



data science glossary

1. Natural Language Processing

- 1.1. Tokenization
- 1.2. Normalization
- 1.3. Stemming
- 1.4. Lemmatization
- 1.5. Corpus
- 1.6. Stop Words
- 1.7. Parts-of-speech (POS) Tagging
- 1.8. Statistical Language Modeling
- 1.9. Bag of Words
- 1.10. n-grams
- 1.11. Regular Expressions
- 1.12. Zipf's Law
- 1.13. Similarity Measures
- 1.14. Syntactic Analysis
- 1.15. Semantic Analysis
- 1.16. Sentiment Analysis
- 1.17. Information Retrieval

2. Internet of Things (IoT)

- 2.1. 6LoWPAN
- 2.2. Advanced Encryption Standard (AES)
- 2.3. Application Programming Interface (API)
- 2.4. Bluetooth Low Energy (BLE)
- 2.5. Embedded Software

2.6. Gateway

2.7. General Packet Radio Service (GPRS)

2.8. Industrial, Scientific, and Medical (ISM) Band

2.9. Link Budget

2.10. Machine to Machine (M2M)

2.11. Media Access Control (MAC)

3. Predictive Analytics

3.1. Predictive Model

3.2. Artificial Intelligence

3.3. Uplift Model

3.4. Vast Search

3.5. Automatic Suspect Discovery (ASD)

4. Database

4.1. Relational Database

4.2. Database Management System (DBMS)

4.3. Primary Key

4.4. Foreign Key

4.5. Structured Query Language (SQL)

4.6. NoSQL

4.7. Metadata

4.8. Consistency

4.9. Data Redundancy

4.10. ACID

4.11. CAP Theorem

4.12. Sharding

4.13. Key-value Store

4.14. Document Store

4.15. Column-oriented Database

4.16. Graph Database

5. Clustering

5.1. Feature Selection

5.2. Expectation Maximization (EM)

5.3. Distance-based Methods

5.4. Density- and Grid-Based Methods

5.5. Matrix Factorization

5.6. Spectral Methods

5.7. Graph-based Techniques

5.8. Streaming scenario

6. Big Data

6.1. Big Data Volume

6.2. Big Data Velocity

6.3. Big Data Variety

6.4. Big Data Veracity

6.5. Big Data Variability

6.6. Big Data Value

6.7. Predictive Analytics

6.8. Descriptive Analytics

6.9. Prescriptive Analytics

6.10. Database

6.11. Data Warehouse

6.12. ETL

6.13. Business Intelligence

- 6.14. Apache Hadoop
- 6.15. Apache Spark
- 6.16. Data lake
- 6.17. Data mining
- 6.18. Data preparation
- 6.19. Data vault
- 6.20. Data munging
- 6.21. Data wrangling
- 6.22. Data governance
- 6.23. Data stewardship
- 6.24. Data visualization
- 6.25. Data Storytelling

7. Machine Learning

- 7.1. Classification
- 7.2. Regression
- 7.3. Clustering
- 7.4. Association
- 7.5. Decision Trees
- 7.6. Support Vector Machines
- 7.7. Neural Networks
- 7.8. Deep Learning
- 7.9. Reinforcement Learning
- 7.10. (k-fold) Cross-validation
- 7.11. Bayesian
- 7.12. Random Forest

8. deep learning

- 8.1. Artificial Neural Networks (ANNs)
- 8.2. Biological Neuron
- 8.3. Perceptron
- 8.4. Multilayer Perceptron (MLP)
- 8.5. Feedforward Neural Network
- 8.6. Recurrent Neural Network
- 8.7. Activation Function
- 8.8. Backpropagation
- 8.9. Cost Function
- 8.10. Gradient Descent
- 8.11. Vanishing Gradient Problem
- 8.12. Convolutional Neural Network
- 8.13. Long Short Term Memory Network (LSTM)

9. Descriptive Statistics

- 9.1. Population
- 9.2. Sample
- 9.3. Parameter
- 9.4. Statistic
- 9.5. Generalizability
- 9.6. Distribution
- 9.7. Mean
- 9.8. Median
- 9.9. Mode
- 9.10. Skew
- 9.11. Range
- 9.12. Variance

9.13. Standard Deviation

9.14. Interquartile Range (IQR)

10. Cloud Computing

10.1. XaaS (Anything-as-a-Service)

10.2. Software-as-a-Service (SaaS)

10.3. Platform-as-a-Service (PaaS)

10.4. Infrastructure-as-a-Service (IaaS)

10.5. Public Cloud

10.6. Private Cloud

10.7. Hybrid Cloud

10.8. AWS

10.9. Amazon EC2 (Elastic Cloud Compute)

10.10. Amazon Simple Storage Service (S3)

10.11. Cloud Sourcing

10.12. Consumer Cloud

10.13. Multi-tenancy

10.14. Vertical Cloud

10.15. Cloud Portability

10.16. Cloud Backup

10.17. Cloud Enablement

10.18. Cloud Migration

10.19. Cloudstorming

10.20. Cloud Broker

11. Hadoop

11.1. MapReduce

11.2. Hadoop Distributed File System (HDFS)

11.3. Yet Another Resource Negotiator (YARN)

11.4. HBase

11.5. Hive

11.6. Apache Pig

11.7. Apache Spark

11.8. Sqoop

11.9. Oozie

11.10. ZooKeeper

11.11. Apache Flume

11.12. Hue

11.13. Mahout

11.14. Ambari

11.15. Hadoop Common

12. Apache

12.1. RDD

12.2. DataFrame

12.3. Dataset

12.4. MLlib

12.5. ML Pipelines

12.6. GraphX

12.7. Spark Streaming

12.8. Structured Streaming

12.9. spark-packages.org

12.10. Catalyst Optimizer

12.11. Tungsten

12.12. Continuous Applications

12.13. In-memory computing