# General Approach to Data Science
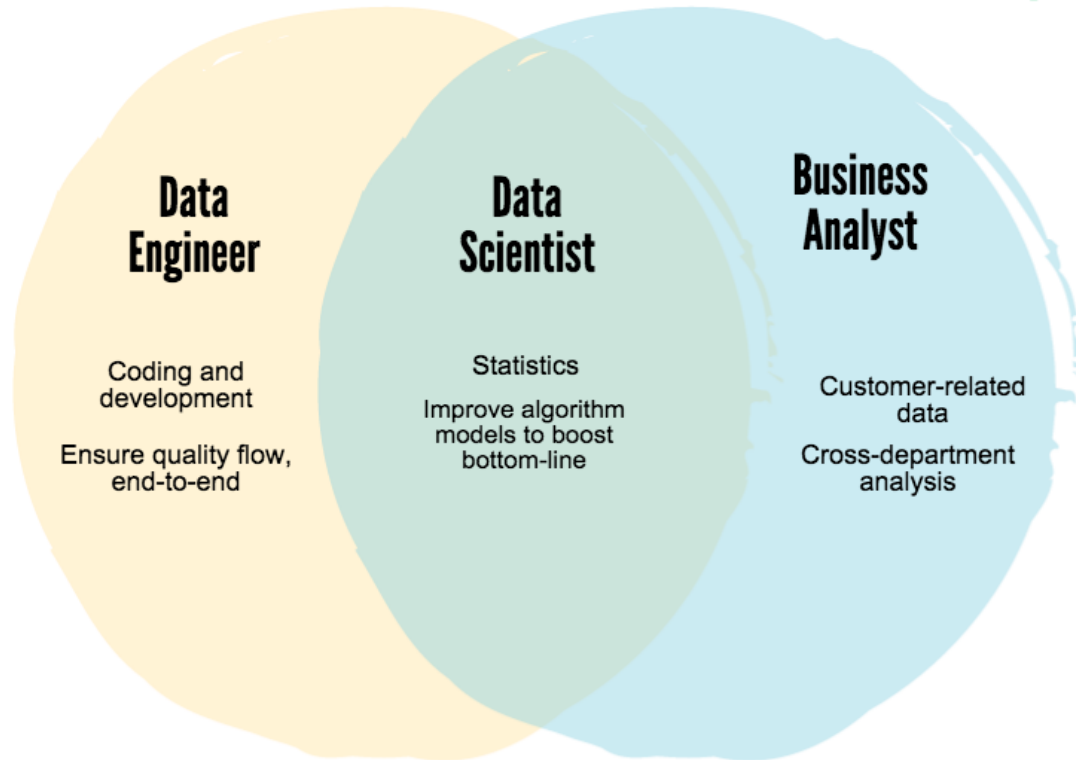
Bhuvan M S

msbhuvanbhuvi@gmail.com
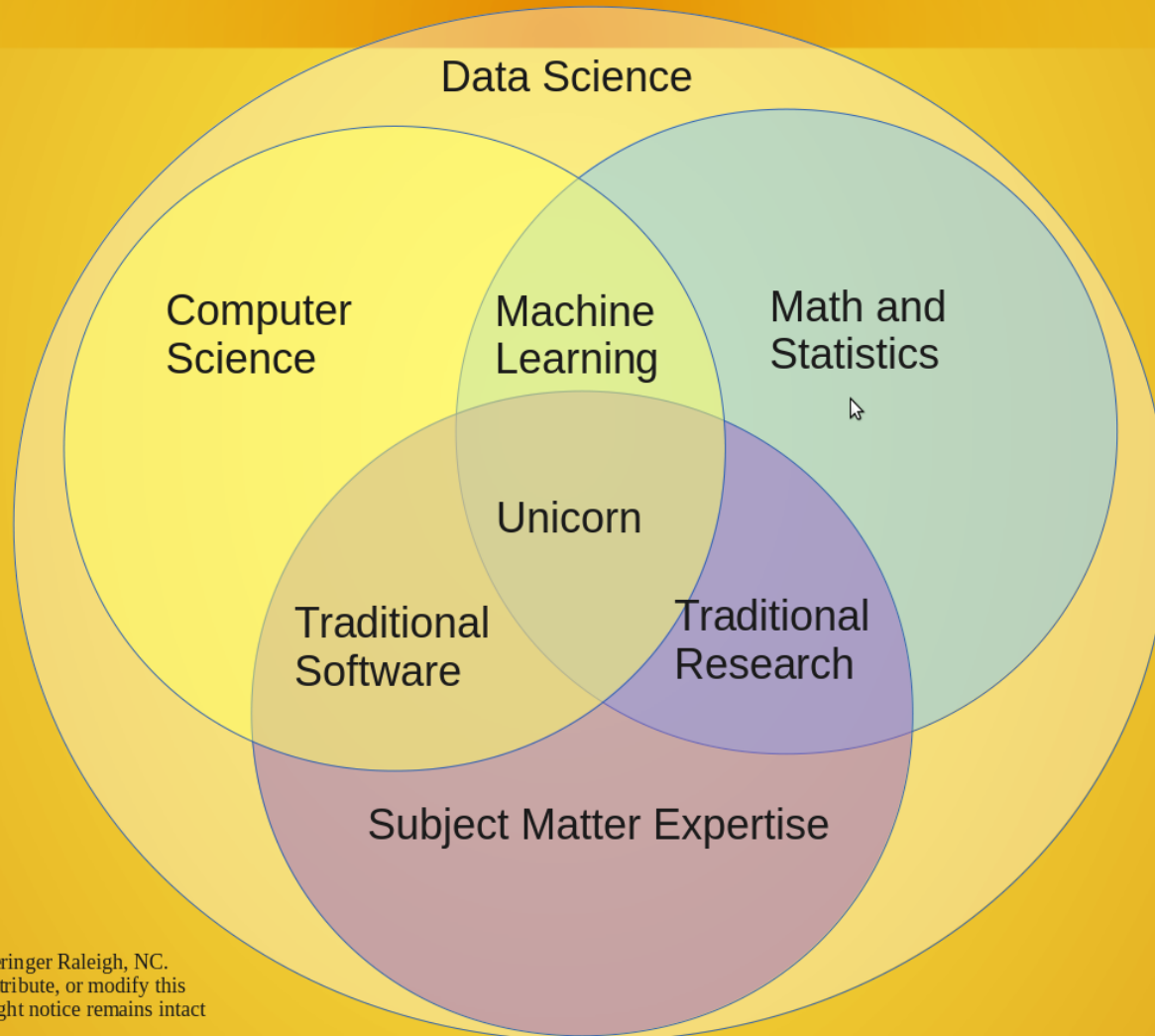
# Contents

- Scope of Data Science
- Skills for a Data Scientist
- How to approach towards a data centric solution to a problem?
  - Analysis Pipeline
- Tools for implementing
- Deployment Framework
- Note on NLP specific approach

# Scope of Data Science

Data Science Venn Diagram v2.0

# Skills of a Data Scientist



Data Scientist Skill Set

# The Process and Workflow

Data Science Process

# Exploratory Data Analysis

► Slicing and Dicing

► Suggest hypotheses about the causes of observed phenomena

► Assess assumptions on which statistical inference will be based

► Support the selection of appropriate statistical tools and techniques

► Provide a basis for further data collection through surveys or experiments.

# Feature Engineering

- Features: The dimensions of the data!
- Data Types: Binary, Numeric, Categorical, Ordinal

- Features Identification

- Features Extraction

# Data Preprocessing

▶ Data Distribution and Data Scale observation

▶ Data Integration

▶ Data cleaning

 ▶ Missing values

 ▶ Noise

▶ Dimensionality reduction

 ▶ PCA

 ▶ Correlation Analysis

▶ Data Transformation: Normalization

 ▶ Data Type specific

# Machine Learning

# ML pipeline

# ML Techniques

- Supervised:
  - Classification

- Unsupervised
  - Rule Based Classification
  - Clustering

- Association Rule Mining

# Feature Importance Mining

- Feature Weights

- Ablation Study

- Random Forests

# Deep Learning

- Neural Networks
  - Classification

- Auto-encoders
  - Automatic Feature Learning

- Self Organizing Maps
  - Clustering



Hidden Layers

Input Layer

Output Layer

# Evaluation and Tuning

- Accuracy
- RMSE
- F1 measure
- Precision and Recall
- Receiver Operating Characteristics (ROC)

# We evaluate performance of the model



non-default cases

default cases

cut-off value

test result value

## ROC curve

TPF, sensitivity

FPF, 1-specificity

Line of no discrimination (Random guess)

# Evaluation of a Campaign

- Confusion Matrix

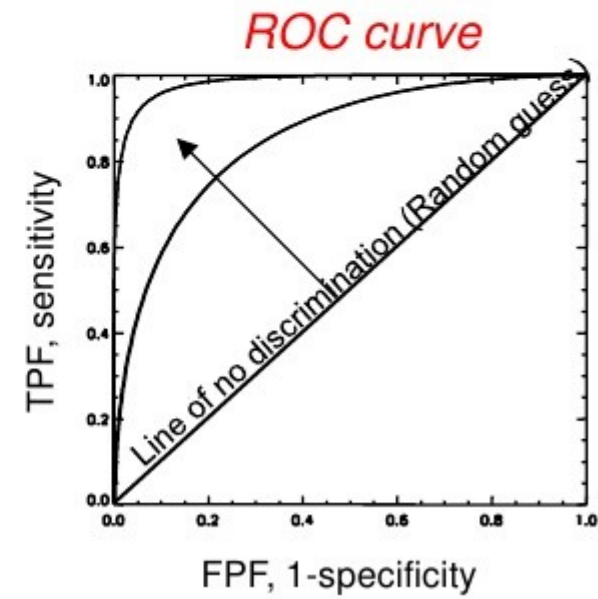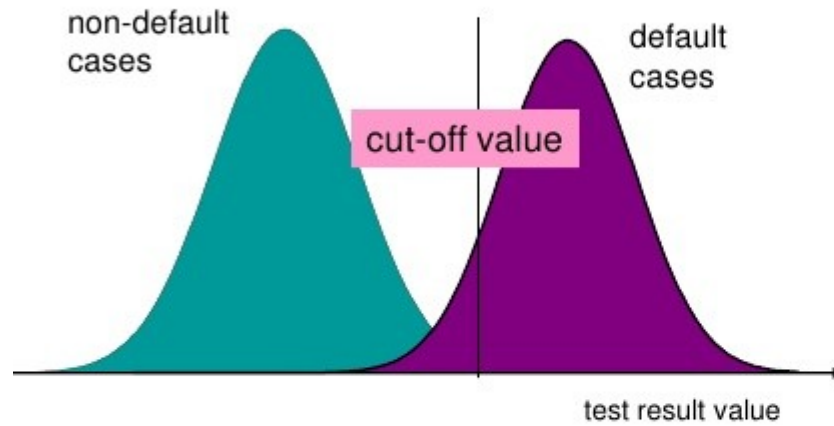| Predicted class ($c_i$) | | True Class ($y_i$) | |
|---|---|---|---|
| | | Churner ($y_i$=1) | Non-Churner($y_i$=0) |
| | Churner ($c_i$=1) | TP | FP |
| | Non-Churner ($c_i$=0) | FN | TN |

- Accuracy $= \dfrac{TP+TN}{TP+TN+FP+FN}$

- Recall $= \dfrac{TP}{TP+FN}$

- Precision $= \dfrac{TP}{TP+FP}$

- F1-Score $= 2\dfrac{Precision * Recall}{Precision+Recall}$

# Tools

- Programming
  - Python: Scikit-Learn, Py-Weka, NLTK, PyBrain
  - PySpark: mllib (distributed)
  - R
  - Matlab, Neural Network Toolkit, Image Processing Toolkit

- Experimentation
  - Weka: Explorer, Experimenter, Knowledge Flow
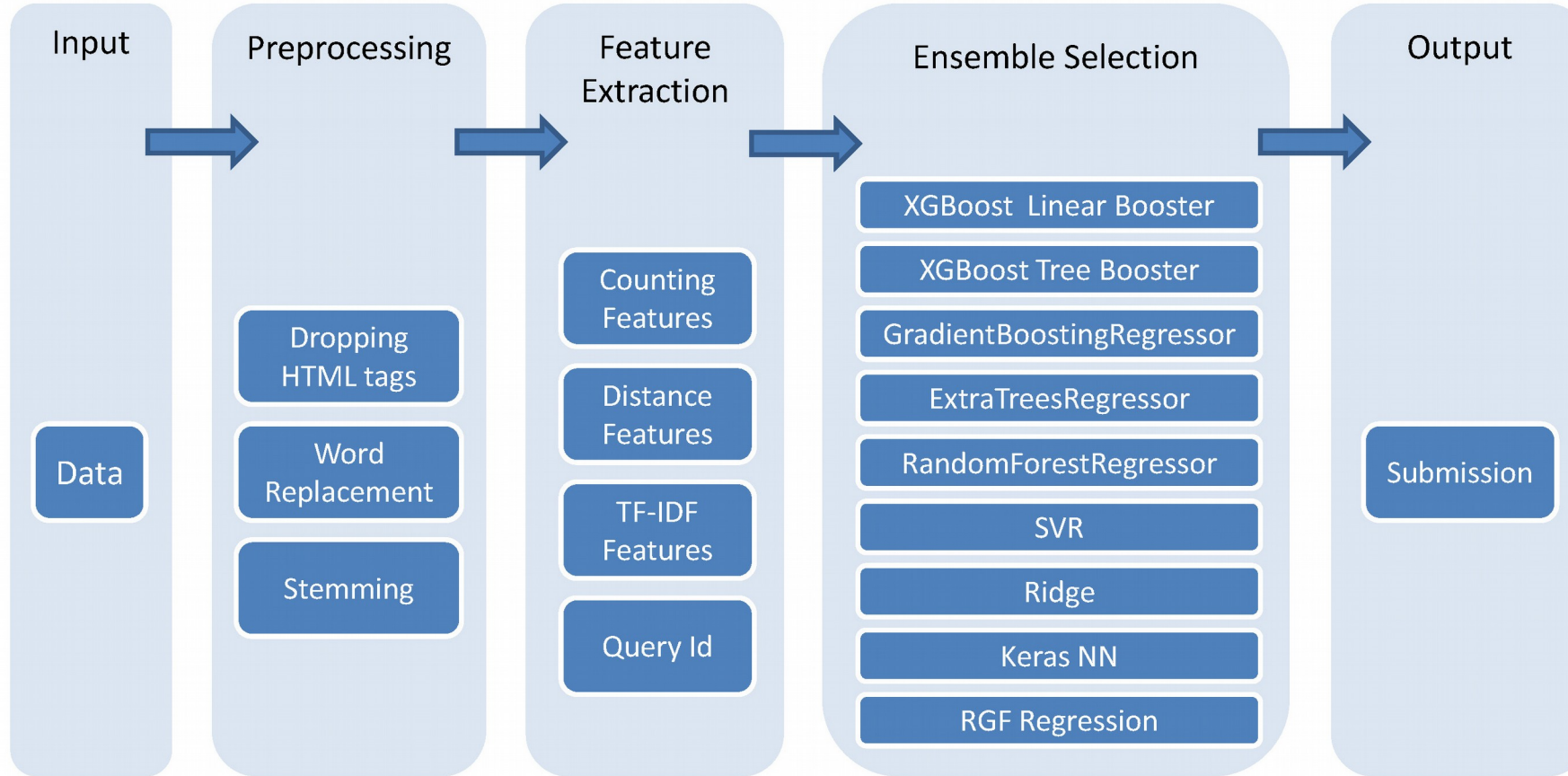
- Visualization
  - Python: matplotlib
  - Tableau
  - D3 js

# Deployment Frameworks

- Database: Scalable, Distributed
  - Graph Based: Neo4j
  - Document Store: MangoDB
  - Other: HDFS, Spark (RDD)

- Handle Big Data – Map Reduce Programming Paradigm
  - Apache Spark
    - MLLIB
    - Streaming
  - Apache Storm
    - Topology – Spout, Bolt

# Questions?

Thank you