df = pd.read_csv('file.csv')
```

Matplotlib

plt.show()

```
Plotting
   import matplotlib.pyplot as plt
   plt.plot(, )
```

# **R Cheat Sheet: Essential Commands and Functions**

## **General Syntax**

Comments

```
# This is a single-line comment
```

# **Data Types**

```
• Numeric: numeric
• Integer: integer
```

• Character: character

• Logical: logical

• Vector: vector

• List: list

• Matrix: matrix

• Data Frame: data.frame

# **Variables**

```
x <- 5
y <- "Hello, World!"
```

# **Basic Operations**

Arithmetic

```
addition <-5+3
subtraction <-5-3
multiplication <-5*3
division <- 5 / 3
remainder <- 5 %% 3
exponentiation <- 5 ^ 3
```

# **Strings**

Concatenation

```
str1 <- "Hello"
str2 <- "World"
result <- paste(str1, str2)
```

## **Vectors**

Initialization

```
fruits <- c("apple", "banana", "cherry")</pre>
```

Add Elements

```
fruits <- c(fruits, "orange")</pre>
```

## Lists

• Initialization

```
student <- list(name = "John", age = 25)</pre>
```

Access Values

```
name <- student$name
```

# Initialization

**Matrices** 

matrix\_data <- matrix(1:6, nrow = 2, ncol = 3)</pre>

```
Data Frames
```

# df <- data.frame(</pre>

Creation

Name = c("John", "Anna", "Peter", "Linda"), Age = c(28, 24, 35, 32)

```
    Read CSV File

    df <- read.csv('file.csv')</pre>
```

**Control Structures** 

} else {

if (x > 0) {

print("Positive")  $}$  else if (x < 0) { print("Negative")

Conditional Statements

```
print("Zero")
Loops
  For Loop
      for (fruit in fruits) {
       print(fruit)
```

count <- 5

While Loop

while (count > 0) { print(count)

count <- count - 1

```
Functions
  Definition
    greet <- function(name) {</pre>
     return(paste("Hello,", name))
```

# greeting <- greet("Alice")</pre>

Libraries

```
    Install and Load

    install.packages("ggplot2")
    library(ggplot2)
```

# Data Manipulation

dplyr

```
library(dplyr)
df <- df %>%
 filter(Age > 30) %>%
  arrange(desc(Age))
```

# ggplot2

Plotting

```
ggplot(df, aes(x = Name, y = Age)) +
 geom_bar(stat = "identity")
```

# **SQL Cheat Sheet: Essential Commands and Functions**

## **Basics**

• Comments

```
-- This is a single-line comment
/* This is a
  multi-line comment */
```

## **Data Definition Language (DDL)** Create Database

```
CREATE DATABASE dbname;
```

Drop Database

```
DROP DATABASE dbname;
```

Create Table

```
CREATE TABLE tablename (
  column1 datatype PRIMARY KEY,
  column2 datatype,
   column3 datatype,
);
```

Drop Table

```
DROP TABLE tablename;
```

Alter Table

```
    Add Column

   ALTER TABLE tablename
   ADD columnname datatype;
```

Drop Column

```
ALTER TABLE tablename
DROP COLUMN columnname;
```

**Data Manipulation Language (DML)** 

INSERT INTO tablename (column1, column2, column3,  $\dots$ )

Insert Data

```
VALUES (value1, value2, value3, ...);

    Update Data
```

UPDATE tablename

```
SET column1 = value1, column2 = value2
   WHERE condition;

    Delete Data
```

DELETE FROM tablename WHERE condition;

```
Data Query Language (DQL)
```

### SELECT column1, column2, ... FROM tablename

WHERE condition

Select Data

```
ORDER BY column1, column2, ... ASC|DESC;
Select All Data
  SELECT * FROM tablename;
```

**Data Control Language (DCL)** 

Grant Privileges

Revoke Privileges

```
GRANT ALL PRIVILEGES ON databasename.* TO 'username'@'host';
```

REVOKE ALL PRIVILEGES ON databasename.\* FROM 'username'@'host';

```
SELECT column1, column2
```

FROM tablename WHERE condition;

AND/OR Clause

**Conditional Clauses** 

• WHERE Clause

```
SELECT column1, column2
FROM tablename
WHERE condition1 AND condition2;
```

# SELECT table1.column1, table2.column2 INNER JOIN table2 ON table1.common\_column = table2.common\_column;

FROM table1

Joins

Inner Join

```
Left Join
   SELECT table1.column1, table2.column2
```

SELECT table1.column1, table2.column2

FROM table1 RIGHT JOIN table2 ON table1.common\_column = table2.common\_column;

Right Join

```
Aggregation Functions
• COUNT
   SELECT COUNT(columnname) FROM tablename;
```

LEFT JOIN table2 ON table1.common\_column = table2.common\_column;

SELECT SUM(columnname) FROM tablename;

SUM

```
o AVG
   SELECT AVG(columnname) FROM tablename;
```

SELECT MIN(columnname) FROM tablename;

MIN

```
• MAX
   SELECT MAX(columnname) FROM tablename;
```

**Group By and Having** 

```
Group By
   SELECT column1, COUNT(*)
   FROM tablename
```

GROUP BY column1;

```
    Having Clause

    SELECT column1, COUNT(*)
    FROM tablename
    GROUP BY column1
   HAVING COUNT(*) > 1;
```

# **Data Science Project Guide Notion Template**



## Template Structure

 $Project \ Overview \cdot Project \ Name \cdot Description \cdot Objectives \cdot Outcomes \cdot Data \ Collection \cdot Data \ Sources \cdot Data \ Descri...$ 

## **Notion Template Content**

## **Project Overview**

Field	Description
Project Name	Enter the name of your project
Description	Provide a brief description of the project
Objectives	List the main objectives or goals of the project
Outcomes	Describe the expected outcomes and potential impact of the project

### **Data Collection**

Field	Description
Data Sources	Describe where the data will be sourced from (e.g., databases, APIs)
Data Description	Provide a high-level overview of the dataset(s)
Methods of Collection	Detail the methods used for data collection (e.g., web scraping, surveys)
Tools and Technologies	List the tools and technologies used for data collection (e.g., Python, SQL)

## **Data Cleaning**

Field	Description
Data Inspection	Describe methods used for initial data inspection (e.g., visual inspection, summary statistics)
Handling Missing Values	Explain the approach to handle missing values (e.g., imputation, removal)
Data Transformation	Detail any transformations applied to the data (e.g., converting data types)
Data Normalization	Describe how data normalization is performed, if applicable

# Data Analysis

Field	Description
Exploratory Data Analysis (EDA)	Outline the exploratory techniques used (e.g., histogram, scatter plot)
Statistical Analysis	Detail any statistical analyses performed (e.g., correlation, hypothesis testing)
Machine Learning Models	List and describe the machine learning models used in the project
Tools and Technologies	List the tools and technologies used for data analysis (e.g., pandas, scikit-learn)

# Data Interpretation

Field	Description
Key Findings	Summarize the key findings from the data analysis
Visualizations	Include visualizations that support the findings
Insights and Recommendations	Provide actionable insights and recommendations based on the findings
Limitations	Discuss any limitations of the analysis

## Appendix

Field	Description
References	List any references used in the project
Glossary	Define key terms and concepts used in the project
Additional Resources	Include links to additional resources (e.g., tutorials, documentation)

## **Template Structure**

### 1. Project Overview

- o Project Name
- Description
- o Objectives
- Outcomes

### 2. Data Collection

- Data Sources
- Data Description
- Methods of Collection
- Tools and Technologies

### 3. Data Cleaning

- Data Inspection
- Handling Missing Values
- Data Transformation
- Data Normalization

### 4. Data Analysis

- Exploratory Data Analysis (EDA)
- Statistical Analysis
- Machine Learning Models
- Tools and Technologies

### 5. Data Interpretation

- Key Findings
- Visualizations
- Insights and Recommendations
- Limitations

### 6. Appendix

- o References
- Glossary
- Additional Resources

### ↑ Data Science Guide

### **Extras**



### Python Code Reference Guide for Data Science

Table of Contents · Introduction · Python Basics · Data Types · Control Structures · Functions · Libraries for Data Sci...



### NLP

What is Natural Language Processing? · Natural Language Processing (NLP) is a sub-field of artificial intelligence (Al...



### Regex

What is Regex? · Regex, short for Regular Expression, is a sequence of characters that forms a search pattern. It is...



### Basic Stats, Analysis, Data Science Overview

Basic Statistics · Mean, Median, Mode: Measures of central tendency. · Mean: Average of all data points. · Median: M...



### API

 $Advanced \ API \ Concepts \cdot Authentication: \cdot \ API \ Key: \ A \ simple \ token \ that \ identifies \ the \ calling \ program. \cdot Example: \ Aut...$ 



### Glossary

Algorithms · An algorithm is a set of instructions we give a computer so it can take values and manipulate them into...

## **Python Code Reference Guide for Data Science**

## **Table of Contents**

- 1. Introduction
- 2. Python Basics
  - 1. Data Types
  - 2. Control Structures
  - 3. Functions
- 3. Libraries for Data Science
  - 1. NumPy
  - 2. Pandas
  - 3. Matplotlib
  - 4. Seaborn
  - 5. Scikit-Learn
  - 6. SciPy
  - 7. Statsmodels
- 4. Data Manipulation
  - 1. Loading Data
  - 2. Cleaning Data
  - 3. Transforming Data
- 5. Data Visualization
- 6. Machine Learning
  - 1. Supervised Learning
  - 2. Unsupervised Learning
  - 3. Model Evaluation and Selection
- 7. Advanced Topics
  - 1. Natural Language Processing
  - 2. Deep Learning
  - 3. <u>Time Series Analysis</u>
- 8. Conclusion
- 9. References

# Introduction

Python is a versatile programming language widely used in data science due to its simplicity, readability, and the extensive ecosystem of libraries and tools. This comprehensive guide covers essential Python code and libraries used in data science, from basic data manipulation to advanced machine learning techniques.

# **Python Basics**

# **Data Types**

```
# Integers
a = 10
# Floats
b = 20.5
# Strings
c = "Hello, World!"
# Lists
d = [1, 2, 3, 4, 5]
# Tuples
e = (1, 2, 3, 4, 5)
# Dictionaries
f = {"name": "Alice", "age": 25}
# Sets
g = \{1, 2, 3, 4, 5\}
```

# if a > b:

**Control Structures** 

# If-else statement

```
print("a is greater than b")
 else:
    print("a is not greater than b")
 # For loop
 for i in range(5):
   print(i)
 # While loop
 i = 0
 while i < 5:
   print(i)
    i += 1
Functions
```

# return x + y

def add(x, y):

```
result = add(5, 3)
 print(result)
Libraries for Data Science
NumPy
```

# import numpy as np

# Statistical operations mean = np.mean(arr)

df = pd.DataFrame(data)

# Reading data from CSV df = pd.read\_csv('data.csv')

# DataFrame operations df['Age'] = df['Age'] + 1

import matplotlib.pyplot as plt

plt.plot([1, 2, 3, 4], [1, 4, 9, 16])

## # Creating arrays arr = np.array([1, 2, 3, 4, 5])

arr2 = arr \* 2print(arr2)

# Array operations

NumPy is a fundamental package for numerical computing in Python.

```
std_dev = np.std(arr)
 print(mean, std_dev)
Pandas
Pandas is a powerful library for data manipulation and analysis.
 import pandas as pd
 # Creating DataFrame
 data = {
     'Name': ['Alice', 'Bob', 'Charlie'],
    'Age': [25, 30, 35]
```

## # Displaying DataFrame print(df)

Matplotlib

# Line plot

plt.xlabel('X-axis')

import seaborn as sns

# Load example dataset

tips = sns.load\_dataset("tips")

print(df)

Matplotlib is a plotting library for creating static, animated, and interactive visualizations.

```
plt.ylabel('Y-axis')
plt.title('Line Plot')
plt.show()
# Bar plot
plt.bar(['A', 'B', 'C'], [10, 20, 30])
plt.xlabel('Categories')
plt.ylabel('Values')
plt.title('Bar Plot')
plt.show()
Seaborn
Seaborn is a statistical data visualization library based on Matplotlib.
```

## # Scatter plot sns.scatterplot(x="total\_bill", y="tip", data=tips) plt.show()

# Box plot

plt.show()

Scikit-Learn

sns.boxplot(x="day", y="total\_bill", data=tips)

Scikit-Learn is a machine learning library for Python.

model.fit(X\_train, y\_train)

predictions = model.predict(X\_test)

print(f"Mean Squared Error: {mse}")

mse = mean\_squared\_error(y\_test, predictions)

SciPy is a library used for scientific and technical computing.

# Make predictions

# Evaluate the model

from scipy import stats

# Statistical functions data = [1, 2, 3, 4, 5]mean = stats.tmean(data) variance = stats.tvar(data)

print(mean, variance)

```
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LinearRegression
from sklearn.metrics import mean_squared_error
# Load dataset
from sklearn.datasets import load_boston
boston = load_boston()
X = boston.data
y = boston.target
# Split data into training and test sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
random_state=42)
# Train a linear regression model
model = LinearRegression()
```

# Statsmodels is a library for estimating and testing statistical models. import statsmodels.api as sm

**Statsmodels** 

# Load dataset

y = data['mpg']

SciPy

# Print model summary print(model.summary())

data = sm.datasets.get\_rdataset("mtcars").data

# Fit a linear regression model

X = data[['hp', 'wt']]

 $X = sm.add\_constant(X)$ model = sm.OLS(y, X).fit()

**Data Manipulation** 

# Load data from a CSV file df = pd.read\_csv('data.csv')

# Load data from an Excel file df = pd.read\_excel('data.xlsx')

# Load data from a SQL database

filtered\_df = df[df['column\_name'] > 10]

plt.plot(df['column1'], df['column2'])

plt.hist(df['column\_name'], bins=10)

sns.heatmap(df.corr(), annot=True) plt.title('Correlation Heatmap')

# Train a linear regression model

predictions = model.predict(X\_test)

print(f"Mean Squared Error: {mse}")

mse = mean\_squared\_error(y\_test, predictions)

model = LinearRegression() model.fit(X\_train, y\_train)

# Make predictions

# Evaluate the model

grouped\_df = df.groupby('column\_name').sum()

merged\_df = pd.merge(df1, df2, on='common\_column')

**Loading Data** 

import sqlite3

```
conn = sqlite3.connect('database.db')
 df = pd.read_sql_query('SELECT * FROM table_name', conn)
Cleaning Data
 # Handling missing values
 df.fillna(0, inplace=True) # Replace NaNs with 0
 df.dropna(inplace=True)  # Drop rows with NaNs
 # Removing duplicates
 df.drop_duplicates(inplace=True)
 # Renaming columns
 df.rename(columns={'old_name': 'new_name'}, inplace=True)
Transforming Data
 # Filtering data
```

```
Data Visualization
# Line plot
```

plt.show()

# Histogram

plt.show()

# Heatmap

plt.show()

plt.xlabel('X-axis') plt.ylabel('Y-axis') plt.title('Line Plot')

plt.xlabel('Value') plt.ylabel('Frequency') plt.title('Histogram')

# Grouping data

# Merging data

```
Machine Learning
Supervised Learning
 from sklearn.model_selection import train_test_split
 from sklearn.linear_model import LinearRegression
 from sklearn.metrics import mean_squared_error
 # Load dataset
 from sklearn.datasets import load_boston
 boston = load_boston()
 X = boston.data
 y = boston.target
 # Split data into training and test sets
 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
 random_state=42)
```

```
import matplotlib.pyplot as plt
import numpy as np
```

# Plot clusters

plt.scatter(X[:, 0], X[:,

**Unsupervised Learning** 

```
from sklearn.cluster import KMeans
# Generate sample data
X = np.random.rand(100, 2)
# Fit KMeans model
kmeans = KMeans(n_clusters=3)
kmeans.fit(X)
# Predict clusters
labels = kmeans.predict(X)
```

### **NLP**

## What is Natural Language Processing?

Natural Language Processing (NLP) is a sub-field of artificial intelligence (AI) that focuses on the interaction between computers and humans through natural language. The ultimate goal of NLP is to enable computers to understand, interpret, and generate human language in a way that is both meaningful and useful.

## Importance and Applications of NLP

NLP has a broad range of applications, including but not limited to:

- **Language Translation**: Google Translate uses advanced NLP techniques to translate text from one language to another.
- Sentiment Analysis: Analyzing customer reviews on products and services to gauge public opinion.
- **Chatbots and Virtual Assistants**: Siri, Alexa, and Google Assistant use NLP to understand and respond to user queries.
- Text Summarization: Summarizing large documents or articles to extract the main points.

## **Core Concepts**

### **Text Preprocessing**

Before diving into complex NLP techniques, it is crucial to preprocess the text data. This chapter covers essential text preprocessing steps.

### Tokenization

```
import nltk
nltk.download('punkt')
from nltk.tokenize import word_tokenize

text = "Natural Language Processing is fascinating."
tokens = word_tokenize(text)
print(tokens)
```

### Stop Words Removal

```
from nltk.corpus import stopwords
nltk.download('stopwords')

stop_words = set(stopwords.words('english'))
filtered_tokens = [word for word in tokens if not word in stop_words]
print(filtered_tokens)
```

## • Lemmatization and Stemming

```
from nltk.stem import WordNetLemmatizer
nltk.download('wordnet')

lemmatizer = WordNetLemmatizer()
lemmatized_words = [lemmatizer.lemmatize(word) for word in filtered_tokens]
print(lemmatized_words)
```

## Regex

## What is Regex?

Regex, short for Regular Expression, is a sequence of characters that forms a search pattern. It is mainly used for string matching and manipulation. Regular expressions are used in various programming languages, text editors, and utilities to search, match, and manipulate text based on specific patterns.

## **Basic Components of Regex**

- 1. Literals: Match the exact characters. For example, cat matches the string "cat".
- 2. Metacharacters: Special characters that have specific meanings:
  - .: Matches any single character except a newline.
  - Asserts the position at the start of a line.
  - \$: Asserts the position at the end of a line.
  - ``: Matches 0 or more repetitions of the preceding element.
  - +: Matches 1 or more repetitions of the preceding element.
  - ?: Matches 0 or 1 repetition of the preceding element.
  - []: Matches any one of the characters inside the brackets.
  - I : Acts as an OR operator.
  - ( ) : Groups expressions and captures the matched sub-expressions.
  - \\: Escapes a metacharacter.

## **Examples of Regex Patterns**

## 1. Simple Match:

- o Pattern: cat
- Matches: "cat", "concatenate", "catapult"

# 2. Character Classes:

Pattern: [aeiou]

Matches: Any single vowel (a, e, i, o, u)

# 3. Quantifiers:

- Pattern: a\*
- Matches: "", "a", "aa", "aaa", ...

## 4. Anchors:

- Pattern: ^cat
- Matches: "cat" at the beginning of a string

## 5. **Alternation**:

- Pattern: cat | dog
- Matches: "cat" or "dog"

## 6. Grouping and Capturing:

- Pattern: (cat|dog)s Matches: "cats", "dogs"

# **Using Regex in Different Programming Languages**

# **Python**

```
Сору
import re
pattern = r'\\bcat\\b'
text = "The cat sat on the mat."
match = re.search(pattern, text)
if match:
   print("Found:", match.group())
   print("Not found")
```

# **JavaScript** Copy

```
const pattern = /\bcat\\b/;
const text = "The cat sat on the mat.";
const match = text.match(pattern);
if (match) {
    console.log("Found:", match[0]);
} else {
    console.log("Not found");
}
Java
```

### Copy import java.util.regex.\*;

```
public class RegexExample {
    public static void main(String[] args) {
        String pattern = "\\\bcat\\\b";
        String text = "The cat sat on the mat.";
        Pattern compiledPattern = Pattern.compile(pattern);
        Matcher matcher = compiledPattern.matcher(text);
        if (matcher.find()) {
            System.out.println("Found: " + matcher.group());
        } else {
            System.out.println("Not found");
        }
    }
 }
Testing and Validating Regex
```

# To test and validate regular expressions, you can use online tools such as:

RegExr

and RegexPal can be used to test and validate regular expressions.

<u>RegexPal</u>

Regex101

- These tools provide a platform to write, test, and debug regular expressions interactively.
- **TLDR**

Regular expressions are powerful tools for text processing and manipulation. By understanding the basic components and syntax, you can create patterns to match specific text sequences efficiently. Whether

you are working in Python, JavaScript, Java, or any other language, regex can significantly enhance your

ability to handle string data. Regular expressions (regex) are sequences of characters that form search patterns, used for string matching and manipulation. Key components include literals, metacharacters, and various patterns such as simple match, character classes, quantifiers, anchors, alternation, and grouping. Regex can be used in different programming languages like Python, JavaScript, and Java. Online tools like Regex101, RegExr,

## **Basic Stats, Analysis, Data Science Overview**

### **Basic Statistics**

- 1. **Mean, Median, Mode**: Measures of central tendency.
  - Mean: Average of all data points.
  - Median: Middle value in a sorted list of numbers.
  - Mode: Most frequently occurring value.
- 2. Standard Deviation and Variance: Measures of spread.
  - Standard Deviation: How much data varies from the mean.
  - Variance: Square of the standard deviation.
- 3. **Correlation Coefficient (r)**: Measures the strength and direction of a linear relationship between two variables.
  - Range from -1 to 1.
  - Close to 1 or -1 indicates a strong relationship, 0 indicates no linear relationship.

## Data Analytics/Science

### 1. Regression Analysis:

- **Coefficients**: Estimate the change in the dependent variable for a one-unit change in the independent variable.
- **R-squared**: Proportion of variance in the dependent variable that can be predicted from the independent variable(s).
- **P-value**: Statistical significance of the coefficients. A p-value < 0.05 often considered statistically significant.

### 2. ANOVA (Analysis of Variance):

- Tests differences between means of different groups.
- Look at F-statistic and p-value for significance.

### 3. Confusion Matrix (Classification):

- True Positives (TP): Correctly predicted positive observations.
- True Negatives (TN): Correctly predicted negative observations.
- False Positives (FP): Incorrectly predicted positive observations.
- False Negatives (FN): Incorrectly predicted negative observations.
- Metrics like Accuracy, Precision, Recall, F1-score are derived from this.

## **Machine Learning**

## 1. Precision and Recall:

- Precision: TP / (TP + FP). Precision is about being precise, i.e., how accurate your predictions are.
- **Recall**: TP / (TP + FN). Recall tells you how many of the actual positives your model captures through labeling it as positive.

## 2. ROC Curve and AUC:

- Receiver Operating Characteristic curve: plots the true positive rate against the false positive rate at various threshold settings.
- AUC (Area Under Curve): Measures the entire two-dimensional area underneath the entire ROC curve. Higher AUC represents a better model.

## **Big Data Analytics**

## 1. Hadoop Ecosystem Components:

• HDFS for storage, YARN for resource management, MapReduce for processing.

## 2. Spark:

• Faster processing than Hadoop MapReduce, in-memory processing.

## **Tools**

## 1. Python/R Libraries:

- Python: Pandas for data manipulation, Scikit-learn for machine learning, Matplotlib/Seaborn for plotting.
- R: dplyr for data manipulation, ggplot2 for plotting, caret for machine learning.

## 2. **SQL**:

 $\circ \quad \text{Basic commands for data retrieval and manipulation (SELECT, INSERT, UPDATE, DELETE, JOIN)}.$ 

# Interpreting Outputs

- Always look at the context: For example, a high R-squared value in regression is good for prediction but does not imply causation.
- Check assumptions: Many statistical tests have underlying assumptions (normality, homoscedasticity, independence).
- Data Quality: Ensure data is clean and representative.

## **API**

## **Advanced API Concepts**

- 1. Authentication:
  - $\circ\quad$  API Key: A simple token that identifies the calling program.
    - Example: Authorization: Bearer <api-key>
  - **OAuth**: A more secure method, often used for allowing third-party applications to access user data without exposing credentials.
  - JWT (JSON Web Token): A compact, URL-safe means of representing claims to be transferred between two parties.

### 2. Rate Limiting:

- Many APIs restrict the number of requests a client can make in a given time frame to prevent abuse.
- Headers like X-RateLimit-Limit, X-RateLimit-Remaining, and Retry-After are used to communicate this information.

### 3. Pagination:

- Used when dealing with large datasets to break the data into manageable chunks.
- Common methods include:
  - Offset and Limit:
    - Example: GET /users?offset=0&limit=50
  - Cursor-based:
    - Example: GET /users?cursor=abc123

## 4. Versioning:

- Ensures backward compatibility as APIs evolve.
- Common methods:
  - URL Path: https://api.example.com/v1/users
    - **Headers**: Accept: application/vnd.example.v1+json

### 5. Error Handling:

- Good APIs provide meaningful error messages.
- Common fields in error responses:
  - code: A numeric or string code representing the error.
  - message : A human-readable message explaining the error.
  - details: Additional information to help debug the issue.

## **Practical Tips for Using APIs**

- 1. **Documentation**: Always refer to the API's documentation. It provides details on endpoints, methods, required parameters, authentication, and examples.
- 2. **Testing**: Use tools like Postman or Insomnia to test API calls before integrating them into your code.
- 3. **Libraries**: Use libraries in your preferred programming language to simplify making API requests. For example:

```
• Python: requests
```

- JavaScript: axios, fetch
- Ruby: rest-client

# 4. Error Handling in Code:

• Example in Python:

```
import requests

try:
    response = requests.get('<https://api.example.com/v1/users>',
headers={'Authorization': 'Bearer <token>'})
    response.raise_for_status() # Raises HTTPError for bad responses
except requests.exceptions.HTTPError as err:
    print(f"HTTP error occurred: {err}")
except Exception as err:
    print(f"An error occurred: {err}")
else:
    print("Success!", response.json())
```

# Workflow

Let's put it all together with a real-world example. Suppose you want to interact with a hypothetical task management API to create, update, retrieve, and delete tasks.

- 1. Set Up Authentication:
  - Get API key
- 2. Create a Task:

```
curl -X POST <https://api.example.com/v1/tasks> \\\
-H "Authorization: Bearer <api-key>" \\\
-H "Content-Type: application/json" \\\\
-d '{"title": "Buy groceries", "due_date": "2024-05-25"}'
```

3. Retrieve Tasks:

```
curl -X GET <https://api.example.com/v1/tasks> \\\
-H "Authorization: Bearer <api-key>"
```

4. Update a Task:

```
curl -X PUT <https://api.example.com/v1/tasks/1> \\\
-H "Authorization: Bearer <api-key>" \\\\
-H "Content-Type: application/json" \\\\
-d '{"title": "Buy groceries and cook dinner", "due_date": "2024-05-25"}'
```

5. **Delete a Task**:

```
curl -X DELETE <https://api.example.com/v1/tasks/1> \\
-H "Authorization: Bearer <api-key>"
```

# Glossary

# **Algorithms**

An algorithm is a set of instructions we give a computer so it can take values and manipulate them into a usable form. This can be as easy as finding and removing every comma in a paragraph, or as complex as building an equation that predicts how many home runs a baseball player will hit in 2018.

**Back End** 

The back end is all of the code and technology that works behind the scenes to populate the front end with useful information. This includes databases, servers, authentication procedures, and much more. You can think of the back end as the frame, the plumbing, and the wiring of an apartment.

# Big data is a term that suffers from being too broad to be useful. It's more helpful to read it as, "so much

**Big Data** 

data that you need to take careful steps to avoid week-long script runtimes." Big data is more about strategies and tools that help computers do complex analysis of very large (read: 1+ TB) data sets. The problems we must address with big data are categorized by the 4 V's: volume, variety, veracity, and velocity.

Classification Classification is a supervised machine learning problem. It deals with categorizing a data point based on its similarity to other data points. You take a set of data where every item already has a category and look at common traits between each item. You then use those common traits as a guide for what category the new item might have.

As simply as possible, this is a storage space for data. We mostly use databases with a Database Management System (DBMS), like PostgreSQL or MySQL. These are computer applications that allow us to interact with a database to collect and analyze the information inside.

A data warehouse is a system used to do quick analysis of business trends using data from many sources. They're designed to make it easy for people to answer important statistical questions without a Ph.D. in

# database architecture.

**Data Warehouse** 

**Front End** The front end is everything a client or user gets to see and interact with directly. This includes data dashboards, web pages, and forms.

**Fuzzy Algorithms** Algorithms that use fuzzy logic to decrease the runtime of a script. Fuzzy algorithms tend to be less precise than those that use Boolean logic. They also tend to be faster, and computational speed sometimes outweighs the loss in precision.

# and 1. That is, fuzzy logic allows statements like "a little true" or "mostly false."

**Greedy Algorithms** 

A greedy algorithm will break a problem down into a series of steps. It will then look for the best possible solution at each step, aiming to find the best overall solution available. A good example is Dijkstra's

## A process where a computer uses an algorithm to gain understanding about a set of data, then makes predictions based on its understanding. There are many types of machine learning techniques; most are classified as either supervised or unsupervised techniques.

Machine Learning

Overfitting happens when a model considers too much information. It's like asking a person to read a

sentence while looking at a page through a microscope. The patterns that enable understanding get lost in

is a value that we calculate or infer from data. We get the median (a statistic) of a set of numbers by using techniques from the field of statistics. **Training and Testing** 

This is part of the machine learning workflow. When making a predictive model, you first offer it a set of

Statistics (plural) is the entire set of tools and methods used to analyze a set of data. A statistic (singular)

# Underfitting

**Fields of Focus** As businesses become more data-focused, new opportunities open up for people of various skill sets to

Underfitting happens when you don't offer a model enough information. An example of underfitting would be asking someone to graph the change in temperature over a day and only giving them the high and low. Instead of the smooth curve one might expect, you only have enough information to draw a straight line.

## work in AI centers on using machine awareness to solve problems or accomplish some task. In case you didn't know, AI is already here: think self-driving cars, robot surgeons, and the bad guys in your favorite

**Artificial Intelligence (AI)** 

video game. **Business Intelligence (BI)** Similar to data analysis, but more narrowly focused on business metrics. The technical side of BI involves

learning how to effectively use software to generate reports and find important trends. It's descriptive,

This discipline is the little brother of data science. Data analysis is focused more on answering questions

## about the present and the past. It uses less complex statistics and generally tries to identify patterns that can improve an organization.

Data engineering is all about the back end. These are the people that build systems to make it easy for data scientists to do their analysis. In smaller teams, a data scientist may also be a data engineer. In larger groups, engineers are able to focus solely on speeding up analysis and keeping data well organized and easy to access.

# **Data Science**

**Data Journalism** 

the discipline of using data and advanced statistics to make predictions. Data science is also focused on creating understanding among messy and disparate data. The "what" a scientist is tackling will differ greatly by employer. **Data Visualization** 

The art of communicating meaningful data visually. This can involve infographics, traditional plots, or even full data dashboards. Nicholas Felton is a pioneer in this field, and Edward Tufte literally wrote the book.

either recommend or make trading decisions based on huge amounts of data, often on the order of

Given the rapid expansion of the field, the definition of data science can be hard to nail down. Basically, it's

# **Quantitative Analysis** This field is highly focused on using algorithms to gain an edge in the financial sector. These algorithms

data. These are some of the most basic and vital statistical tools to help you get started. Correlation

Correlation is the measure of how much one set of values depends on another. If values increase together, they are positively correlated. If values from one set increase as the other decreases, they are negatively correlated. There is no correlation when a change in one set has nothing to do with a change in the other.

A calculation that gives us a sense of a "typical" value for a group of numbers. The mean is the sum of a list of values divided by the number of values in that list. It can be deceiving used on its own, and in

# practice we use the mean with other statistical values to gain intuition about our data. Median

Mean (Average, Expected Value)

Outlier An outlier is a data point that is considered extremely far from other points. They are generally the result of exceptional cases or errors in measurement and should always be investigated early in a data analysis workflow.

A set of data is said to be normalized when all values have been adjusted to fall within a common range. We normalize data sets to make comparisons easier and more meaningful. For instance, taking movie

ratings from different websites and adjusting them so they all fall on a scale of 0 to 100.

# **Standard Deviation** The standard deviation of a set of values helps us understand how spread out those values are. This statistic is more useful than variance because it's expressed in the same units as the values themselves.

symbol sigma (σ).

**Statistical Significance** 

**Summary Statistics** 

over months or temperature throughout days.

it's only 18 degrees, we have an error (residual) of 2 degrees.

between individual values and their mean for that set.

patterns meaningfully for decision-making improvement.

completion occurs within teams' projects/workflows.

opinions of all Greece.

Summary statistics are measures we use to communicate insights about our data simply. Examples include mean, median, and standard deviation. **Time Series** 

A time series is a set of data that's ordered by when each data point occurred. Think stock market prices

The residual measures how much real value differs from some statistical value we calculated based on the set of data. Given a prediction that it will be 20 degrees Fahrenheit at noon tomorrow when noon hits and

Variance measures how spread out values are within a set. Mathematically, it's the average difference

A result is statistically significant when we judge that it probably didn't happen due to chance. It's highly

used in surveys and statistical studies but not always an indication of practical value.

# insights from data. **Data Exploration**

**Data Pipelines** 

**Data Wrangling (Munging)** 

**ETL (Extract Transform Load)** 

Key warehouses describe stages bringing numerous places raw forms screens ready analysis gifted engineers running behind scenes generally ETL systems operate smoothly! Web Scraping

Process involves pulling website source codes usually scripting identifying information pulling files later

Taking raw form taming until consistent larger datasets replacing/removing affecting analysis/performance

# Clustering Collect/categorize sufficiently similar close points measuring distance complexity increases features added problem space grows larger accordingly clustering techniques apply here!

**Feature Engineering** 

# training testing models measuring relevance predicting target variables choosing subsets highperformance achieved accordingly feature selection applies here too!

Branching questions observations predicting target values tend over-fit models datasets grow large random forests type decision tree algorithm designed reduce over-fitting issues arise frequently otherwise encountered commonly faced scenarios! **Deep Learning** Models use very large neural networks deep nets solving complex problems facial recognition layers

accordingly deep learning applies here too! Translating human knowledge quantitative values computers understand visually representing mug images pixel intensities feature engineering translates effectively understandable formats computers utilize efficiently accordingly feature engineering applies here too!

# Loosely based neural connections brains connected nodes segmented layers input output hidden layers heavy lifters making predictions filtering connections next layers final outputs given predictions made

**Neural Networks** 

**Supervised Machine Learning** 

# **Unsupervised Machine Learning**

Techniques building understandings unlabeled datasets looking patterns classifying shared traits accordingly unsupervised machine learning applies here too!

**Database** 

# **Fuzzy Logic** An abstraction of Boolean logic that substitutes the usual True and False for a range of values between 0

algorithm, which looks for the shortest possible path in a graph.

# **Overfitting**

the noise.

Regression

Regression is another supervised machine learning problem. It focuses on how a target value changes as other values within a data set change. Regression problems generally deal with continuous variables, like how square footage and location affect the price of a house. Statistic vs. Statistics

## training data so it can build understanding. Then you pass the model a test set, where it applies its understanding and tries to predict a target value.

become part of the data community. These are some of the areas of specialization that exist within the data science realm.

A discipline involving research and development of machines that are aware of their surroundings. Most

# **Data Analysis**

Data Engineering

rather than predictive.

This discipline is all about telling interesting and important stories with a data-focused approach. It has come about naturally with more information becoming available as data. A story may be about the data or

informed by data. There's a full handbook if you'd like to learn more.

picoseconds. Quantitative analysts are often called "quants."

**Statistical Tools** There are several statistics data professionals use to reason and communicate information about their

# In a set of values listed in order, the median is whatever value is in the middle. We often use the median along with the mean to judge if there are values that are unusually high or low in the set. This is an early

**Normalize** 

Sample The sample is the collection of data points we have access to. We use the sample to make inferences about a larger population. For instance, a political poll takes a sample of 1,000 Greek citizens to infer the

Mathematically, it's the square root of variance for a set of values. It's often represented by the Greek

Residual (Error)

Parts of a Workflow

**Variance** 

The part where scientists ask basic questions that help them understand their dataset's context during exploration will guide more in-depth analysis later. **Data Mining** The process involves cleaning and organizing data; analyzing it for patterns; communicating these

A collection passes along scripts/functions until appropriately cleaned/transformed suitable task

later stages interchangeably used wrangling/munging terms apply here too!

While every workflow is different, these are some general processes that data professionals use to derive

# **Machine Learning Techniques** The field has grown large enough positions exist exclusively Machine Learning Engineers terms below broad overview common techniques used machine learning today!

**Decision Trees** 

analyses purposes!

starting simple patterns building complexity nuanced understandings accurately classifying predicting final sets outputs given predictions made successfully ideally achieved desired results expected outcomes

**Feature Selection** Identifying valuable traits building models helpful large datasets fewer features decrease time complexity

# successfully ideally achieved desired results expected outcomes accordingly neural networks apply here

Techniques giving well-defined labeled columns knowing exactly looking similar professors handing syllabuses telling expectations finals accordingly supervised machine learning applies here too!