



The Present & Future of Data Analytics

Dr. Siddhartha Roy

Goals of this presentation

- Provide Analytics overview
- Discuss various types of Analytical models
- Highlight current and future trends
- Analytics project Lifecycle and Workflow
- Analytics Landscape

Table of Contents

- ❖ Introduction to Analytics
- ❖ Data Analytics Few Facts
- ❖ Analytical Models
- ❖ Analytics Workflow
- ❖ Modeling Journey
- ❖ Job Opportunities
- ❖ Analytics Landscape
- ❖ Benefits of R
- ❖ Logistic Regression
- ❖ Clustering Analysis



What is **Analytics**?

Analytics is defined as a scientific process which helps organizations in:



Answering business questions and taking fact-based decisions



Identifying patterns



Discover hidden relationships



Forecasting and Predictions



Understand market trends



Define customer preferences



What do Experts say...

“ Without data, you are just another person with an opinion”

– W. Edward Deming eminent researcher & scholar

“Data are becoming the new raw material of business.”

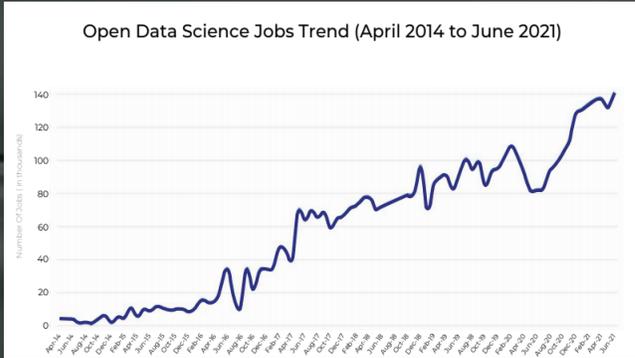
– Craig Mundie Senior Advisor to the CEO at Microsoft

“We're making this analogy that AI is the new electricity. Electricity transformed industries: agriculture, transportation, communication, manufacturing.”

– Andrew Ng Co-founder and Head of Google Brain

Facts about Data Analytics Jobs in India

Job Opportunities

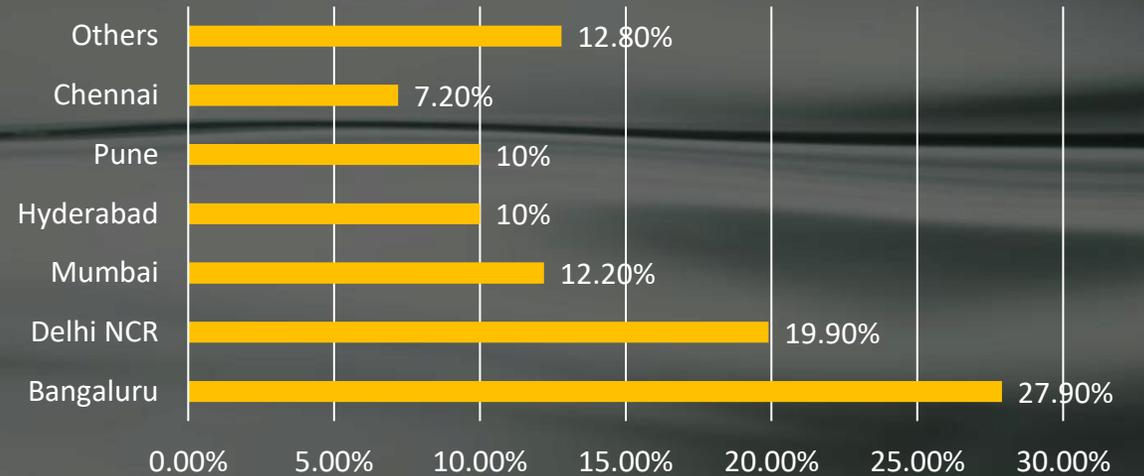


47.1% increase in open jobs as compared to last year. Around 137K open positions as per Jun-21



India contributed to **9.4%** of the total global analytics job openings

Share of Open Jobs by City



Jobs by Sector/Industry



Top Sectors with share of Open Jobs June 2021

1. BFSI **27.8%**
2. Energy & Utilities **17.3%**
3. Media & Entertainment **17.3%**
4. Pharma & Healthcare **16.7%**
5. Ecommerce **11.6%**

Jobs by Tools & Technologies



- Python, R, and Tableau are key skill-sets required in Analytics jobs
- AWS, Azure, & Google Cloud are offering cloud platform jobs
- Sharp increase in demand for Spark & Hadoop skills



Where can we apply **Analytics**?

- Pharma & Health care
- Retail
- Ecommerce
- Logistics
- Agriculture
- Education
- Defense
- Telecom
- Automobile
- Energy & Utilities
- Travel & Hospitality
- Media & Entertainment
and many other..





Types of data **Analytics**



Descriptive

- Trends analysis
- Summarization
- Exploratory data analysis
- Historical summary



Diagnostic

- Root cause analysis, why it happened
- how many times it happened
- Evaluating different scenarios



Predictive

- Predict outcomes
- Forecasting



Prescriptive

- Finding optimal solution
- Identify possible actions



Cognitive

- Advanced form of analytic
- Artificial Intelligence
 - Machine Learning



Analytical Models



Default rate model

What are the influential factors which are associated with default rates?



Event Modeling

Can I predict consistent models that accurately reflect the business world



Collection Analysis

Can I segment my delinquent portfolio basis DPD, Amt overdue, area of needs & past transactions?



Optimization

How to arrive at the optimum premium amount of products to meet profit targets?



Risk Based Pricing

What is optimum premium or price to be charged to customers basis their profiles & transaction patterns?



Fraud Analytics

How can we detect potential fraudulent transactions and prevent it?



Cross-sell model

How can we cross-sell new products to the existing customers?



Performance Analysis and Scorecard

How to develop performance ranking system for insurance agent?



Analytical Models Continued..



Conversion Analysis

Can I segment customer base basis their conversion rate?



Claim Analytics

How can we analyze claims patterns and process claims efficiently?



Portfolio Modeling/Analysis

What is the revenue break up by Customer loyalty segment?



ROI analysis

Which are my different project investment and what is my profitability?



Operational Analytics

Can I take out topmost redundant statements in calling script to reduce AHT?



Response Rate /Activation Modeling

How to optimize promotional cost by predicting high response segment to any promotion?



Spend/Usage Analysis

What is the expenditure data with the purpose for reducing procurement costs, improving efficiency and monitoring compliance?



Retention /Churn Model

Who are the customers that are going to churn in next cycle & what are the corrective measures to avoid the same?



Analytical Models Continued..



Market Basket Analysis

What are the purchasing patterns of the customers



Customer Satisfaction Analysis

Can I derive key factors influencing CSAT score basis verbatim analysis of chat with my customers?



Survey Analytics

How can I make an effective questionnaire for my survey?



Cross-sell analysis

How can I design a cross-sell strategy for product portfolio?



Product Conversion Rate Analysis

Can I segment my URL basis their conversion rate?



Pricing Optimization

How to arrive at the optimum price of products to meet profit targets?



Customer Life-Time value

What is the transition probability of customers patronizing me basis their transaction pattern?

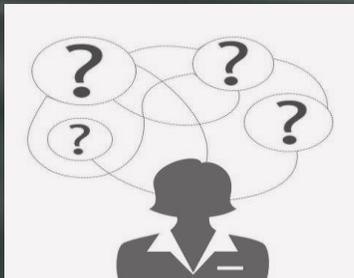


Customer Segmentation

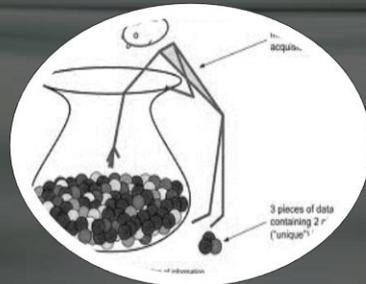
How many discrete customer groups that share similar characteristics basis my profit, demographics etc.



Analytics Workflow



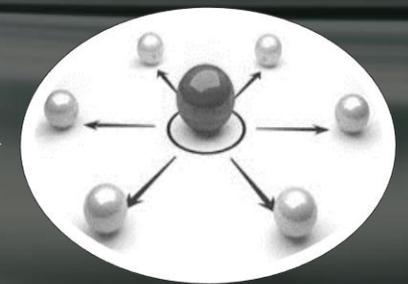
Problem Statement



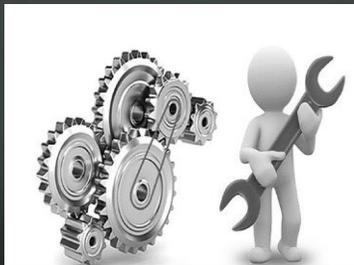
Data Sampling & Consolidation



Data Cleaning



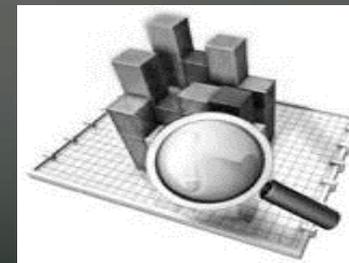
Master Data Repository



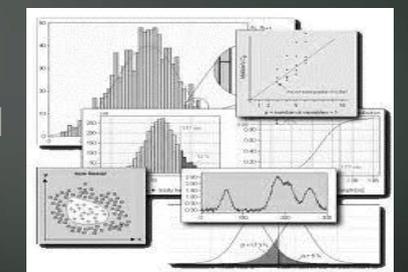
Model Enhancement



Recommendation



Data Analysis and Modeling



Exploratory Data Analysis



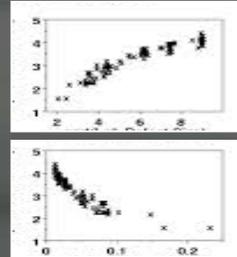
Modeling Journey

LeadKey	Litigation	LossDate	LossLocaf	LossName	MANAME	MovementDate	Movement	Orig	Cur	OrigCur	OrigUnits	OrigUnits	OtherNat	Participant	Reinsured	Reporting	Std	Service	Key
ESTERN UN	N	#####	BELGIUM	Catlin Und		8/3/2014 0:00	AC	334	EUR		2.01E+12		Lead	Catlin Underwrt	BB	Standard			
N	#####	PALES OF	BANCO A	Talbot Und		8/3/2014 0:00	AB	775	USD	BANCO CC	2.01E+12		Lead	MAPFRE	STalbot Underwrt	BB	Standard		
N	7/2/2001	GLUCKESTER	Vibe Synd			8/3/2014 0:00	AC	676	USD	SATYANM	2.00E+12		Lead	LUMBER	Erie Syndicate	NA	Standard		
N	#####	TBA		Catlin Und		8/3/2014 0:00	AB	801	USD	GREGORY	2.01E+12		Lead	MEDAM	Catlin Underwrt	BT	Standard		
ALEXANDR	N	9/4/2010	20TH STREET	NEAR	Novae Sy	8/3/2014 0:00	AI	751	CAD		2.01E+12		Lead	CAROLYN		NA	Standard		
N	#####	S.A. Meac				8/3/2014 0:00	AD	1108	USD	BANKERS	2.00E+12		Lead	THE TRAV	S.A. Meacock	BOO	Technical		
N	#####	MALVINA	IATM ATTA	Chauver S		8/3/2014 0:00	AB	1230	USD	BANCO DE	2.01E+12		Lead	SANCOR	Chauver Syndic	BB	Standard		
ACQUE N	#####	SICILY		Arch Unde		8/3/2014 0:00	AB	799	EUR		2.01E+12		Lead	Arch Underwrt	NA	Standard			
N	---	TBA	FALSE AD	ANV Synd		8/3/2014 0:00	AC	180	USD	GENCOR	2.01E+12		Lead	NURFACE	ANV Syndicate	UC	Standard		
ETROLEUM	N	#####	CANADA	Mitsui Sur		8/3/2014 0:00	AE	751	CAD		2.01E+12		Lead	Mitsui Sumit	NA	Standard			
N	#####	S.A. Meac				8/3/2014 0:00	AC	1108	USD	BANKERS	2.00E+12		Lead	THE TRAV	S.A. Meacock	BOO	Technical		
MENDES N	---	USA		Beasley Fl		8/3/2014 0:00	AG	823	USD		2.01E+12		Lead	Beasley Fur	12	Standard			
N	---	KUWAIT	DUBAI BR	ORE Unde		8/3/2014 0:00	AC	309	GBP	AL AHJ BA	2.01E+12		Lead	KUWAIT	ORE Underwrt	BB	Standard		
N	7/9/2011	SICILY		Arch Unde		8/3/2014 0:00	AD	799	EUR		2.01E+12		Lead	Arch Underwrt	NA	Standard			
N	---	USA		Brit Syndic		8/3/2014 0:00	AE	618	USD		2.01E+12		Lead	Brit Syndicate	16	Standard			
HDES N	#####	14 CAMPBELL	SQA AJ	Novae Sy		8/3/2014 0:00	AE	595	USD		2.01E+12		Lead	Novae Syndic	NA	Standard			
TANISH MEN	---	MANITOBA	CLINIC	Chauver S		8/3/2014 0:00	AN	823	CAD		2.01E+12		Lead	ORE	DATHE		Standard		
CATALON N	#####	MONTHLY	FIRE IN	Vire Synd		8/3/2014 0:00	AD	111	EUR		2.01E+12		Lead	Fire In W		Standard			
NV	N	---	USA	Novae Sy		8/3/2014 0:00	AD	823	USD		2.01E+12		Lead	CLINTON		Standard			
PLUS IMAE N	#####	MZS		GBE Unde		8/3/2014 0:00	AE	301	GBP		2.01E+12		Lead	TBA		Standard			
ACQUE N	---	SICILY		Arch Unde		8/3/2014 0:00	AE	799	EUR		2.01E+12		Lead	Arch Underwrt	NA	Standard			
MCCOLLUM	N	---	---	Brit Syndic		8/3/2014 0:00	AC	375	USD		2.01E+12		Lead	MORISON		Standard			
COOP N	#####	FLORIDA	Canopus			8/3/2014 0:00	AD	446	USD		2.01E+12		Lead	Omega Unde	1A	Standard			

80 %

20 %

TRANSFORMATION



FEATURE SELECTION

LeadKey	Litigation	LossDate	LossLocaf	LossName	MANAME	MovementDate	Movement	Orig	Cur	OrigCur	OrigUnits	OrigUnits	OtherNat	Participant	Reinsured	Reporting	Std	Service	Key	
ESTERN UN	N	#####	BELGIUM	Catlin Und		8/3/2014 0:00	AC	334	EUR		2.01E+12									
N	#####	PALES OF	BANCO A	Talbot Und		8/3/2014 0:00	AB	775	USD	BANCO CC	2.01E+12									
N	7/2/2001	GLUCKESTER	Vibe Synd			8/3/2014 0:00	AC	676	USD	GATEWAY	2.00E+12									
N	#####	TBA		Catlin Und		8/3/2014 0:00	AB	801	USD	GREGORY	2.01E+12									
ALEXANDR	N	9/4/2010	20TH STREET	NEAR	Novae Sy	8/3/2014 0:00	AI	751	CAD		2.01E+12									
N	#####	S.A. Meac				8/3/2014 0:00	AD	1108	USD	BANKERS	2.00E+12									
N	---	MALVINA	IATM ATTA	Chauver S		8/3/2014 0:00	AB	1230	USD	BANCO DE	2.01E+12									
ACQUE N	#####	SICILY		Arch Unde		8/3/2014 0:00	AB	799	EUR		2.01E+12									
N	---	TBA	FALSE AD	ANV Synd		8/3/2014 0:00	AC	180	USD	GENCOR	2.01E+12									
ETROLEUM	N	#####	CANADA	Mitsui Sur		8/3/2014 0:00	AE	751	CAD		2.01E+12									
N	#####	S.A. Meac				8/3/2014 0:00	AC	1108	USD	BANKERS	2.00E+12									
MENDES N	---	USA		Beasley Fl		8/3/2014 0:00	AG	823	USD		2.01E+12									
N	---	KUWAIT	DUBAI BR	ORE Unde		8/3/2014 0:00	AC	309	GBP	AL AHJ BA	2.01E+12									

If Not Good

Further refinement and treatment of outliers and variables

DIAGNOSTICS CHECK

Coefficient of Determination → $R^2 = \frac{SSR}{SST} = 1 - \frac{SSE}{SST}$

Sum of Squares Total → $SST = \sum (y - \bar{y})^2$

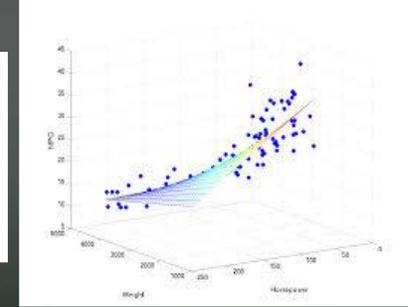
Sum of Squares Regression → $SSR = \sum (y' - \bar{y})^2$

Sum of Squares Error → $SSE = \sum (y - y')^2$

Adjusted $R^2 = 1 - (1 - R^2) \left(\frac{N-1}{N-k-1} \right)$

AIC, Cook's distance, Studentized residuals, Standard Residuals

MODELING



PREDICTION ON TEST DATA

Key	LeadKey	Litigation	LossDate	LossLocaf	LossName	MANAME	MovementDate	Movement	Orig	Cur	OrigCur	OrigUnits	OrigUnits	OtherNat	Participant	Reinsured	Reporting	Std	Service	Key	
1	14	CAMPBELL	SQA AJ	Novae Sy		8/3/2014 0:00	AE	595	USD		2.01E+12										
2	14	CAMPBELL	SQA AJ	Novae Sy		8/3/2014 0:00	AE	595	USD		2.01E+12										
3	14	CAMPBELL	SQA AJ	Novae Sy		8/3/2014 0:00	AE	595	USD		2.01E+12										
4	14	CAMPBELL	SQA AJ	Novae Sy		8/3/2014 0:00	AE	595	USD		2.01E+12										
5	14	CAMPBELL	SQA AJ	Novae Sy		8/3/2014 0:00	AE	595	USD		2.01E+12										
6	14	CAMPBELL	SQA AJ	Novae Sy		8/3/2014 0:00	AE	595	USD		2.01E+12										
7	14	CAMPBELL	SQA AJ	Novae Sy		8/3/2014 0:00	AE	595	USD		2.01E+12										
8	14	CAMPBELL	SQA AJ	Novae Sy		8/3/2014 0:00	AE	595	USD		2.01E+12										
9	14	CAMPBELL	SQA AJ	Novae Sy		8/3/2014 0:00	AE	595	USD		2.01E+12										
10	14	CAMPBELL	SQA AJ	Novae Sy		8/3/2014 0:00	AE	595	USD		2.01E+12										

If Good

If accuracy in Train Data ≈ Test Data & Diagnostics are acceptable

PREDICTION ON NEW SET OF DATA

MODEL MONITORING



Analytics Project – Key Phases

1



**Opportunity
Identification**

2



**Proof of
Concept**

3



Delivery

4



**Automation,
Optimization
and Efficiency**

5



**Large scale
Implementation**

6



Productization

7



**Maintenance &
Enhancement**

8



**Transformation
& Innovation**



Analytics Landscape

 Business Intelligence	 Analytics/ Data Science	 Big Data Analytics	 Artificial Intelligence & Machine Learning
<ul style="list-style-type: none">✓ Data cleaning & management✓ Data reports✓ Visual Reporting✓ Creating dashboards	<ul style="list-style-type: none">✓ Data Cleaning and management✓ Exploratory data analysis✓ Statistical Hypothesis Testing✓ Predictive Modeling✓ Forecasting✓ Identifying patterns✓ Optimization✓ Model Diagnostics & Accuracy Testing	<ul style="list-style-type: none">✓ Data cleaning & management✓ Huge amount of data✓ Social networking sites✓ Sensors✓ Structured, semi-structured, unstructured and✓ Uses parallel processing✓ Identify fraud analytics✓ Weather forecasting✓ Used in designing products✓ Geographic analysis✓ Customer Experiences	<ul style="list-style-type: none">✓ Data Cleaning & management✓ Neural networks✓ Natural Language Processing✓ Natural language Generation✓ Image Analysis and Modeling✓ Sensory Perception✓ Cognition✓ Computer Vision✓ Speech Analytics✓ Robotics✓ Expert System✓ Recommendation Engine✓ Knowledge Engineering

Job Opportunities

- ❖ Data Analyst
- ❖ Analytics Associate
- ❖ Data Scientist
- ❖ Data Engineer
- ❖ Researcher
- ❖ Business Analyst
- ❖ Statistical Analyst/Statistician
- ❖ Analytics Consultant
- ❖ Analytics Specialist
- ❖ Data Analytics Architect
- ❖ Data Solution Architect
- ❖ Machine Learning Engineer
- ❖ Data Analytics Manager
- ❖ Marketing Research Analytics
- ❖ Operations Analytics Manager
- and many more...



Machine Learning





Types of Machine Learning

Supervised Learning

Labeled data is being used for model training both inputs & outputs

- Classification {Categorical Target variable}
- Regression {Continuous Target variable}



Semi-Supervised Learning

Smaller set of labeled data with a large amount of unlabeled data during model training

- Classification {Categorical Target Variable}
- Clustering {Categorical Target Variable}



Unsupervised Learning

Unlabeled input data, identify structure and patterns from input data

- Clustering {Target variable not available}
- Association {Target variable not available}



Reinforcement Learning

Identify an optimal way to accomplish a particular goal

Learns a series of action

- Exploration {Target variable not pre-defined}





Types of Machine Learning

SUPERVISED LEARNING

Classification:

Identifying set of categories sub-populations an observation or observations belongs to

- Logistic Regression
- Decision Trees
- SVM Support-Vector Machines
- Naive Bayes
- K-NN K-Nearest Neighbour

Regression:

Estimate the relationships between a dependent variable often called the 'outcome' or 'response' variable and one or more independent variables

- Linear Regression
- Polynomial Regression
- Ridge/Lasso Regression

SEMI-SUPERVISED LEARNING

Speech Analytics:

Internet Content Classification:

Protein Sequence classification:

[Labelling is time intensive]

UNSUPERVISED LEARNING

Clustering:

Grouping a set of objects in such a way that objects in the same group

- K-means
- Mean-Shift
- Fuzzy C-means

Association Analysis:

Association analysis is the task of finding interesting relationships in large datasets.

- Apriori
- FP-Growth

Dimension Reduction:

Reduction in the number of features or number of observations or both

- Principal Component Analysis PCA
- Partial Least Square Regression PLSR
- Multidimensional Scaling MDS
- Principal Component Regression PCR
- Linear Discriminant Analysis LDA

REINFORCEMENT LEARNING

Q-learning:

Values based learning algo

- Industrial Automation
- Gaming
- Recommendation system
- Chemical reactions
- Self driving cars



Introduction to R Programming

- ❖ Founded by Professors **Robert Gentleman & Ross Ihaka University of Auckland**
- ❖ Open-source programming language
- ❖ A whole bouquet of packages
- ❖ High quality plots
- ❖ Versatile reporting modules
- ❖ Free tutorials & learning materials
- ❖ Ease of data preparation
- ❖ Packages and its applications
 - tidyverse <Data Science>
 - ggplot <Data visualization>
 - dplyr <Data frames>
 - tidyr <Data management>
 - purrr <Functions & Vectors>
 - stringr <String operations>

Packages in R

Library	Description
<code>library(readxl)</code>	To load data: These packages help you read and write Microsoft Excel files from R.
<code>library(dplyr)</code>	To manipulate data: Provides function subsetting, summarizing, rearranging, and joining together data sets
<code>library(tidyr)</code>	To manipulate data: Provide tools for changing the layout of your data sets. Uses the gather and spread functions to convert data into the tidy format
<code>library(gtools)</code>	To manipulate data: Provide functions to assist in R programming like calculate the logit and inverse logit transformations, test if a value is missing, empty or contains only NA and NULL values, define macros, sort strings containing both numeric and character components, enumerate permutations and combinations, etc.
<code>library(class)</code>	Functions for Classification: Various functions for classification, including k-nearest neighbor, Learning Vector Quantization and Self-Organizing Maps.
<code>library(recipes)</code>	Preprocessing and Feature Engineering Steps for Modeling: Prepares data for modeling by providing extensible framework for pipeable sequences of feature engineering steps Statistical parameters can be estimated from an initial data set and then applied to other data sets. The resulting processed output can then be used as inputs for statistical or machine learning models.
<code>library(caret)</code>	To model data: Tools for training regression and classification models
<code>library(gmodels)</code>	To model data: Programming tools for model fitting.
<code>library(magrittr)</code>	Data wrangling: A Forward-Pipe Operator for R that provides a mechanism for chaining commands with a new forward-pipe operator, %>%. This operator will forward a value, or the result of an expression, into the next function call/expression. There is flexible support for the type of right-hand side expressions. This decrease development time and improve readability and maintainability of code.
<code>library(broom)</code>	Convert Statistical Objects into Tidy Tibbles: Easy to report results, create plots and works with large numbers of models at once. tidy() summarizes information about model components such as coefficients of a regression. glance() reports information about an entire model, such as goodness of fit measures like AIC and BIC. augment() adds information about individual observations to a dataset, such as fitted values or influence measures
<code>library(plotROC)</code>	To visualize data: Generate Useful ROC Curve Charts for Print and Interactive Use
<code>library(ggplot2)</code>	To visualize data: Package for making graphics and uses the grammar of graphics to build layered, customizable plots
<code>library(ROCR)</code>	Visualizing the Performance of Scoring Classifiers: Tool for creating cutoff-parameterized 2D performance curves by freely combining two from over 25 performance measures



Logistic regression

- ❖ Supervised algorithm used for classification
- ❖ Used for binary classification & can be extended to multi-class classification
- ❖ Its equivalent to Linear regression using a sigmoid function
- ❖ Models the probability of discrete outcome
- ❖ Computes the decision boundaries among the classes
- ❖ Used extensively in industry
- ❖ Logistic function is represented as

$$F(x) = \frac{1}{1 + e^{-(\beta_0 + \beta_1 x)}}$$



Background, Data & Objective

Background:

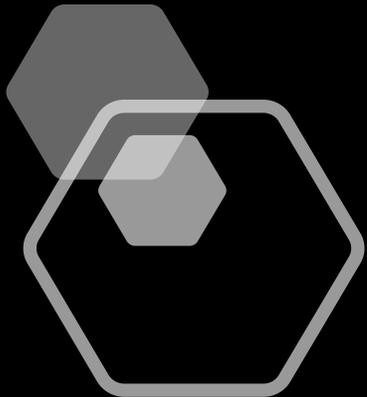
- ❖ Heart diseases are leading causes of morbidity and mortality
- ❖ Around 17.9 million lives are lost, 85% deaths are due to heart attack & stroke
- ❖ Can be attributed to various states – coronary artery diseases, heart rhythm issues, & heart defects
- ❖ Cardiovascular disease also refers to narrowed or blocked blood vessels which may result in heart attack, angina or stroke
- ❖ A lot of studies related to clinical data analysis involves prediction of cardiovascular diseases
- ❖ It is critical to detect cardiovascular disease as early as possible

Data:

- ❖ The dataset used for this case study is Cleveland heart disease data set
- ❖ Original dataset had 76 features but only 14 are chosen here

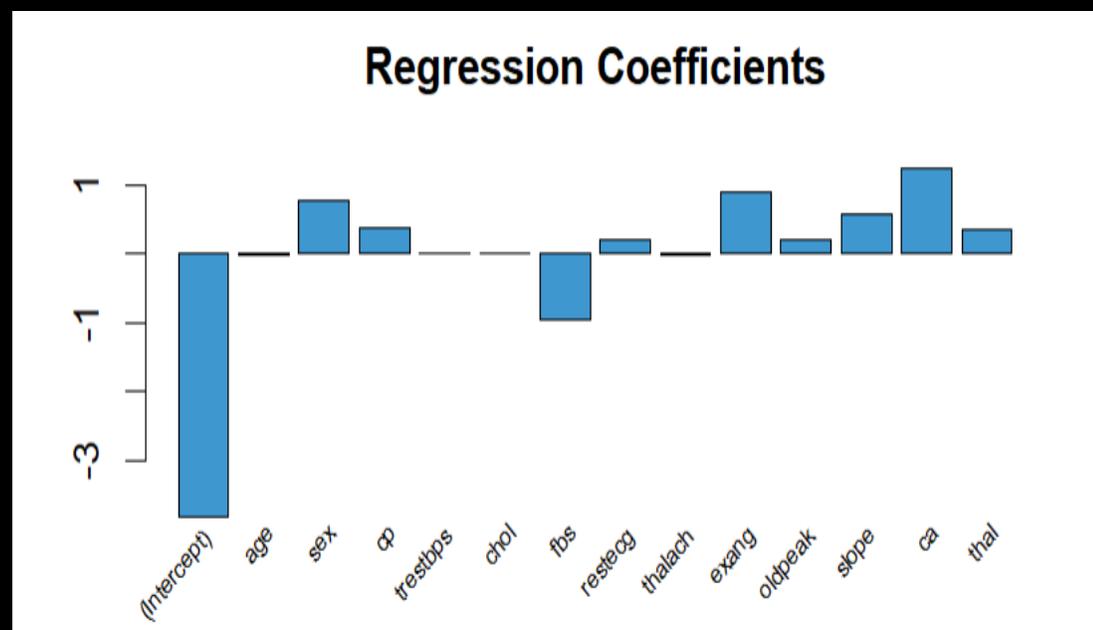
Objective:

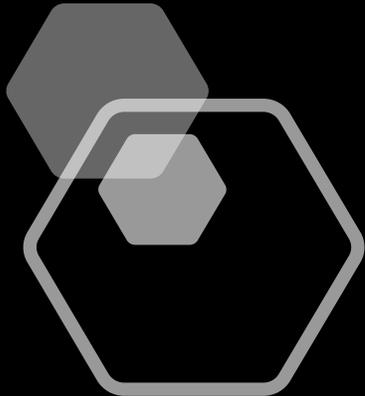
- ❖ Build a modeling solution which assesses the impact of various factors on heart disease



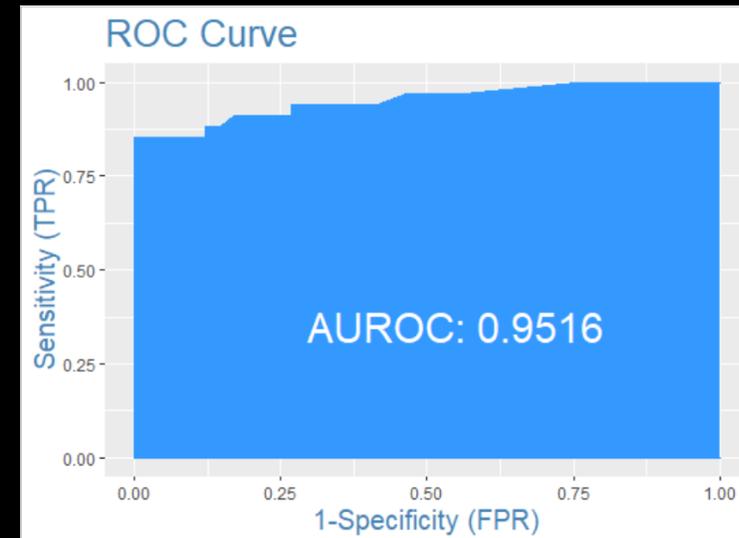
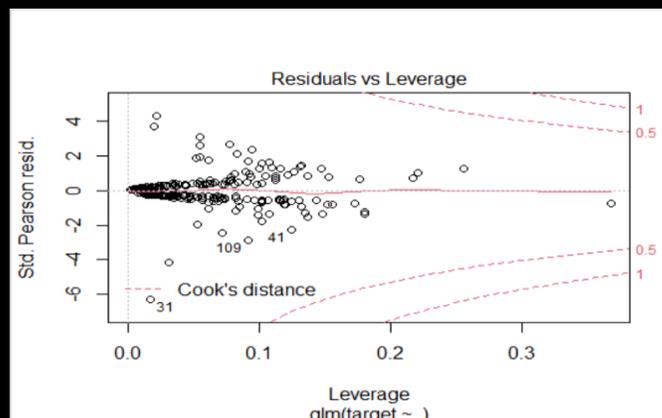
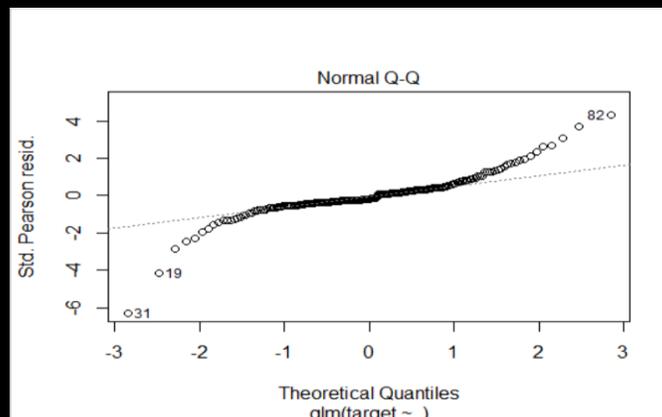
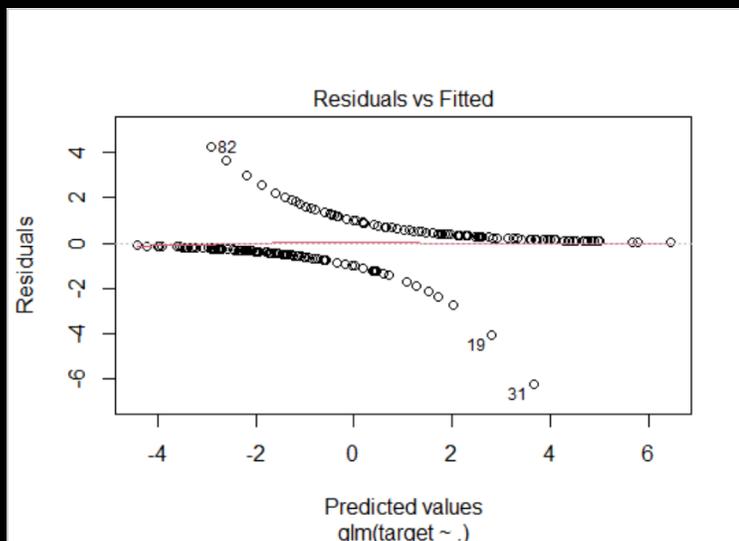
Logistic Regression Results

	Estimate	Std..Error	z.value	Pr...z..
(Intercept)	-3.8139086594	3.210250591	-1.1880408	2.348173e-01
age	-0.0271278536	0.028565033	-0.9496875	3.422711e-01
sex	0.7824952931	0.550298666	1.4219466	1.550418e-01
cp	0.3830443844	0.207779892	1.8435104	6.525453e-02
trestbps	0.0162952430	0.012314322	1.3232757	1.857437e-01
chol	0.0006635756	0.004316101	0.1537442	8.778114e-01
fbs	-0.9638660088	0.586327017	-1.6439052	1.001958e-01
restecg	0.2023828626	0.208143136	0.9723254	3.308887e-01
thalach	-0.0198482303	0.011195339	-1.7729013	7.624507e-02
exang	0.8932838663	0.452713203	1.9731783	4.847525e-02
oldpeak	0.2181737222	0.229754532	0.9495949	3.423182e-01
slope	0.5817356136	0.390271310	1.4905928	1.360684e-01
ca	1.2596492928	0.294218497	4.2813396	1.857716e-05
thal	0.3583134114	0.111283940	3.2198124	1.282745e-03





Logistic Regression Results (contd..)



Cluster Analysis





Clustering Analysis

- Unsupervised algorithm
- Used in multiple domains- documents grouping, understanding customer behavior, anomaly detection, etc.
- Segregates the data into subsets with high intra subset similarity and inter-subset dissimilarity
- Few clustering approaches are –
 - Hierarchical clustering,
 - Non- Hierarchy clustering
 - Density based clustering,
 - Partition clustering,
 - Grid based clustering,
 - Correlation clustering, etc.
- K-means is one of the most common algorithm and uses Euclidean distance measures to create optimum clusters



Types of Clustering Techniques

Hierarchical Clustering

{Divides a data set into a sequence of nested partitions – bottom up & top down}

Agglomerative: Clustering starts with every single object in a single cluster. It merges to the closest pair of clusters as per similarity criteria.

- The Single-link Method
- The Complete Link Method
- The Group Average Method

Divisive: Clustering starts with all objects in one cluster and splits large cluster into small pieces

- DIANA (Divisive Analysis of hierarchical clustering)

Partitioning Clustering

{Classifies data into non-overlapping subsets by selecting centroids}

Fuzzy: Associate each data point with each cluster using membership function

- Fuzzy k-means
- Fuzzy k-modes
- The c-means Method

Centre-based: Cluster have convex shapes and each cluster is represented by a centre. Its tries to minimize its objective function.

- k-means
- k-Medoids
- k-modes

Graph- based clustering

{Clustering a set of graphs & clustering nodes/edges of a single graph}

Used in image segmentation and complex network analysis

- Chameleon
- CACTUS
- ROCK

Grid- based Clustering

{Create a grid structure, and the comparison is performed on grids}

Used for a multi-dimensional data set.

- CLIQUE
- STING
- CLARANS
- BIRCH

Density- based Clustering

{These clusters are areas of higher density than the remainder of the data set}

Uses SCAN method and requires density parameter as termination condition. It is used to identify clusters of arbitrary size and manage noise in data clusters

- Density-Based Spatial Clustering of Applications with Noise (DBSCAN)
- BRIDGE (integrates k-means and DBSCAN)

Model- based Clustering

{Algorithms designed for modelling unknown distribution as a mixture of simpler distributions}

Assumes that the data were generated by a model and tries to recover the original model from the data

- Dirichlet
- COOLCAT
- STUCCO



Clustering Analysis (R code)

K means (X, C, N)

X: data matrix

C: distinct number of cluster centers

N: seed value