# Data Science Introduction

Ivan Chen

Feb 2020

https://www.linkedin.com/in/yaohua-ivan-chen-ph-d-018b0b14/

**Agenda**

- What is Data Science?
  - Data Science vs. Data Analytics (BI)
  - Data Scientist vs. Data Analyst
- What is Data Science Life Cycle (DSLC)?
- What Programming Language do we choose, Python or R?
- What is Machine Learning? Why needs Machine Learning?
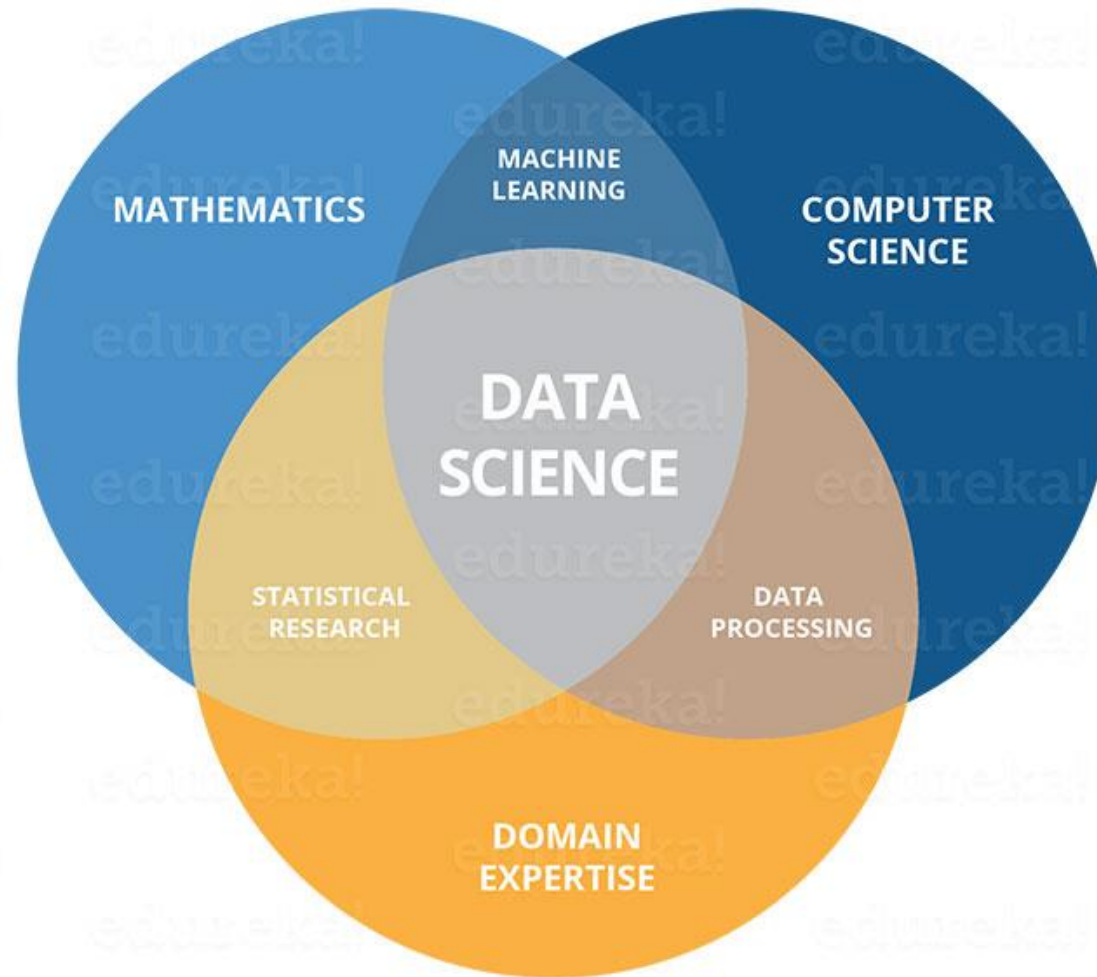- How to do Machine Learning?
- How to Evaluate a Model?
- Q & A

# WHAT IS DATA SCIENCE?

# What is Data Science?
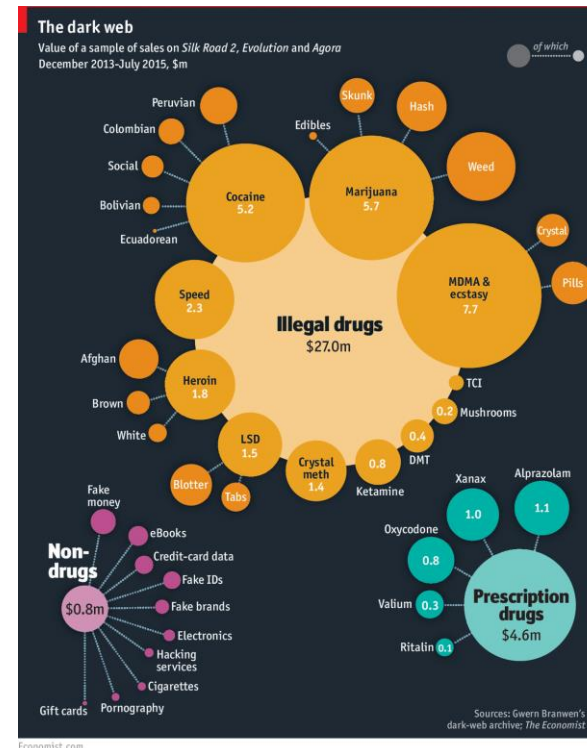
Data science comprises three distinct and overlapping areas:
- The skills of a *statistician* who knows how to model and summarize datasets
- The skills of a *computer scientist* who can design and use algorithms to efficiently store, process, and visualize this data
- The *domain expertise* who formulate the right questions and to put their answers in context.
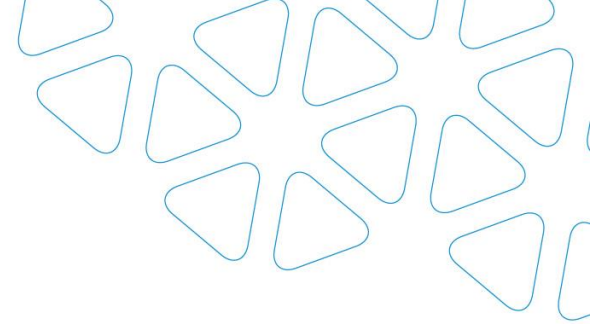


Data Science is to use Data as base, Programming as legs, Machine learning as backbone and Business logics as heart, to discover useful hidden patterns or insights from the data.

# Key Components of Data Science

- Business Understanding

- Data Mining or Discovery

- Data Exploration

- Data Engineering

- Feature Engineering

- Modeling & Evaluation

- Visualization

- Software Implementation

- Product Deployment & Monitoring



**The dark web**
Value of a sample of sales on *Silk Road 2, Evolution* and *Agora*
December 2013-July 2015, $m

# Data Science vs. Data Analytics (BI)

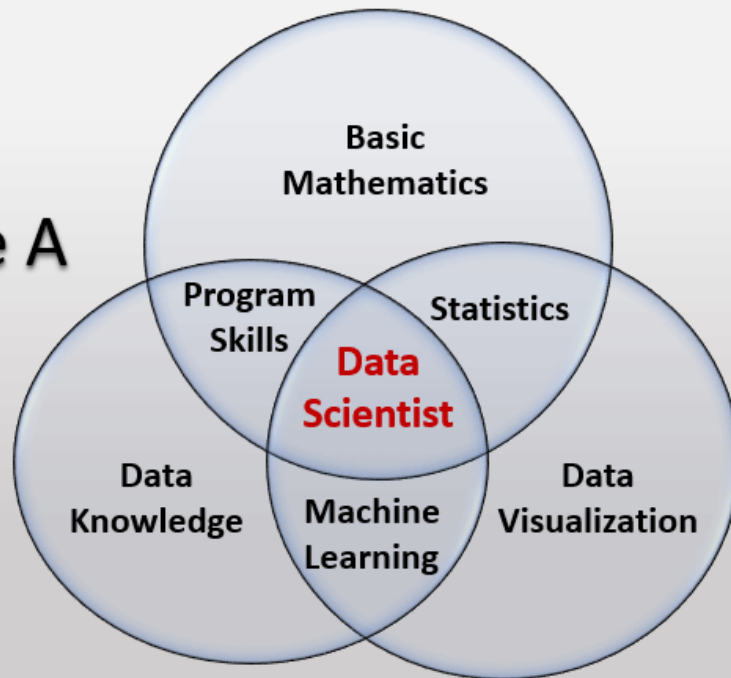| Factors | Business Intelligence | Data Science |
|---|---|---|
| Concept | Deals with data analysis on the business platform. | Consists of several data operations in various domains. |
| Scope | BI analyzes past data | Past data is analyzed for future predictions. |
| Data | Handling static and structured data | Both structured & unstructured data that is also dynamic. |
| Data Storage | Data stored mostly in data-warehouses | Data utilized is distributed in real time clusters. |
| Procedure | BI helps companies to solve questions. | Questions are both curated and solved by data scientists. |
| Tools | MS Excel, SAS BI, Sisense, Microstrategy | Python, R, Hadoop/Spark, SAS, TensorFlow. |

## Data Science Vs. Business Intelligence
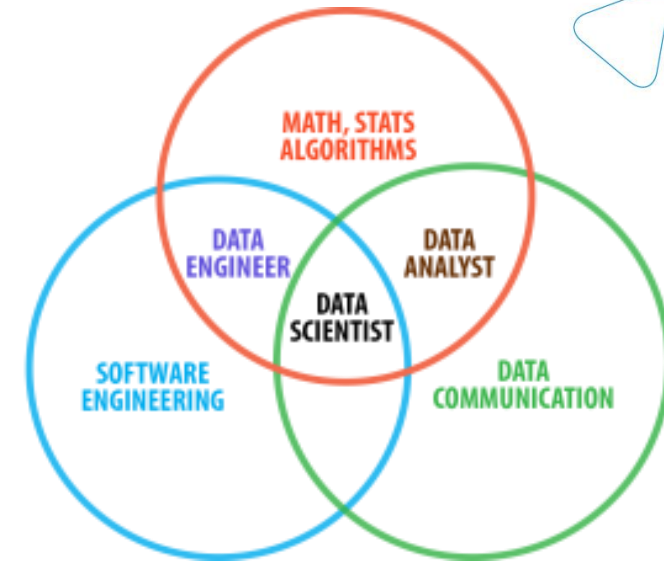
▶ Analytics Spectrum:

| | | |
|---|---|---|
| Descriptive | What happened? | } Traditional BI |
| Diagnostic | Why did it happen? | |
| Predictive | What will happen? | |
| Prescriptive | What should I do? | |

# Data Scientist vs. Data Analyst
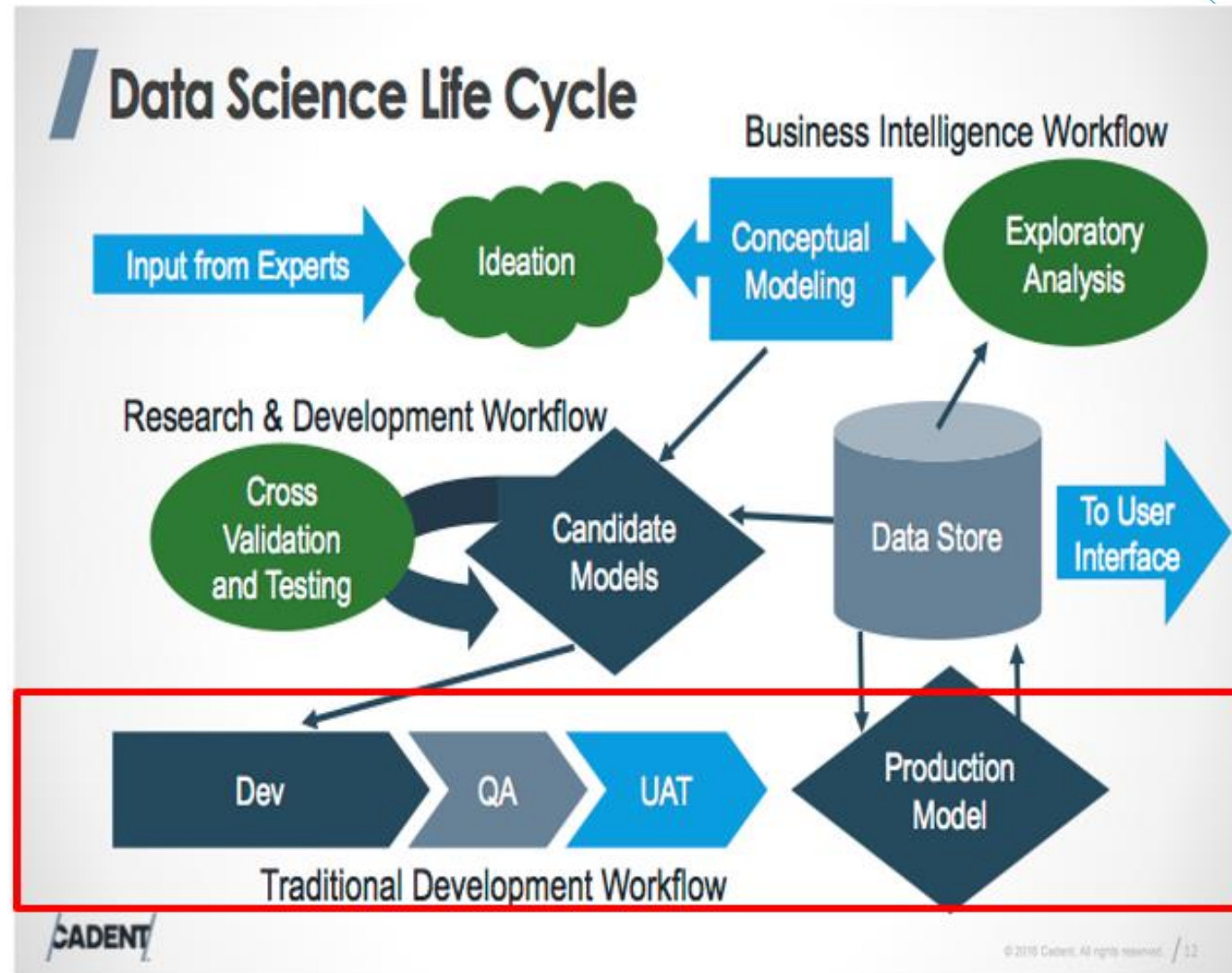
# When you interview a data scientist...

# WHAT IS DATA SCIENCE LIFE CYCLE?

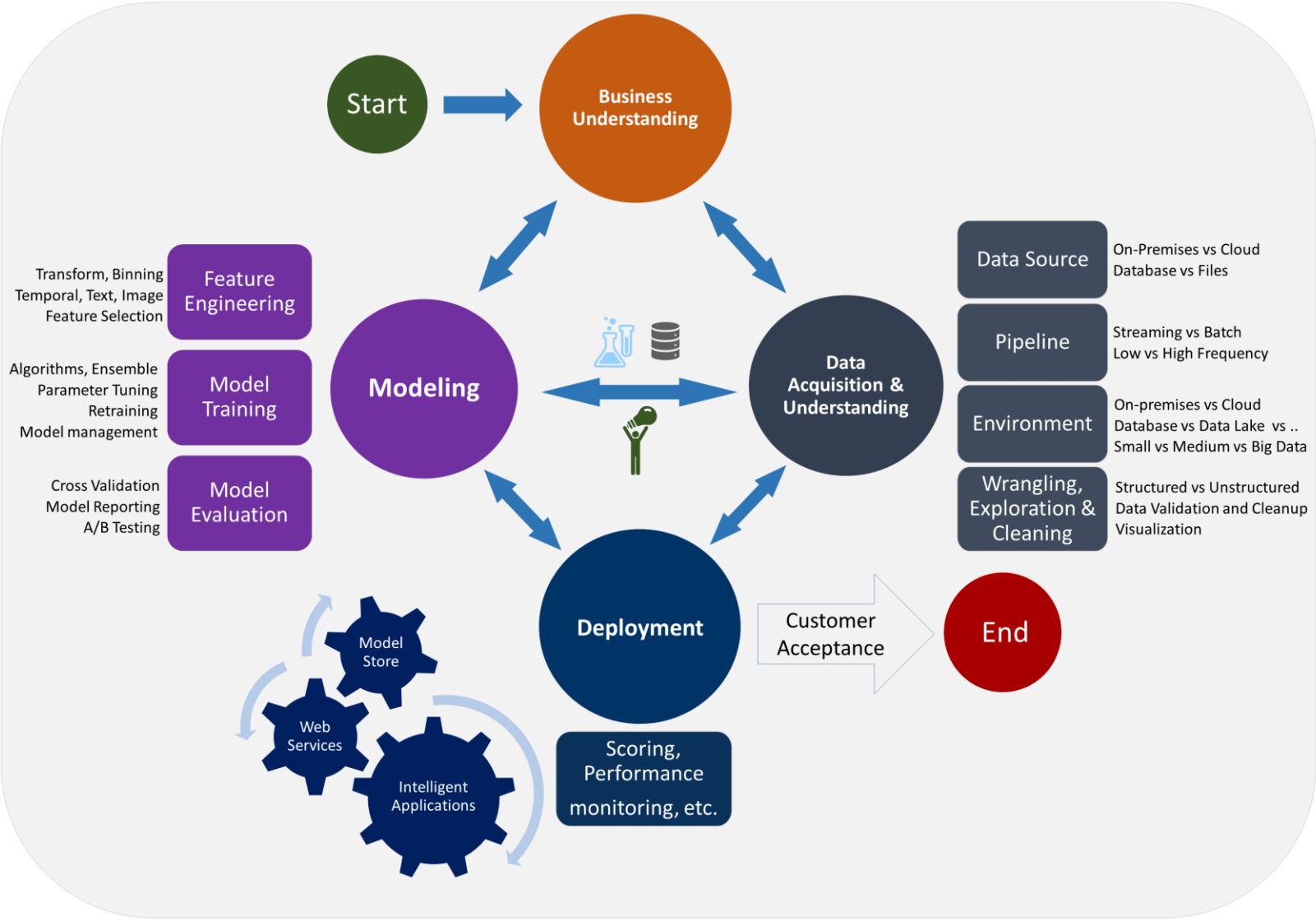**Data Science Life Cycle**

# Data Science Lifecycle

**Data Science Life Cycle**

WHAT LANGUAGE DO WE CHOOSE FOR DATA SCIENCE?

# Programming Language – Python

- Functional Scripting & Object-Oriented Programming
- General Purpose including Software Development, Data Science & Data Engineering
- Huge Open Source Libraries/Packages & Community
- Readability & Maintainability
- Less code base complexity
- Support by most all vendors



Top 20 Technology Skills in Data Scientist Job Listings

Medium



Google Trends, Jan 2012 – Aug 2017

KDnuggets



Projections of future traffic for major programming languages
Future traffic is predicted with an STL model, along with an 80% prediction interval.

Stack Overflow Blog

# Programming Language – Python vs. R



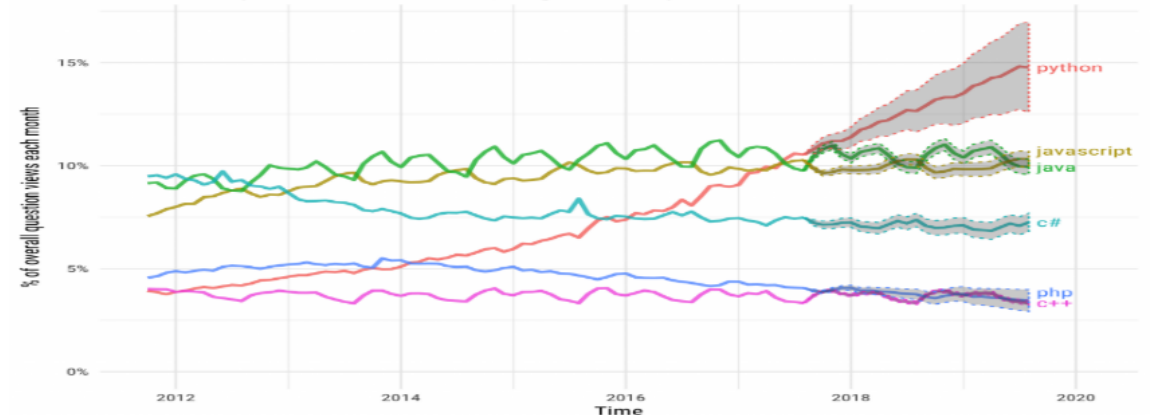| Which Software Should I Choose? | Python | R |
|---|---|---|
| Best for: | General programming; Data analysis; Deep learning; Repeated tasks | Statistical analysis; Data analysis; Single passes of data |
| Availability | Free, open source | Free, open source |
| Easy to learn? | Yes, especially for software engineers | Steep learning curve; Relatively easier if no prior coding experience |
| Advantages | Easy to deploy; General purpose language; Widely used by corporations | Minimal coding required for statistical models |
| Disadvantages | Requires rigorous testing | Very statistics oriented; Not a general-purpose program |

## Comparision of R and Python on following basis

| Score out of 5 [ 1 - lowest ][ 5 - Highest ] | R | Python |
|---|---|---|
| ● Availability / cost | 5 | 5 |
| ● Easy of learning | 3.5 | 4 |
| ● Data handling capabilites | 4 | 4 |
| ● Graphical capabilites | 4.5 | 4.5 |
| ● Advancement in tool | 4.5 | 4.5 |
| ● Job scenario | 4.5 | 4.5 |
| ● Support and community | 3.5 | 3.5 |
| ● Deep Learning support | 3 | 4.5 |

# Python Libraries for Data Science

# WHAT IS MACHINE LEARNING?

# What is Machine Learning?

# A Little History of Machine Learning & Data Science

# Why need Machine Learning in Data Science?

- Problems for which existing solutions require a lot of fine-tuning or long lists of rules: one Machine Learning algorithm can often simplify code and perform better than the traditional approach.

- Complex problems for which using a traditional approach yields no good solution: the best Machine Learning techniques can perhaps find a solution.

- Fluctuating environments: a Machine Learning system can adapt to new data.

- Getting insights about complex problems and large amounts of data.



Machine Learning

Deployment

Application Development

Big Data Processing

Data Storage

ETL

# Machine Learning Landscape



coggle
made for free at coggle.it

**Types of tasks**

Semi-supervised learning
- Recommendation systems
- Natural language generation
- Image generation

Supervised learning
- Image segmentation
- Classification
- Regression

Unsupervised learning
- Clustering
- Anomaly detection
- Dimensionality reduction

Reinforcement learning

**Machine Learning**

**Applications**

Dimensionality reduction
- Image compression
- Feature engineering for ML
- Audio compression

Natural language processing
- Topic modeling
- Text classification
- Sentiment analysis
- Machine translation
- Natural language generation
- Speech recognition
- Text-to-speech
- Text analysis (tagging, parsing etc)
- Summarization
- Entity recognition
- Keywords extraction

Computer vision
- Image classification
- Image segmentation
- Objects detection

Anomaly detection
- Outlier detection
- Novelty detection
- Fraud detection

Time series
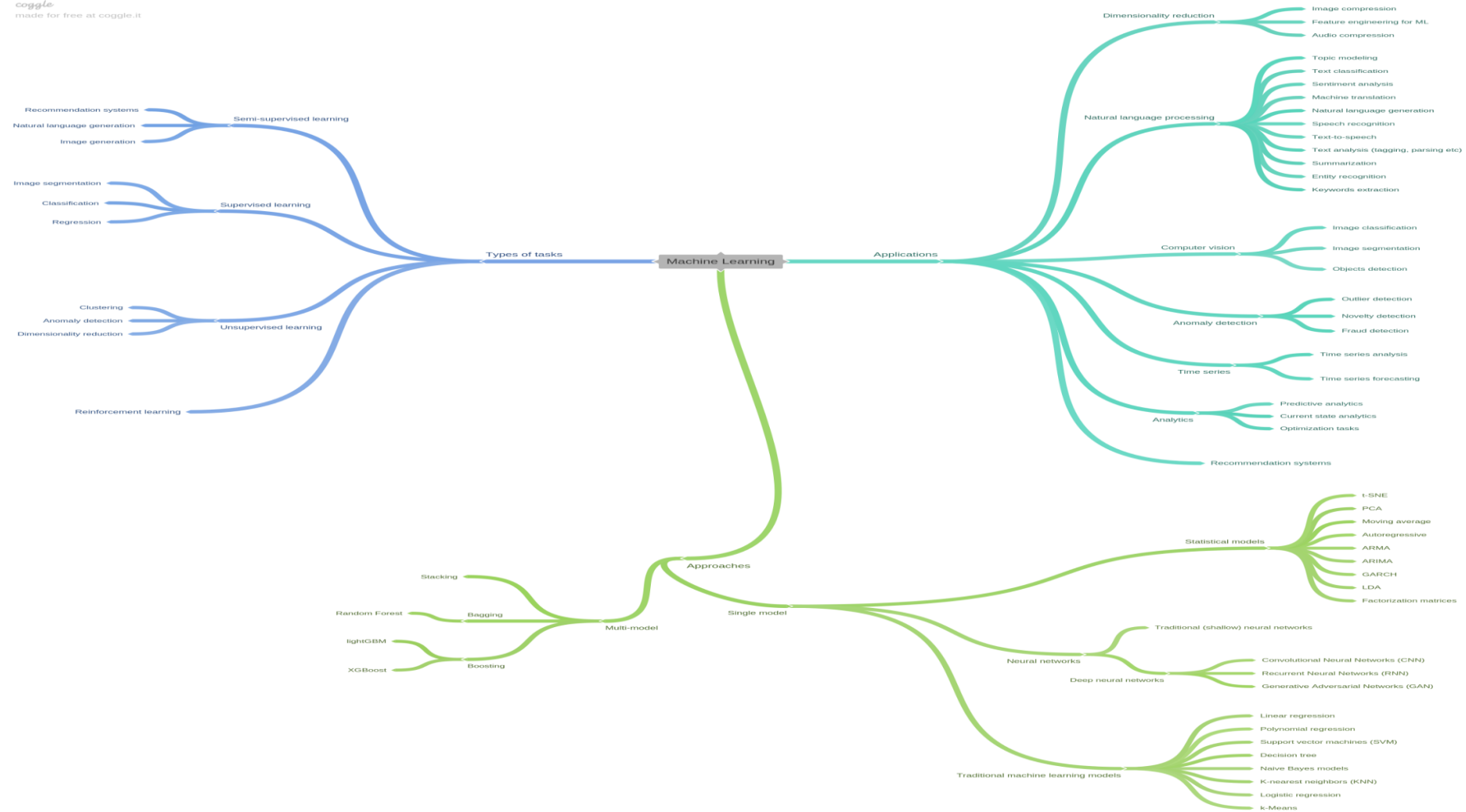- Time series analysis
- Time series forecasting

Analytics
- Predictive analytics
- Current state analytics
- Optimization tasks

Recommendation systems

**Approaches**

Multi-model
- Stacking
- Bagging
  - Random Forest
- Boosting
  - lightGBM
  - XGBoost

Single model

Statistical models
- t-SNE
- PCA
- Moving average
- Autoregressive
- ARMA
- ARIMA
- GARCH
- LDA
- Factorization matrices

Neural networks
- Traditional (shallow) neural networks
- Deep neural networks
  - Convolutional Neural Networks (CNN)
  - Recurrent Neural Networks (RNN)
  - Generative Adversarial Networks (GAN)

Traditional machine learning models
- Linear regression
- Polynomial regression
- Support vector machines (SVM)
- Decision tree
- Naive Bayes models
- K-nearest neighbors (KNN)
- Logistic regression
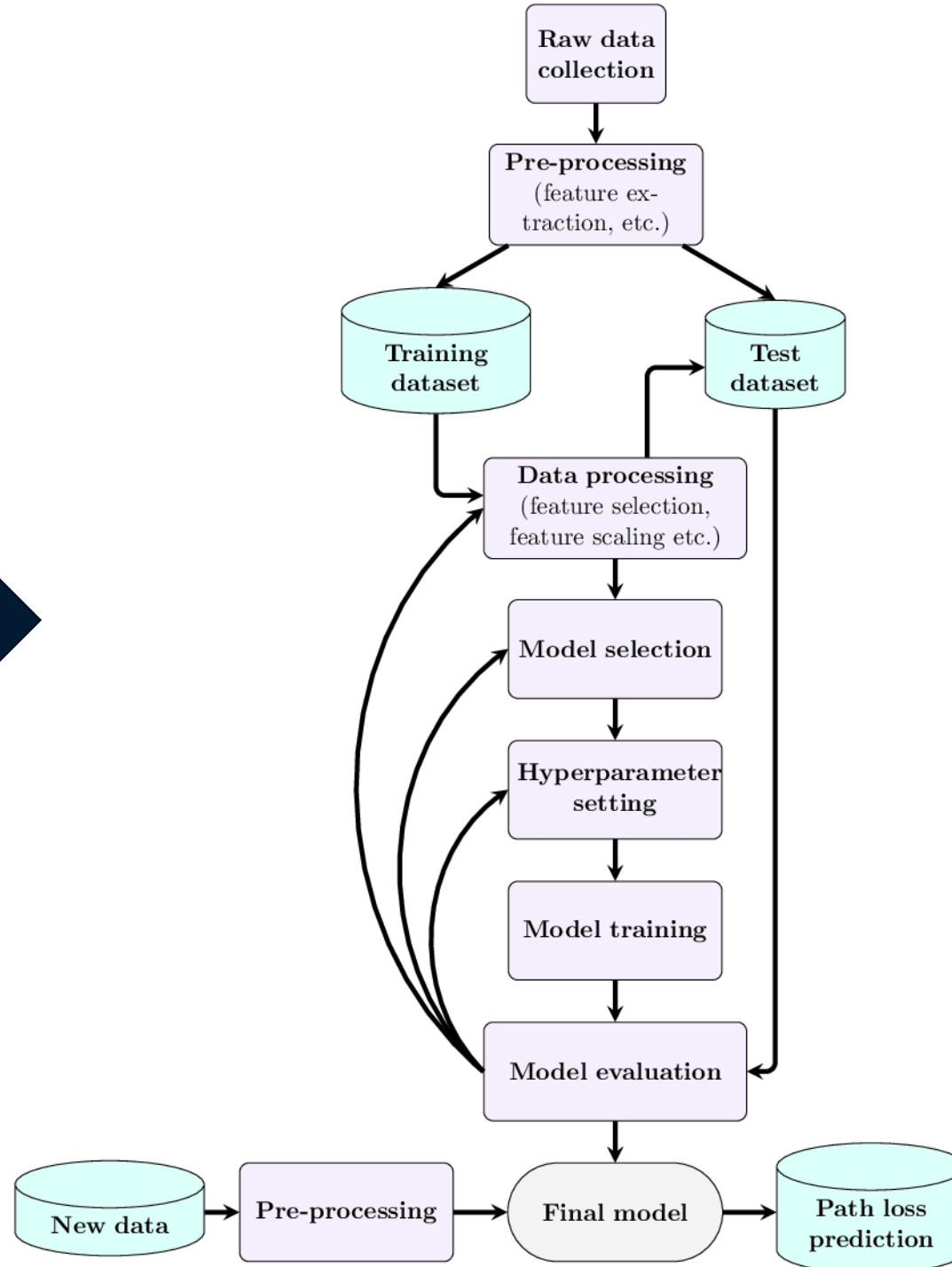- k-Means

# A Simple Machine Learning System Choosing Strategy



Machine Learning Algorithms Cheat-sheet
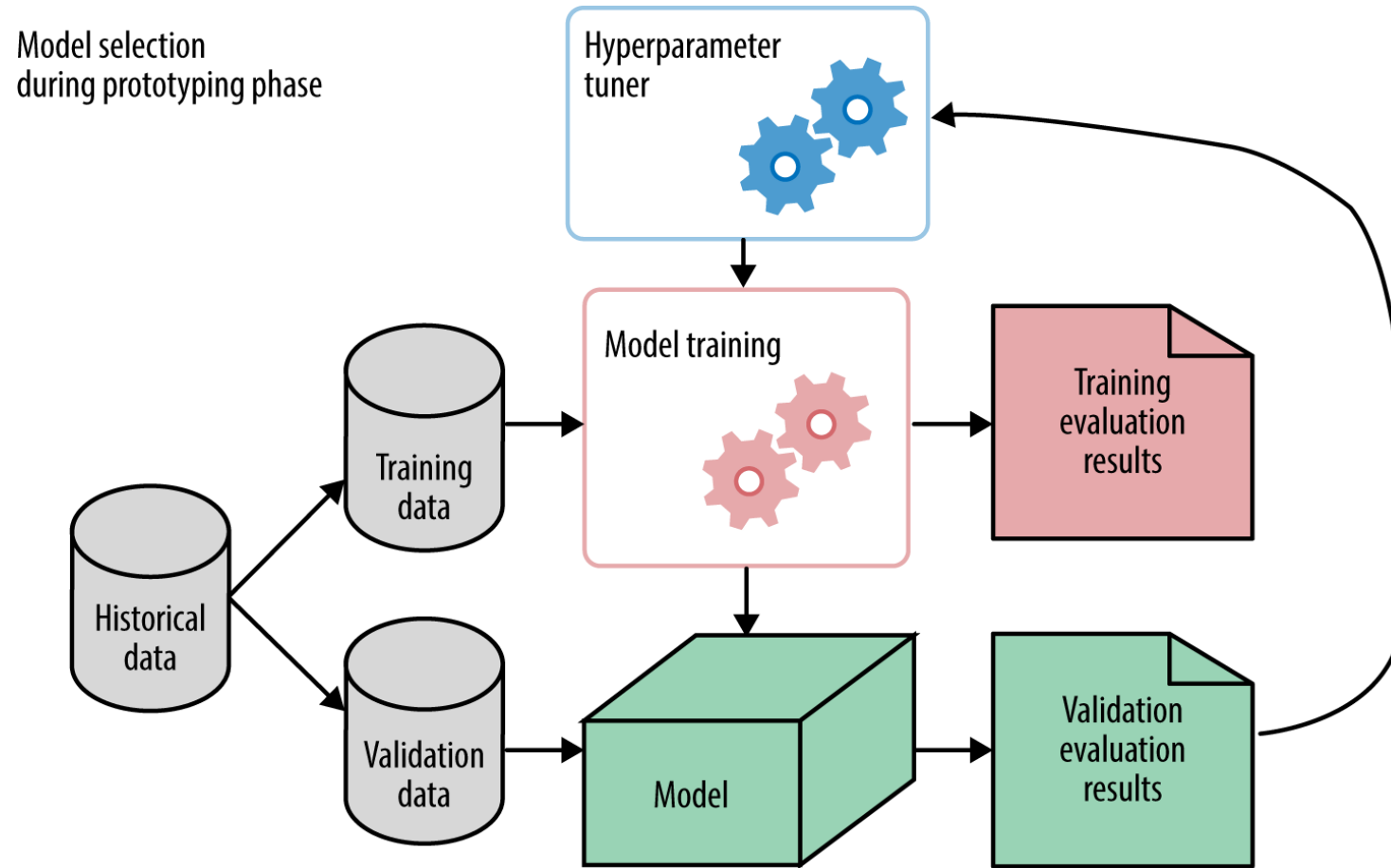
# HOW TO DO MACHINE LEARNING IN DATA SCIENCE?

**Data Science Modeling Work Flow**

Raw data collection

Pre-processing (feature extraction, etc.)

Training dataset

Test dataset

Data processing (feature selection, feature scaling etc.)

Model selection

Hyperparameter setting

Model training

Model evaluation

New data → Pre-processing → Final model → Path loss prediction

# HOW TO EVALUATE A MODEL?

# Model Evaluations



Model selection during prototyping phase

Hyperparameter tuner

Model training

Training evaluation results

Historical data

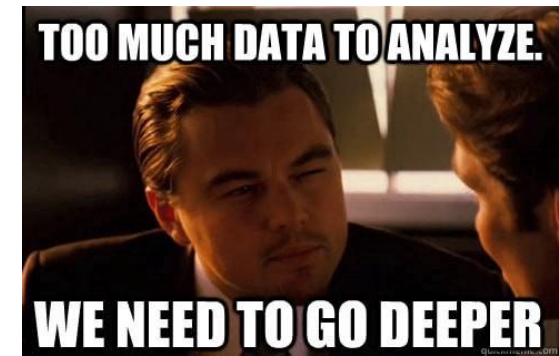Training data

Validation data

Model

Validation evaluation results

# Model Evaluations

- Setup Training (& Development) and Test Sets
  - How large do they need to be? 1000 ~ 10000
  - Must they be on same distribution?
- Setup a Single-number Evaluation Metric to Optimize
  - Mean Absolute Error (MAE) vs. Root Mean Squared Error (RMSE)
  - Precise vs. Recall → F1 Score
  - Accuracy vs. Speed
- Underfit and Overfit - the two big sources of error:
  - 86% accuracy on training set vs. 85% accuracy on testing set
  - 85% accuracy on training set vs. 96% accuracy on testing set
  - 96% accuracy on training set vs. 85% accuracy on testing set
  - 96% accuracy on training set vs. 95% accuracy on testing set

# Techniques for reducing avoidable underfit

- **Increase the model size** (such as number of neurons/layers): This technique should allow you to fit the training set better. If you find that this increases overfitting (variance), then use regularization, which will usually eliminate the increase in overfitting.

- **Reduce or eliminate regularization** (L2 regularization, L1 regularization, dropout): This will reduce underfitting, but increase overfitting.

# Techniques for reducing avoidable overfitting

- **Add more training data** : This is the simplest and most reliable way to address overfitting, as long as you have access to significantly more data and enough computational power to process the data.

- **Add regularization** (L2 regularization, L1 regularization, dropout): This technique reduces overfitting but may increase underfitting.

- **Add early stopping** (i.e., stop gradient descent early): This technique reduces overfitting but may increases underfitting. Early stopping behaves a lot like regularization methods, and some people call it a regularization technique.

- **Feature selection to decrease number/type of input features:** This technique might help with overfitting problems, but it might also increase underfitting.

- **Decrease the model size** (such as number of neurons/layers): *Use with caution.* This technique could decrease overfitting, while possibly increasing underfitting. NOT RECOMMEND this technique for addressing overfitting.

## Techniques for reducing avoidable underfitting and overfitting

- **Modify input features based on insights from error analysis** : Say your error analysis inspires you to create additional features that help the algorithm to eliminate a particular category of errors. These new features could help with both underfitting and overfitting.

- **Modify model architecture** (such as neural network architecture) so that it is more suitable for your problem: This technique can affect both underfitting and overfitting.

# More Questions about data sets

- When you should train and test on different distributions?

- How to decide whether to use all your data?

- How to decide whether to include inconsistent data?

- Can we use artificial data synthesis?

## References:

- Introduction to Machine Learning with Python, Andreas Muller and Sarah Guido, O'Reilly, 2016
- Python Data Science Handbook, Jake VanderPlas, O'Reilly, 2017
- Machine Learning Yearning, Andrew Ng, 2018.
- Hands-on Machine Learning with Scikit-Learn, Keras & Tensorflow, Aurelien Geron, O'Reilly, 2019

# Appendix – Data Science Dev Environment Setup

- Data Science Tool Kit – [New Version on GitHub](#)

- IDEs – Jupyter Notebook, Jupyter Lab and/or Visual Studio Code

- Python Libraries included in Anaconda 2019.10