

به نام آنکه به شماره موجودات آگاه است

آمار و رشته‌ی علوم داده‌ها

عادل محمدپور

دانشگاه صنعتی امیرکبیر (پلی‌تکنیک تهران)

AdelM.ir

adel@aut.ac.ir



Learning Data Science

© Jamie Whitehorn 2014 v1.1



Skills

Curious

Creative

Tenacious

Resourceful

Inventive

Innoviative

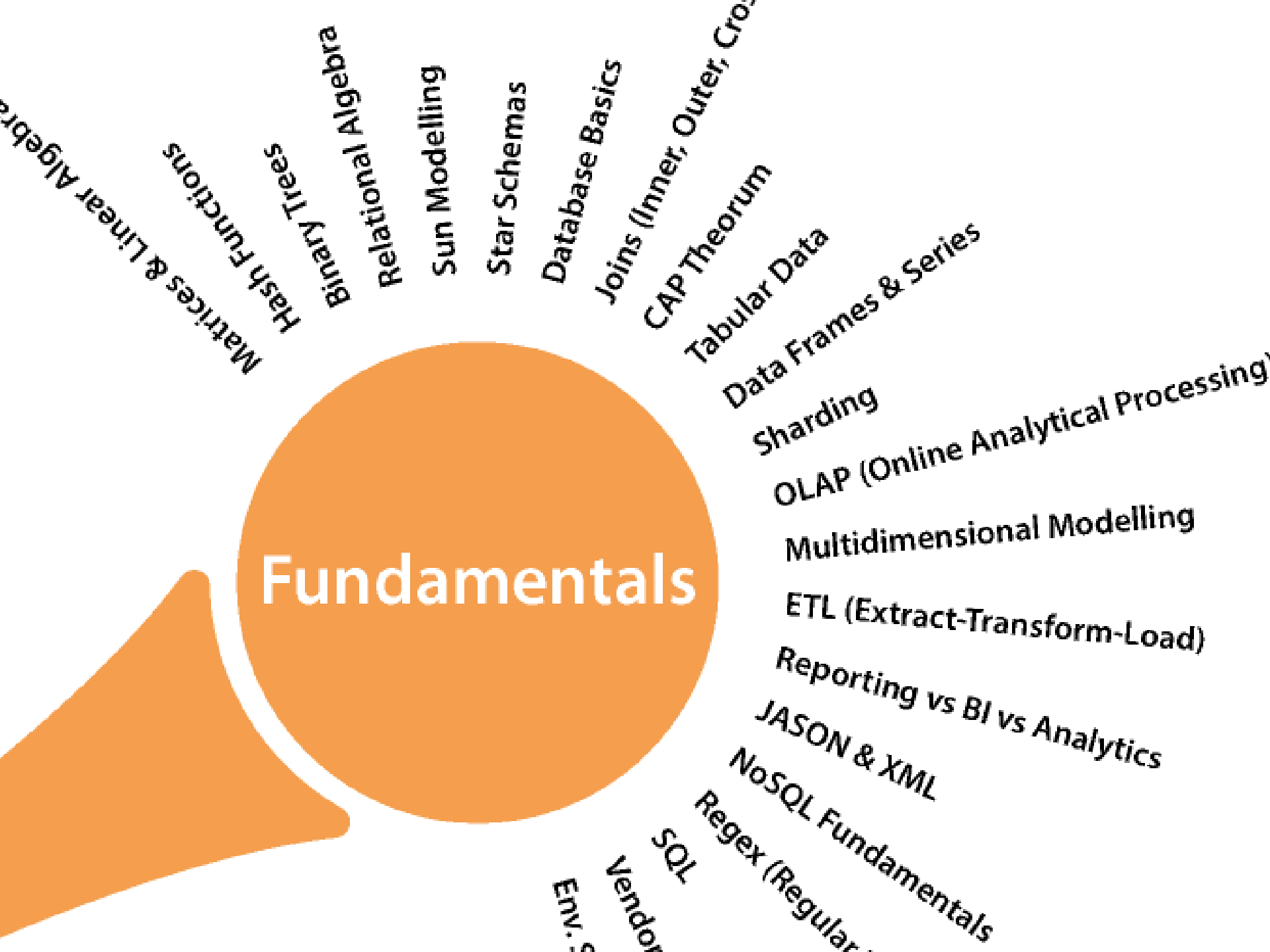
Deep Technical Skills

Communication Skills

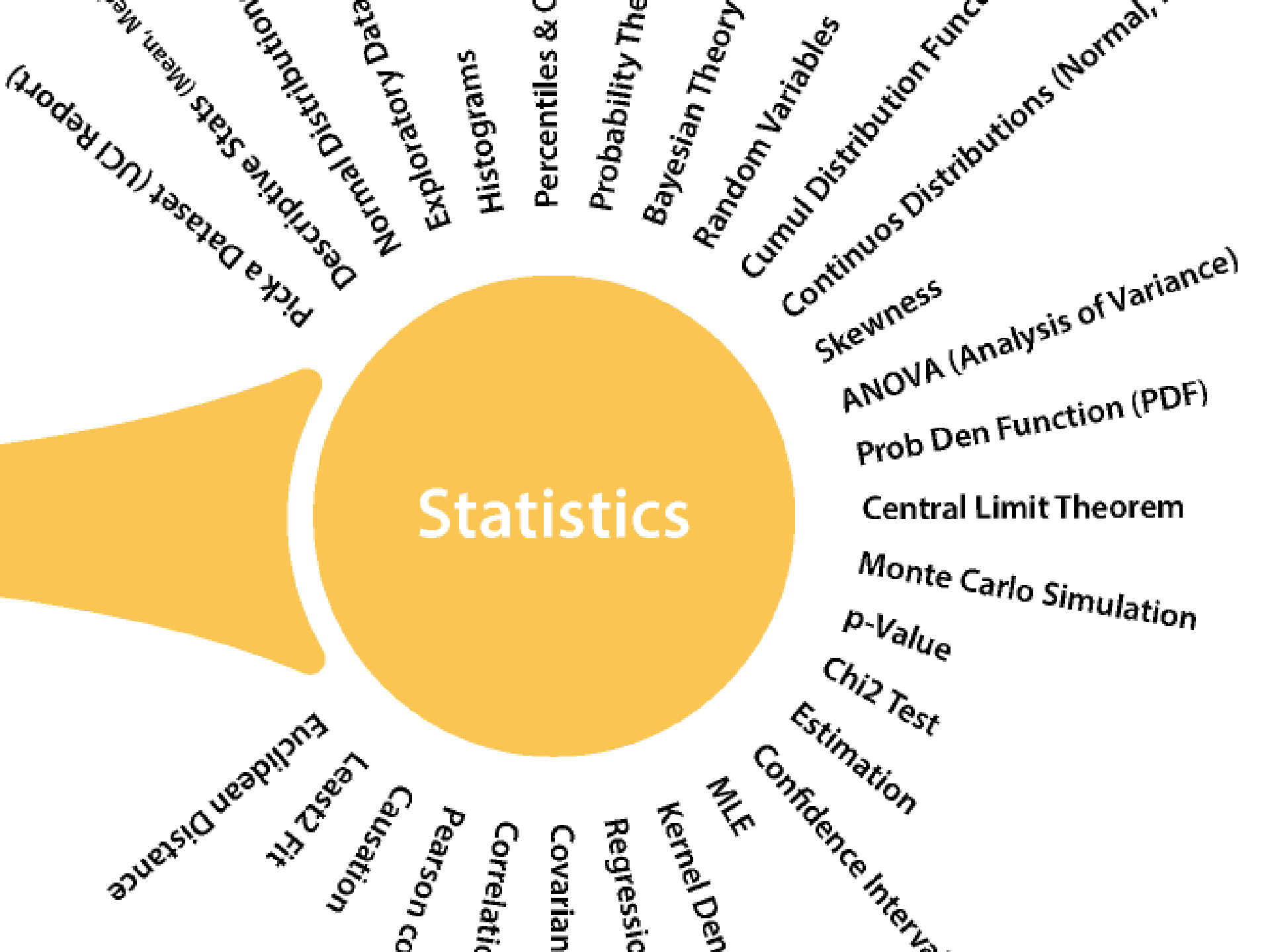
Presentation Skills

Sees beyond the obvious

Fundamentals



Statistics



Programming

Python B

Excel

R Setup / R Stu

R Basics

Expressions

Variables

Vectors

Matrices

Arrays

Factors

Lists

Data Frames

Functional Programming

CSV Data

Raw Data

Subsetting Data

Manipulate Dat

Functions

Factor Analysis

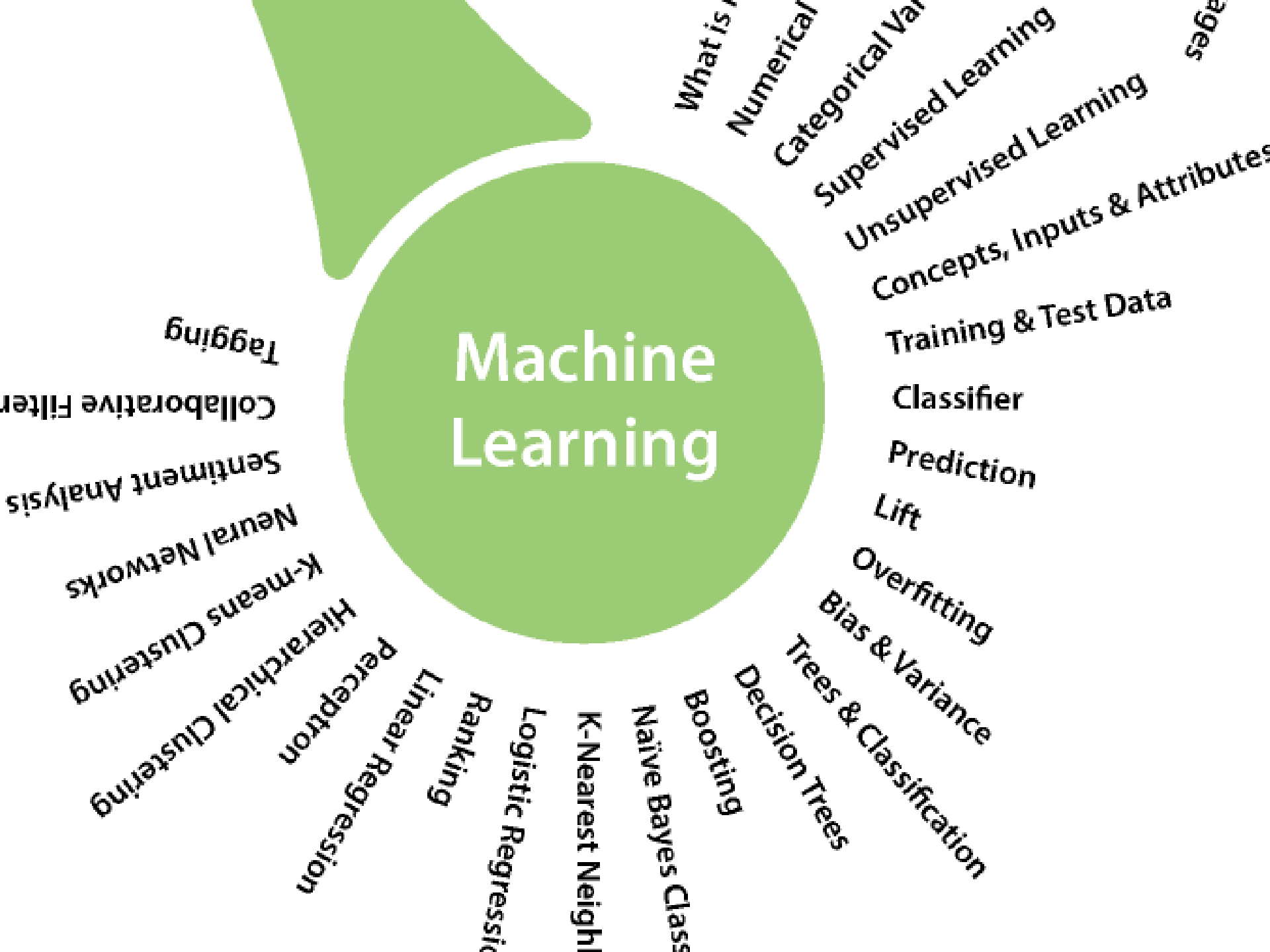
Rapid Packages

Regression

GUI

GUI

Machine Learning



Text Mining / Natural Language Processing

Vocabulary Map

Classify Text

Using NLTK

Using Weka

Using Mahout

Feature Extraction

Market Based Analysis

Association Rules

Support Vector Machines

Term Frequency & Weight

UIMA (Unstructured Information Management Architecture)

Text Analysis

Named Entity Recognition

Corpus

Graphs / Networks

Clustering

Clustering

Visualization

QlikView
Tableau
IBM ManyEyes

InfoVis

D3.js

Decision Tree

Timeline

Survey Plot

Spatial Charts

Line Charts (Bi)

Scatter Plot (Bi)

Tree & Tree Map

Histogram & Pie (Uni)

ggplot2

Data Exploration in R (Hist, Boxplot etc)
Uni, Bi & Multi-variant Visualisation

Guiding

Big Data





Data Ingestion

- Using ETL
- How much Data?
- Google Open Refine
- Data Survey
- Transformation & Enrichment
- Data Fusion
- Data Integration
- Data Sources & Acquisition
- Data Discovery
- Summary of Data Formats

Principal Component Analysis

Stratified Sampling

Sampling

De-noising

Feature Extraction

Binning Sparse Values

Unbiased Estimators

Handling Missing Values

Data Scrubbing

Normalisation

Data Munging

Toolbox

MS Excel with
Java
JavaScript
Python
R, R Studio, Rattle
Weka, Knime, RapidMiner
Hadoop
Spark, Storm
Flume, Scribe, Chukwa
Nutch, Talend, Scraperwiki
Webcrawler, Flume, Sqoop
tm, RWeka, NLTK
RHIPE
D3.js, ggplot2, Shiny
Cassandra
MongoDB
Neo4j



The final step on this journey, is also your first...
There are always new developments and new tools coming
Data Science never finishes, there are always

Keep Looking and Keep Learning!

And if you do find something interesting

to let us know here so we can share it 😊

🌐 www.exploringdatascience.com