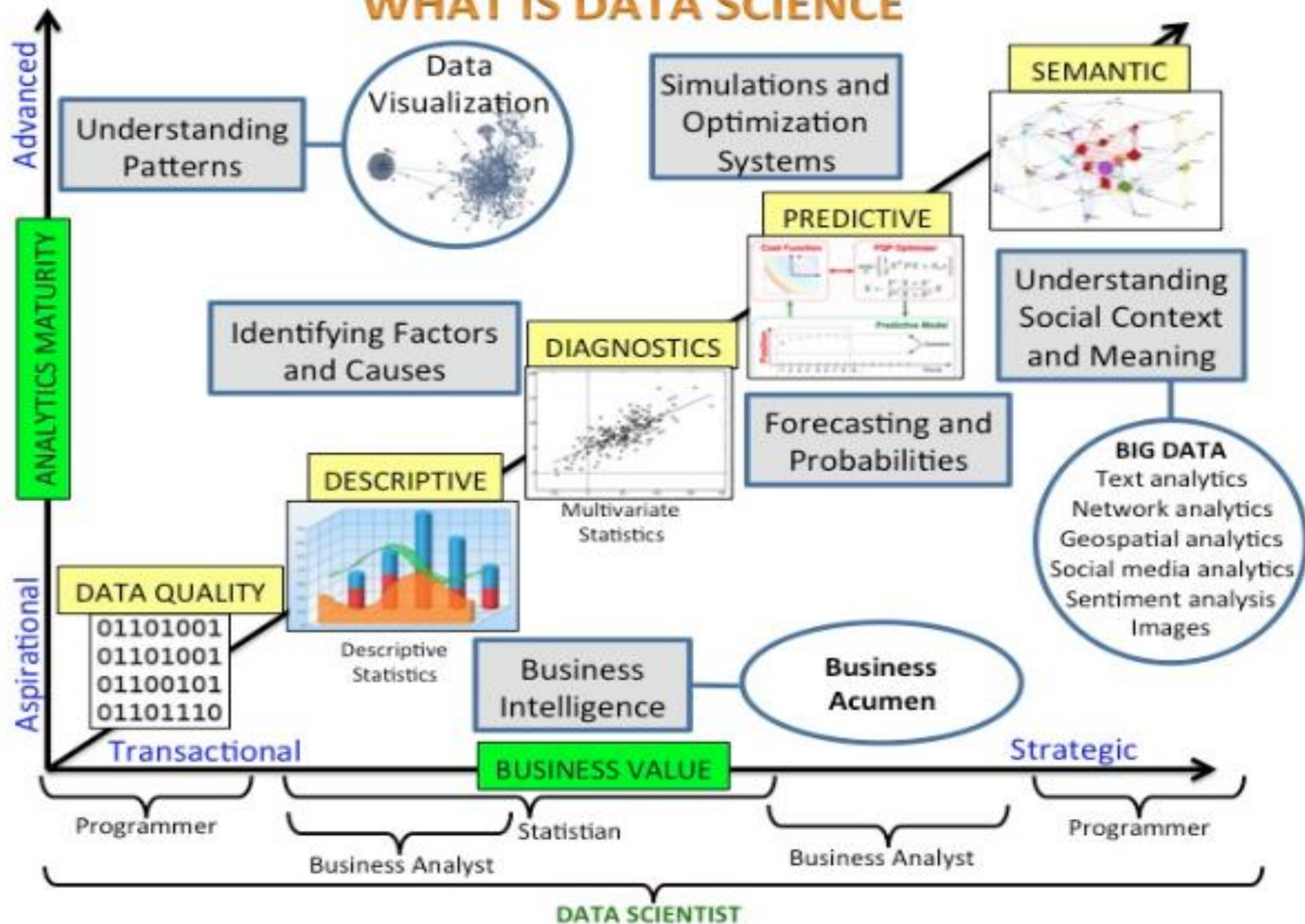


WHAT IS DATA SCIENCE

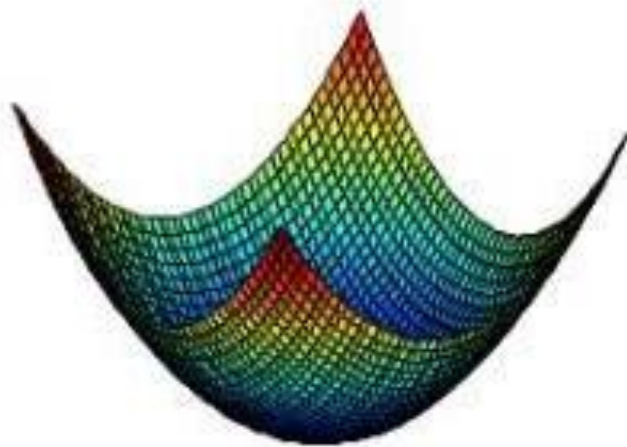


Math for Data Science

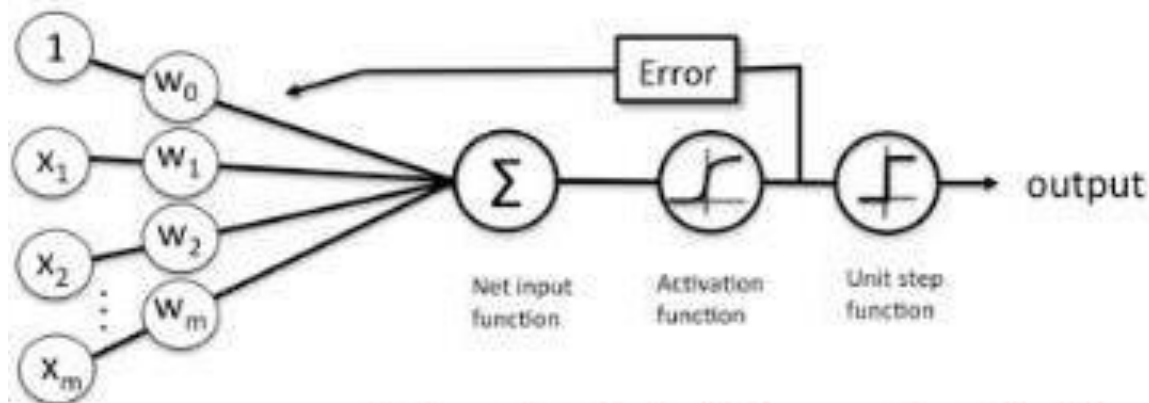


Deep learning vs Machine learning in One Picture

You

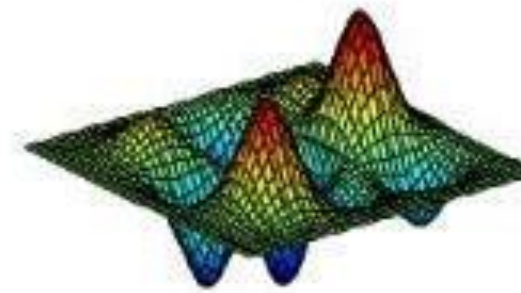


- Unique optimum: global/local.



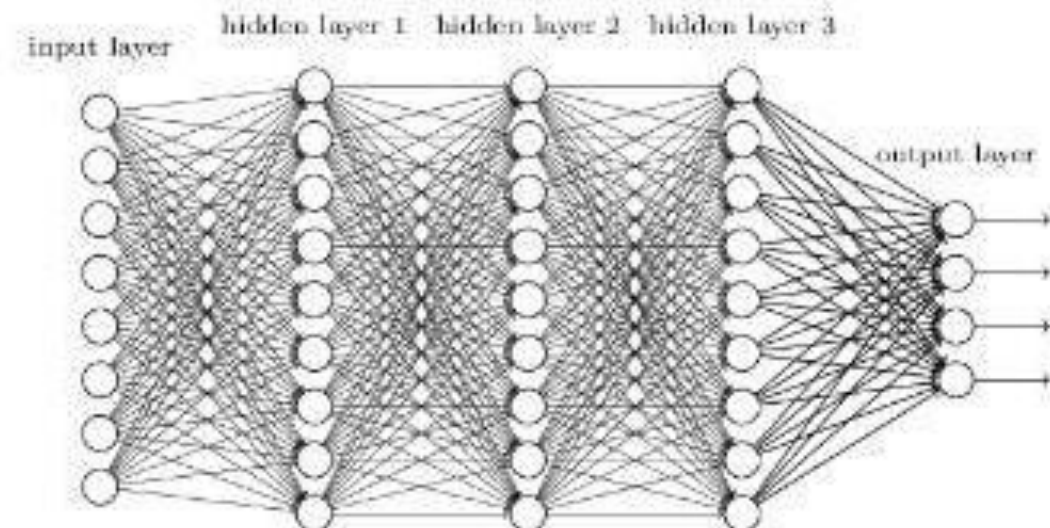
Schematic of a logistic regression classifier.

The guy she tells you
not to worry about



- Multiple local optima
- In high dimensions possibly

Deep neural network



Which Software
Should I Choose?

Best for:

Availability

Easy to learn?

Advantages

Disadvantages

Python

General programming;
Data analysis; Deep
learning; Repeated tasks

Free, open source

Yes, especially for
software engineers

Easy to deploy; General
purpose language; Widely
used by corporations

Requires rigorous testing

R

Statistical analysis; Data
analysis; Single passes of
data

Free, open source

Steep learning curve;
Relatively easier if no
prior coding experience

Minimal coding required
for statistical models

Very statistics oriented;
Not a general-purpose
program

SAS

Statistical analysis; Data
analysis

Paid (free for university
edition); Closed source

Yes, especially if you
already know SQL

Highly reliable, secure
and stable

Relatively expensive

SQL

Database manipulating,
updating, querying;
Extracting, wrangling data

Open and closed source
versions available (free
and paid)

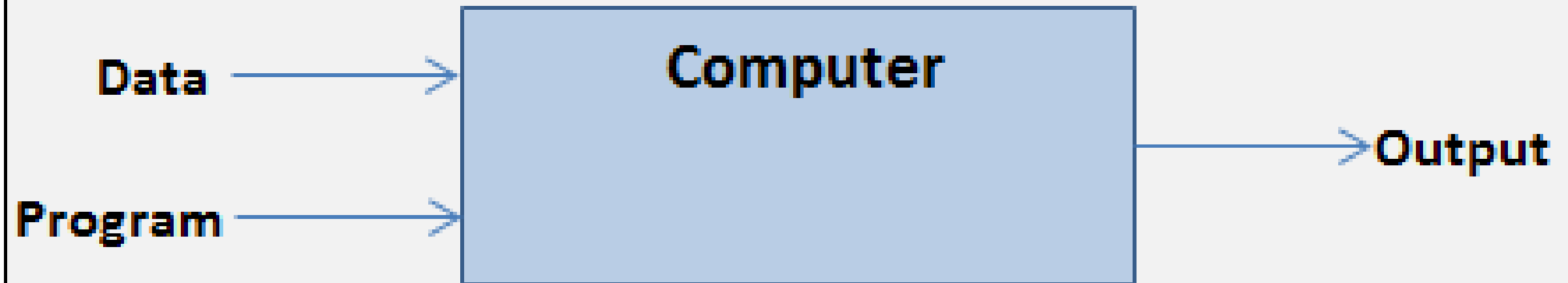
Relatively easy for basic
level; Learning curve for
more complex tasks

Very readable

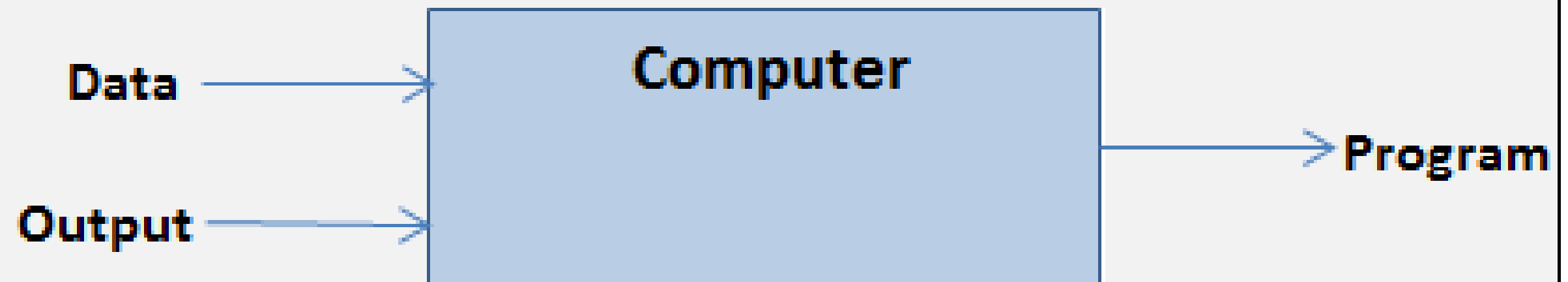
Not general purpose: very
specific, limited capability

Traditional Programming vs ML in One Picture

Traditional Programming



Machine Learning



	Prediction: <i>Predicting Y from X</i>	Inference: <i>Understand relationship between X and Y</i>
Goal	Develop a "best" model (considering all predictors) to predict Y with high accuracy, low error.	Estimate an association between an outcome variable and a predictor variable (while adjusting for confounders).
Answers the question	How can I accurately predict new data points?	What do the relationships between the variables mean?
Example	What mortality levels does the model predict given a certain income and education level?	Which has the biggest impact on mortality: income or education?

I Want to Predict...

Method

Subtypes

**Discrete-valued,
Categorical Outputs**

Example: Given eating habits, predict obese or not obese.

Obese

Not
Obese

Classification

kNN

Decision Trees

Perception
Classifier

**Real-Valued
Outputs**

Example: Given eating habits, predict weight.

Regression

Linear Regression

Ridge Regression

Lasso Regression

100
lbs

200
lbs

220
lbs

150
lbs

110
lbs

241
lbs

1. State the Null Hypothesis

“Null” means that it’s a **commonly accepted fact** that you are working to *nullify*.

Example: the null might be “Data scientists earn an average of \$113,309”

$$H_0: \mu = 133,309$$

2. State the Alternate Hypothesis

The **alternative to the null**.

For example, you think that data scientists earn a lot less than 100k, and you’re going to set out to prove your theory.

$$H_1: \mu < 133,309$$

3. Choose Your Alpha Level

The alpha level α is the **probability of making the wrong decision** when the null hypothesis is true.

Not sure what alpha level to set? 5% is very common.

$$\alpha = 0.05$$

4. Collect Your Data

You may already have the data at hand.

If not, you **have four main choices**:

- Census
- Sample survey
- Experiment
- Observational survey



5. Calculate the Test Statistic

Which test statistic you use depends on which test you’re running. **Choosing the right test** is probably the biggest challenge in hypothesis testing.

Four of the most common:

- Z-Statistic
- T-Statistic
- F-Statistic (ANOVA)
- Chi-Square Statistic

6. Make Accept / Reject Regions

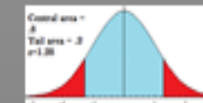
The **rejection region** tells you at what values you should reject the null hypothesis.

The **critical value** (the line that separates accept from reject) is calculated with statistical tables (or software, based on those tables).

7. Accept or Reject the Null Hypothesis

The basic idea is that that your test statistic is going to fall into one of two areas: *accept* or *reject*.

In this graph, you would **reject the null hypothesis** if your test statistic falls into the red area.



Type of Distribution

Number of trials, n

Probability, p

With or without replacement?

When to use it?

Example

Binomial

How many Successes in n trials?

Fixed n

Known p . Probability of Success is constant from trial to trial

With replacement

You know the *exact* probability of an event happening; you want to find the probability of that event happening k times out of n .

Number of defects in a box of 1,000 factory produced widgets

Poisson

Good for rare events (large n , small p)

Unknown n (it is a random variable) and potentially infinite.

Unknown p for each trial (but known average p).

With replacement

You know the *mean* probability of an event and you want to find the probability of n events happening.

Number of innocent people convicted of a crime.

Hypergeometric

Use for small populations, without replacement

Fixed n

Probability for each trial changes (because of no replacement).

Without replacement

Samples are small compared to the population ($n \geq 5\%$). Binomial is easier and provides a good approximation if you have a large population.

30 people (14W, 16M) apply for two jobs. What is the probability both positions are filled by women?

Normal Distribution

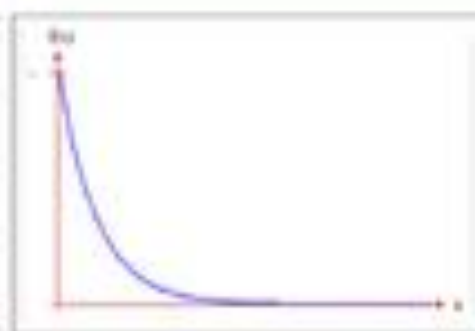
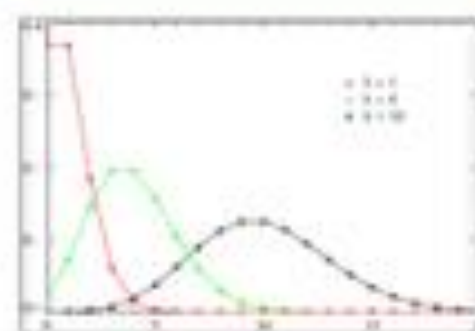
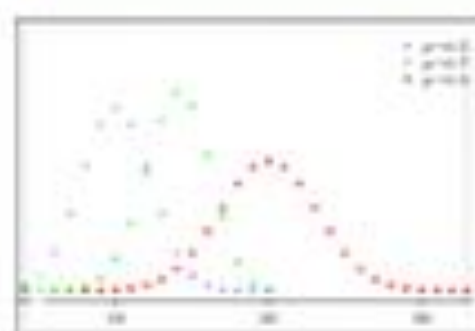
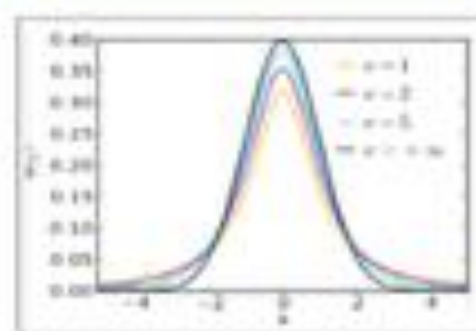
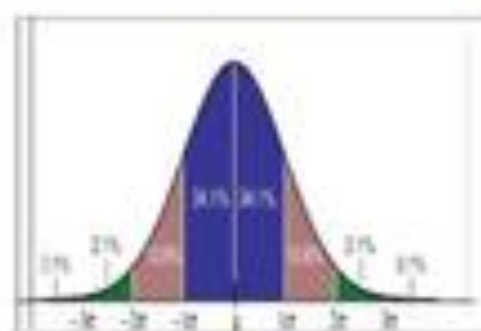
Student's T-Distribution

Binomial Distribution

Poisson Distribution

Exponential Distribution

What does it look like?



Defining Characteristics

Distinctive Bell Shape

Shorter, fatter than the normal distribution.

Two outcomes: Success/Failure

Various shapes, but valid only for integers on the x-axis.

Models Time Between Events

Example of When to Use It

Modeling natural phenomena (height, weight, IQ, test scores etc.)

When you have small samples or don't know the population variance (σ^2).

Coin Toss Probability (Heads, Tails)

Gives probability of number of events in a fixed interval.

"How much time will go by before a major hurricane hits the Atlantic Seaboard?"

Example of DS Application

Least squares fitting or propagation of uncertainty.

Unknown σ^2 is common in real life data, so you'll have to use the T instead of the normal in that case.

Anywhere where binary (yes/no, black/white, vote/don't vote) data is used.

Anywhere there is a waiting time between events.

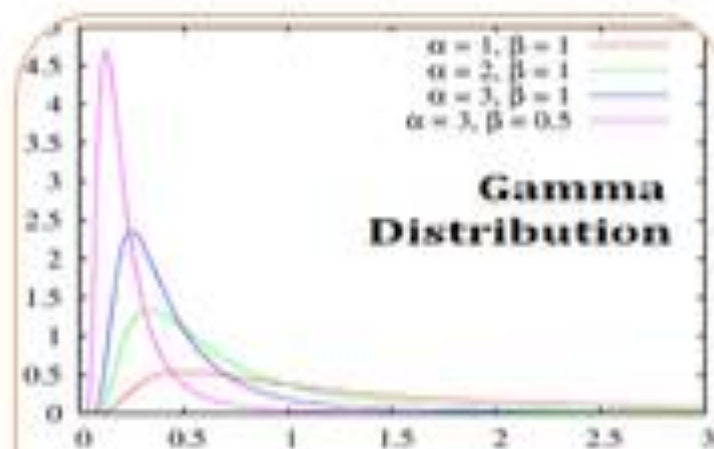
Building continuous-time Markov chains.

Parameters

Overview

Common Uses

When to Use

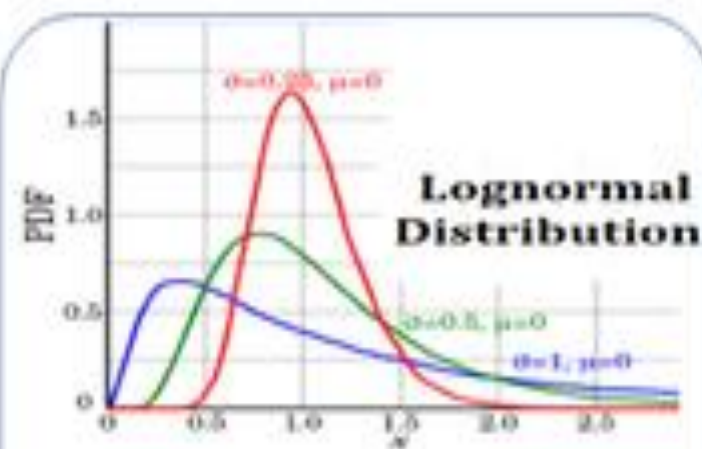


Rate (β); Shape (α)

A family that includes the exponential distribution

Event waiting times when the event process isn't completely random (Mun, 2008)

Most often used as time to a r th hit in a Poisson process

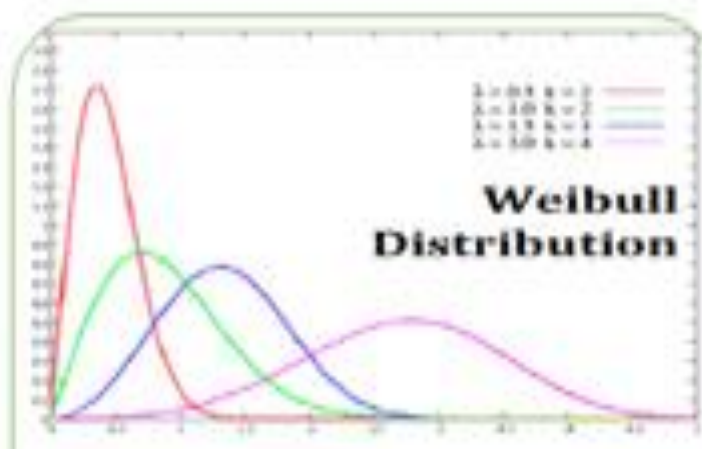


Location (μ); Shape (σ)

Exponential (anti-log) of the normal distribution

Model lifetimes in survival analysis/reliability (Bromideh, 2012)

Use when data has small bumps/perturbations that additively affect the result



Location (μ); Shape (γ); Scale (α)

A natural extension of the exponential distribution

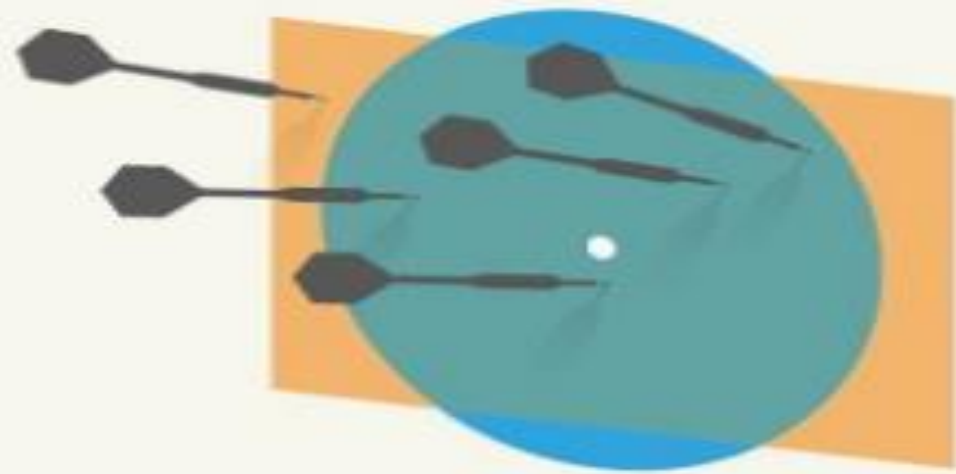
Describe the distribution of lifetime data (Bromideh, 2012) and the time to wait for the next event.

Use when the exponential doesn't quite fit or when you have steep drop-offs.

THE GAUSSIAN CORRELATION INEQUALITY

... in convex geometry:

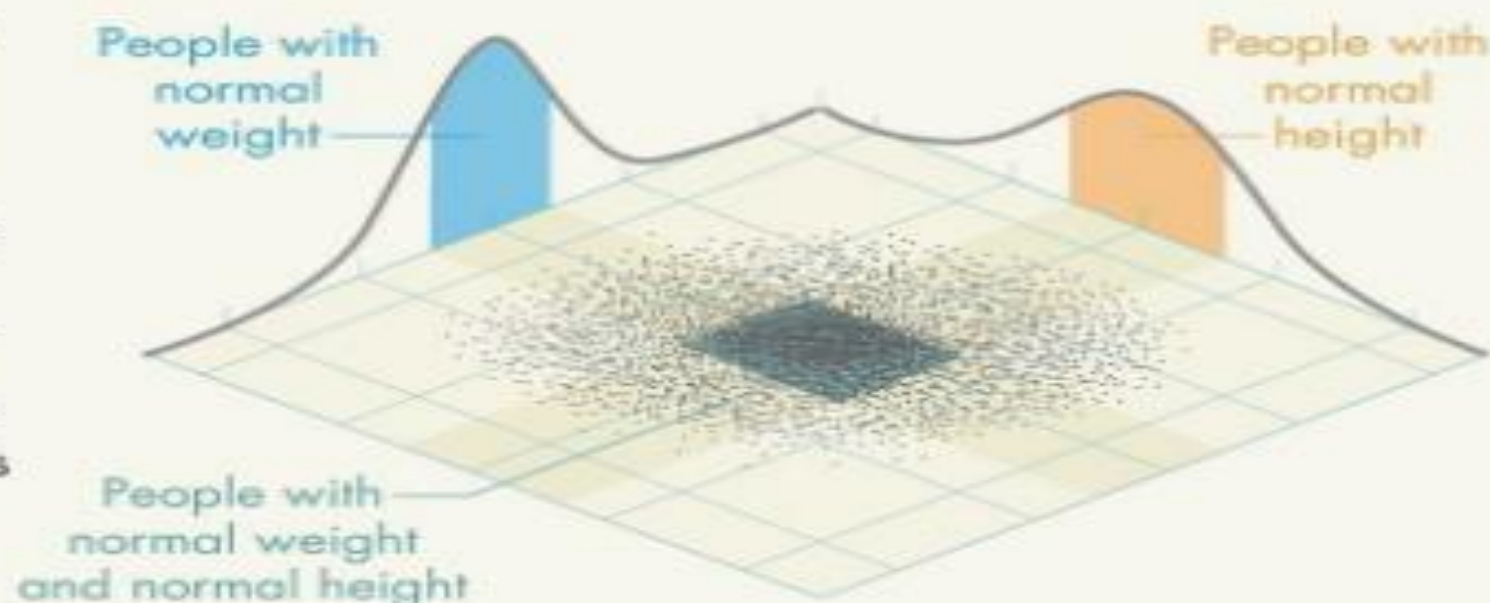
Consider any two convex symmetrical shapes in any number of dimensions that are centered on the same point, which forms a target. Darts thrown at the target will land in a bell curve or "Gaussian distribution" of positions. The overlap of the two shapes increases your probability (P) of striking both.



$$P(\text{Dart lands in both circle + rectangle}) \geq P(\text{Dart lands in circle}) \times P(\text{Dart lands in rectangle})$$

... in multivariate statistics:

Consider a sample of people's weights and heights plotted on x-y axes. Because these variables are correlated, the odds that someone's weight and height will both fall within a combined range is greater than or equal to the product of the independent odds of falling in each range. (The general inequality holds for any number of variables.)



$$P(\text{Person has both normal weight + height}) \geq P(\text{Person has normal weight}) \times P(\text{Person has normal height})$$

Overview

Process

Subtypes

Examples

Supervised Learning

Majority of algorithms. Machine is trained using **well-labeled data**; inputs and outputs are matched.

Mapping function takes inputs and matches to outputs, creating a target function.

Classification,
Regression

Linear regression,
Random forest,
SVM.

Unsupervised Learning

Unlabeled data (inputs only) is analyzed. Learning happens without supervision.

Inputs are used to create a model of the data.

Clustering,
Association.

PCA,
k-Means,
Hierarchical clustering.

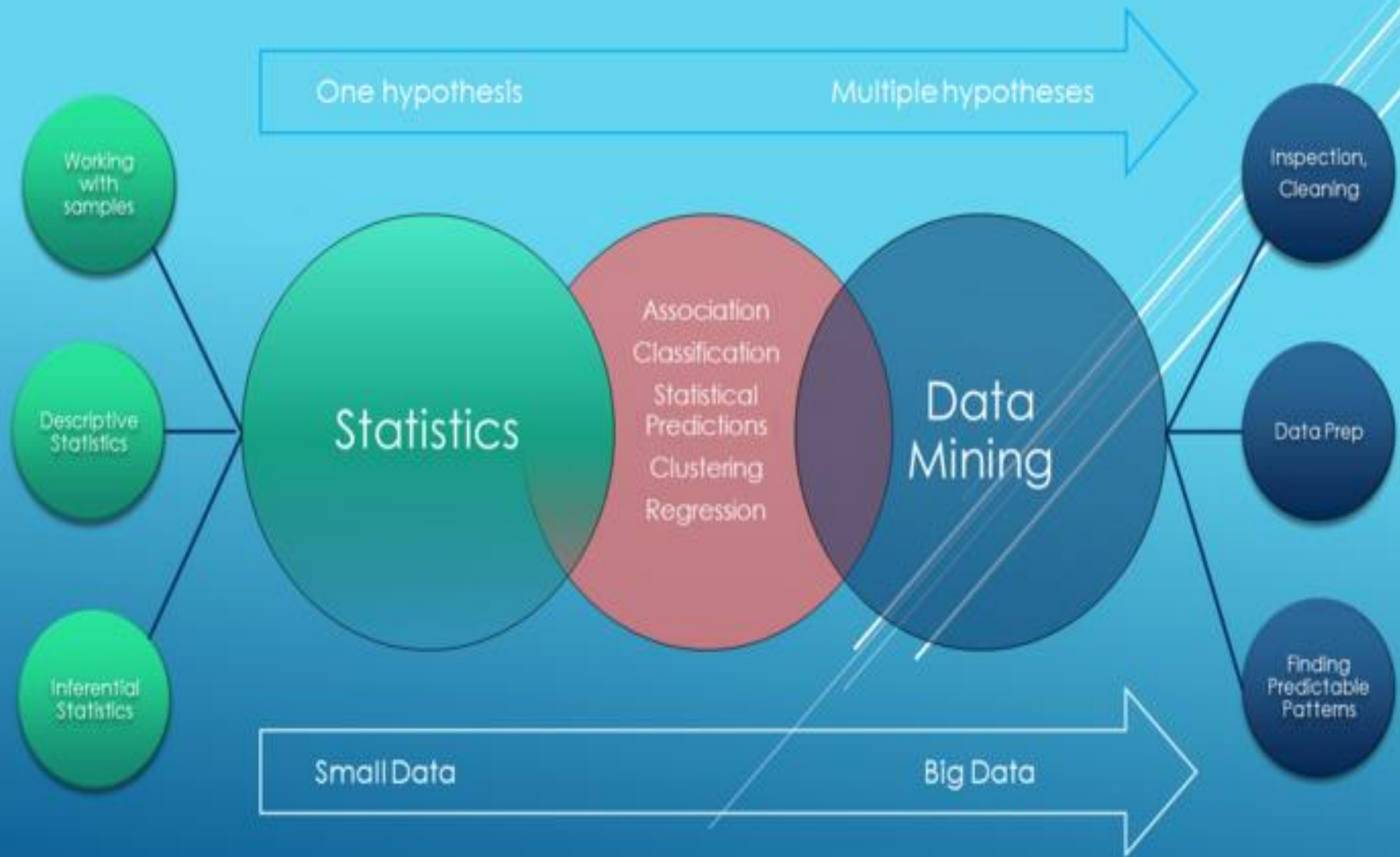
Semi supervised

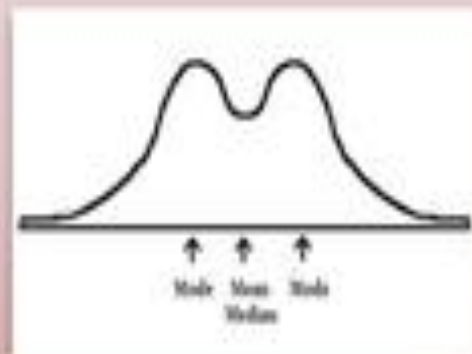
Some data is labeled, some not. Goal: better results than labeled data alone. Good for real world data.

Combination of above processes.

All the above.

Self training,
Mixture models,
Semi-supervised SVM





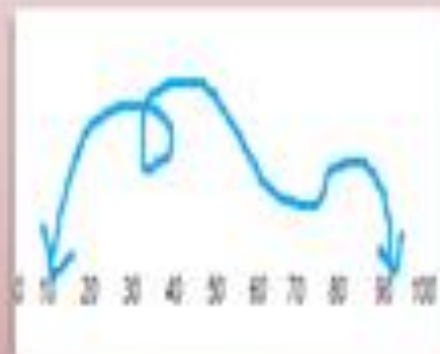
Central Tendency

Mean, Median, Mode, Outliers



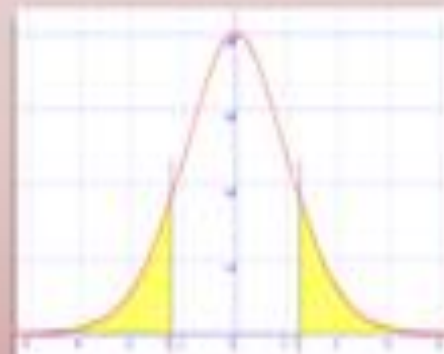
Measures of Spread

Range, Standard deviation, Variance, Quartiles



Percentiles

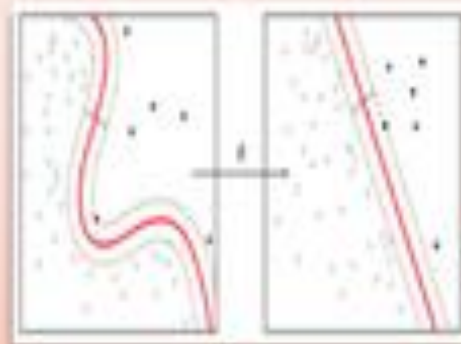
Position of data, percentile rank, percentile range



Probability Distributions:

Uniform, normal (Gaussian), Poisson

Basic Probability and Statistics



Dimensionality reduction

Pruning, PCA



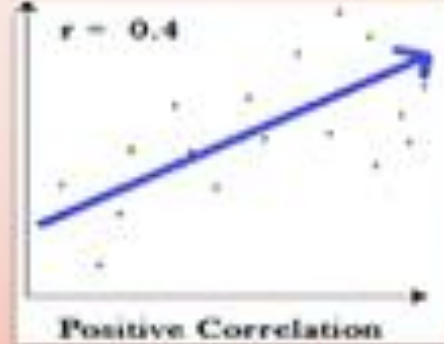
Sampling

SRS, Reservoir, Undersampling, Oversampling,

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

Bayesian statistics

Measuring belief or confidence



Covariance & correlation

How data is related

More Advanced Probability and Statistics

State the hypotheses

- The null hypothesis: An all you can eat pizza chain thinks customers eat an average of 4 slices.
- Alternate: You think the average is much higher.

Collect your data

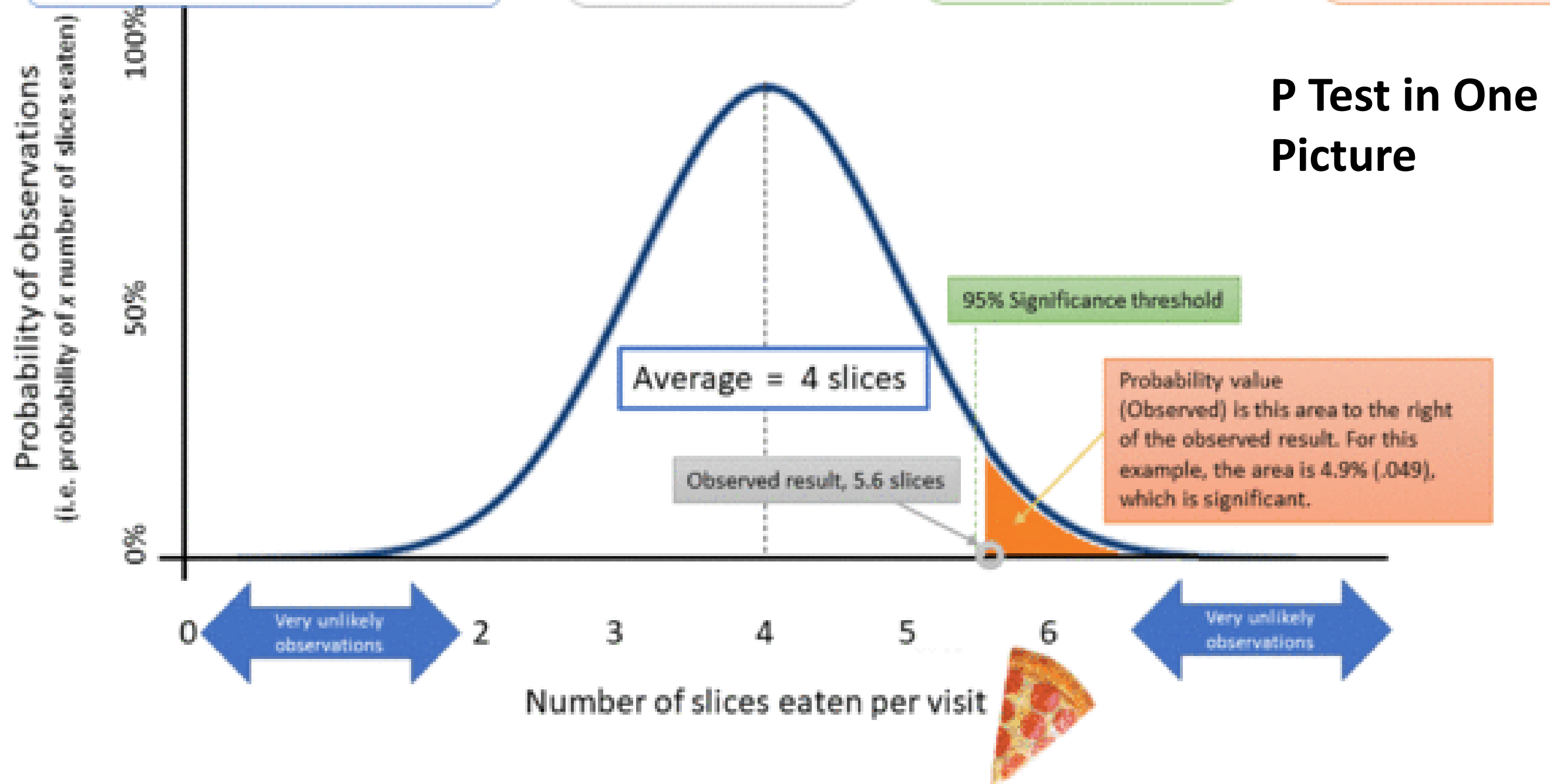
- Collect your observations, making sure they are random. For this example, the observed average is 5.6 slices.

Test the result (avg = 5.6)

- Set your significance level (5% in this example), then run your test. For example: chi-square, T-test or Z-test.

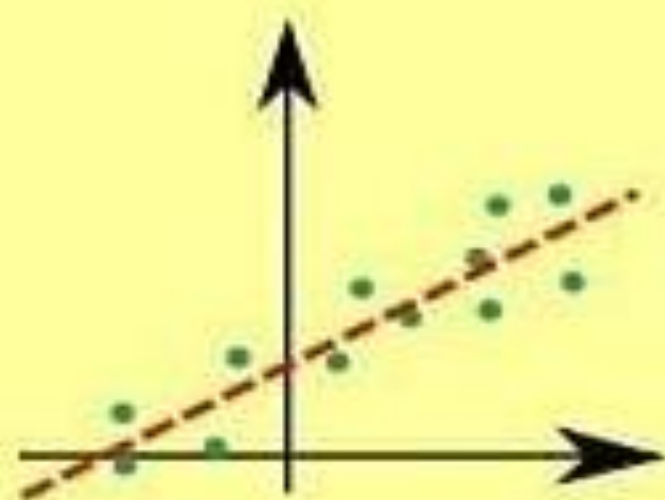
Is the result significant?

- $P > .10$: Not significant
- $p \leq .10$: Marginally significant
- $p \leq .05$: Significant
- If $p \leq .01$ Very significant.



LINEAR REGRESSION

- 1 Econometric modelling
- 2 Marketing Mix Model
- 3 Customer Lifetime Value



Continuous \Rightarrow Continuous

$$y = \alpha_0 + \sum_{i=1}^N \alpha_i x_i$$

`lm(y ~ x1 + x2, data)`

1 unit increase in x
increases y by α

LOGISTIC REGRESSION

- 1 Customer Choice Model
- 2 Click-through Rate
- 3 Conversion Rate
- 4 Credit Scoring



Continuous \Rightarrow True/False

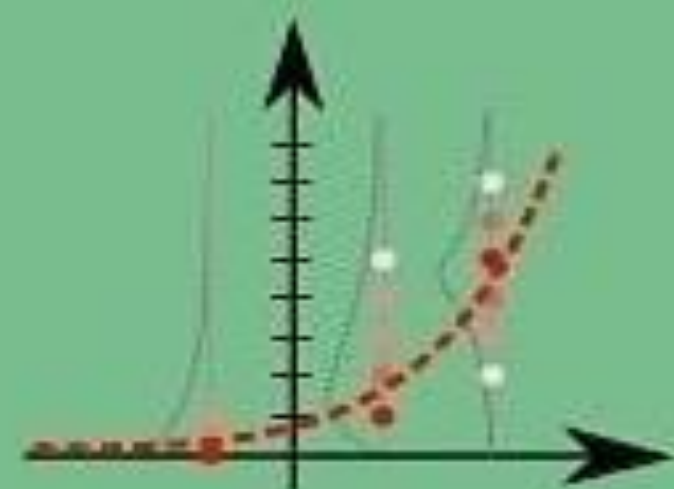
$$y = \frac{1}{1 + e^{-z}}$$
$$z = \alpha_0 + \sum_{i=1}^N \alpha_i x_i$$

`glm(y ~ x1 + x2, data,
family=binomial())`

1 unit increase in x
increases log odds by α

POISSON REGRESSION

- 1 Number of orders in lifetime
- 2 Number of visits per user



Continuous \Rightarrow 0,1,2,...

$$y \sim \text{Poisson}(\lambda)$$
$$\ln \lambda = \alpha_0 + \sum_{i=1}^N \alpha_i x_i$$

`glm(y ~ x1 + x2, data,
family=poisson())`

1 unit increase in x
multiplies y by e^{α}

Differences Between Pearson Correlation and Linear Regression

Use when...

Results

Assumptions (Requirements)

$r = 0.4$



Correlation

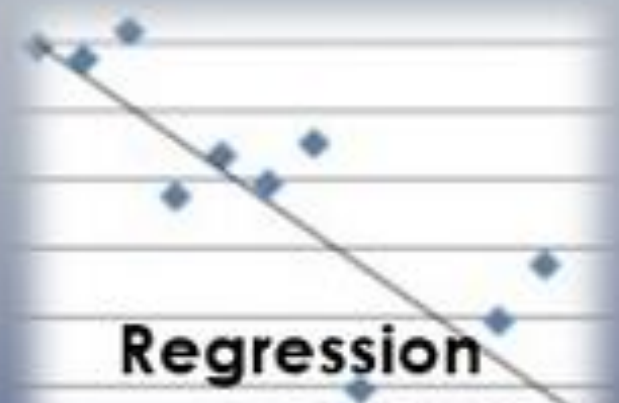
You want to know if there's **an association** between two variables.

Strength & direction of relationship (r):

-1 = perfect negative,
+1 = perfect positive,
0 = no correlation.

Validity: Valid measurements, a good sample, unconfounded comparisons.

Distribution: Linear relationship between the two.



Regression

You want to predict **how one variable will change with another.**

Estimates of parameters for a regression equation:

The B_0 and B_1 in the linear regression equation
 $Y = B_0 + B_1 * X$

All the above, plus:

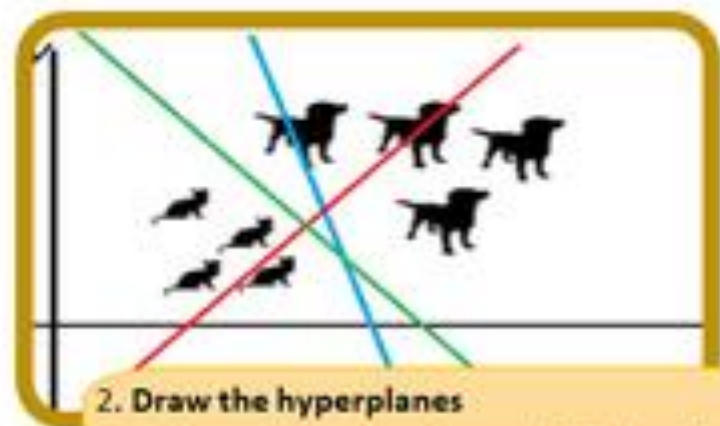
1. Quantitative Data
2. Outlier Condition
3. Independence of Errors
4. Homoscedasticity
5. Normality of Error Distribution

SVM in One Picture



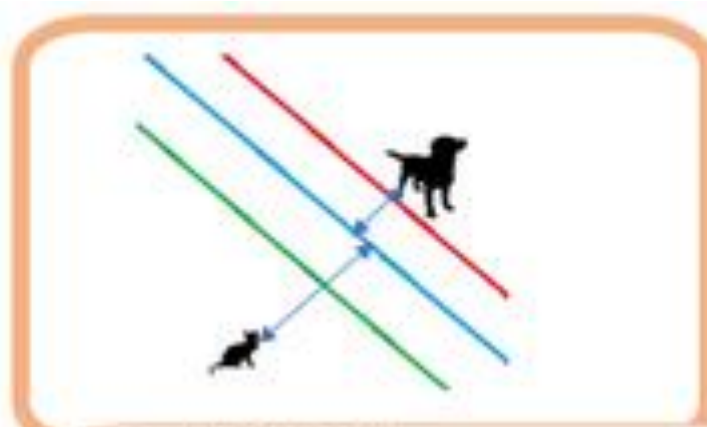
1. Plot the features

- Plot in n -dimensional space (n = number of features). For this example, we have 2 features—cats and dogs—so $n = 2$ (2 dimensional)
- Value of each feature (dogs and cats) = value of coordinate



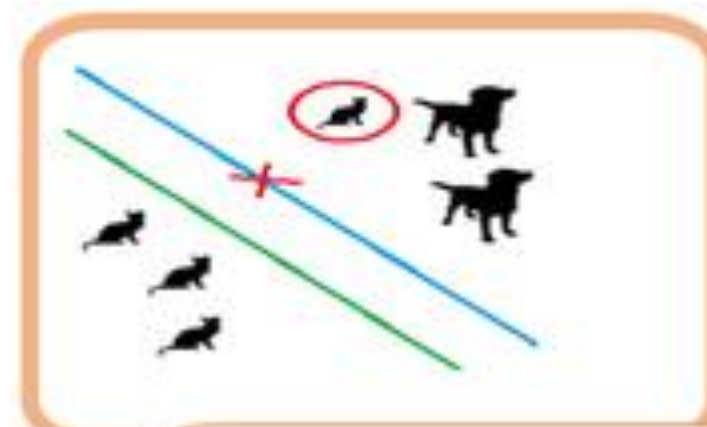
2. Draw the hyperplanes

- The hyperplanes are the possible “best” boundaries to separate features (in this case, cats and dogs).
- Consider all hyperplanes. Here, I plotted 3 for simplicity. In reality, you could have an infinite number of choices!



3a. Consider the Margins.

- A **margin** is the distance between any point and the hyperplane. The points touching this boundary are the Support Vectors (in this simple example, the dog and cat shown are the support vectors).
- Here, the blue line **maximizes the distance**, so it is the “best”.



3b. Look out for Misclassification

- SVM will consider **misclassification** first.
- If the “best” (blue) results in misclassification, the hyperplane is discarded in favor of one that classifies correctly, but with a less-optimal margin (green).



3. Choose the best hyperplane

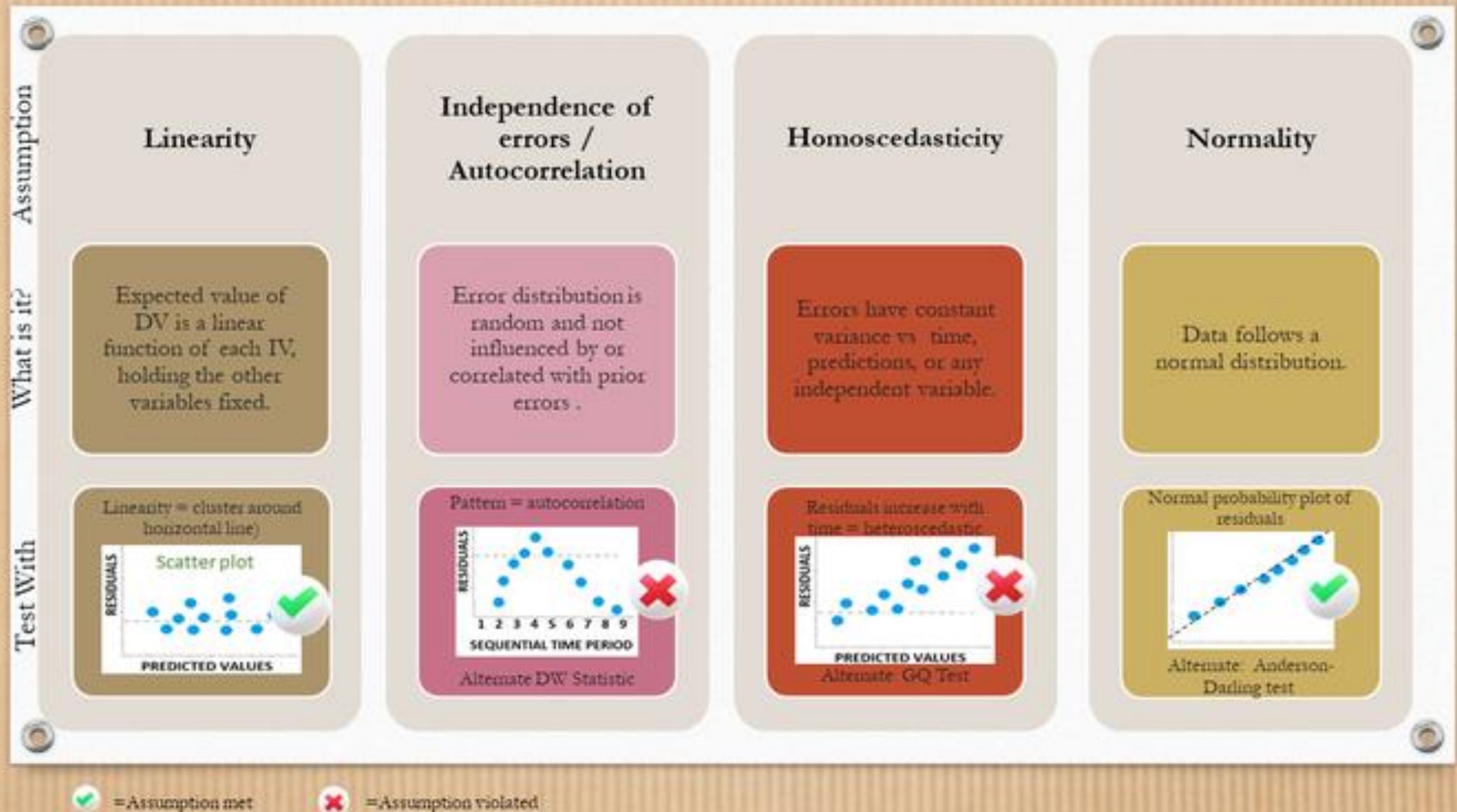
- In this example, the best hyperplane is the green one, which **clearly separates** cats and dogs.
- Take into consideration *margins*, *misclassification*, and *curved/linear* possibilities when making your choice.



3c. Curved or Linear?

- The best choice might not be a straight line; it may be a **curved one**.
- SVM considers all possible hyperplanes (straight and curved) before choosing.

Assumption of Logistic Regression in One Picture



Logistic Regression in One Picture

1. Collect Observations

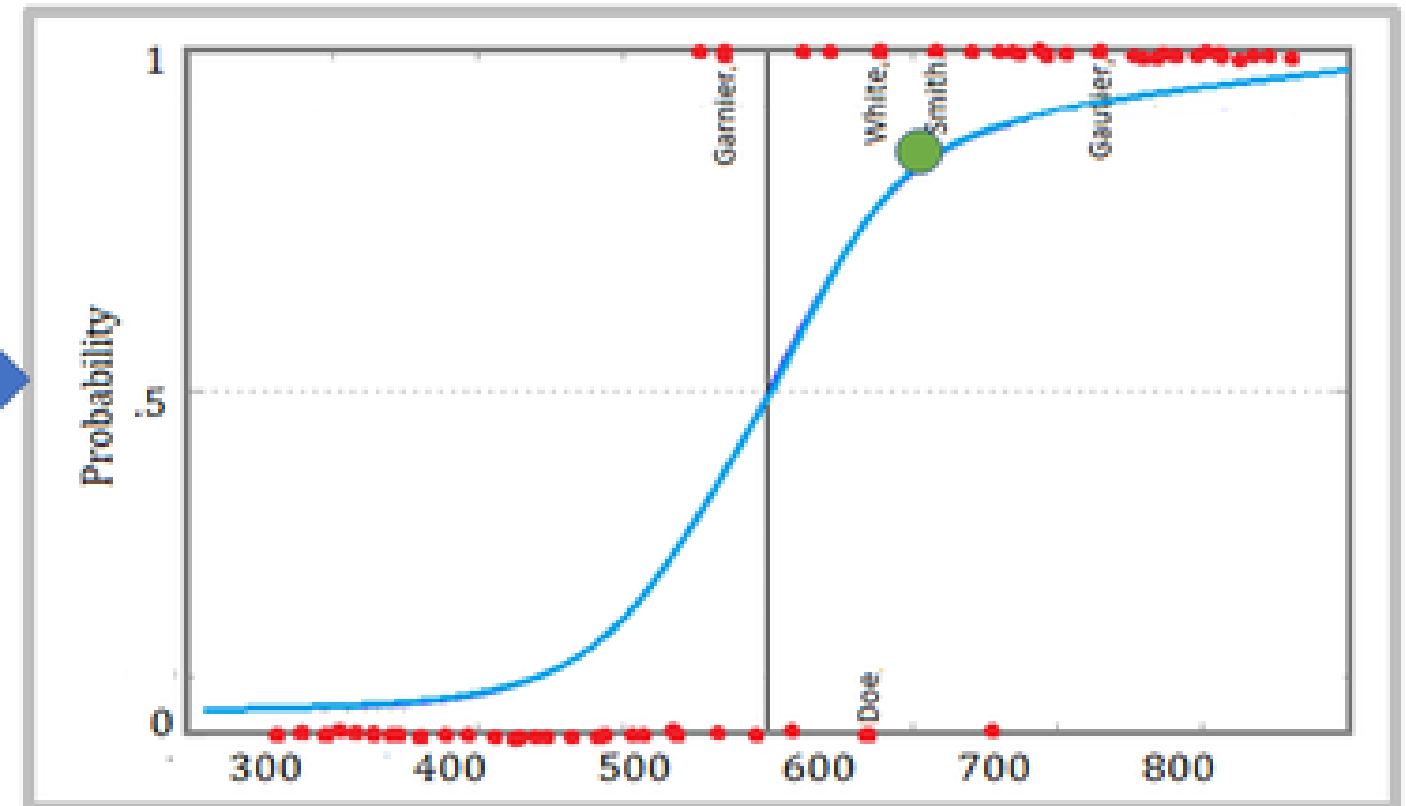
Observations should be classified by success (1) or failure (0). For example, data on 1,000 people, their credit scores, and whether they got approved (1) or denied (0) for a car loan.

N = 1,000

Smith, J., 660, approved (1)
Doe, J., 630, denied (0)
Garnier, B., 550, approved, (1)
Gautier, P., 750, approved, (1)
White, Z., 640, approved (1) ...

2. Build Model

Data points (shown in red in this example) are plotted on a traditional scatter plot. They are binary, so will be clustered on either the "0" line at the bottom, or the "1" line at the top. A regression equation is calculated, which typically follows an S-shape (shown in blue).



3. Read model's output

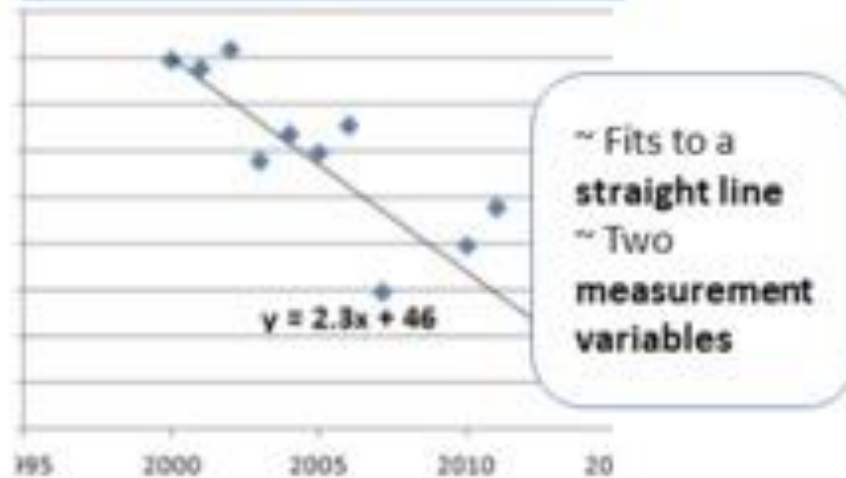
For example, "If I feed a credit score of 660 into the model, what is the probability the person will get approved for a car loan?"

Score of 660 (green dot on graph) = .8

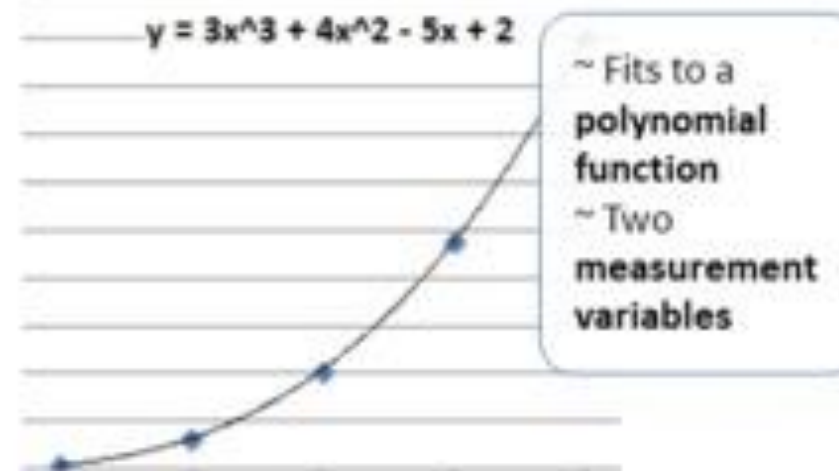
Regression Analysis in One Picture

Basic Fitting

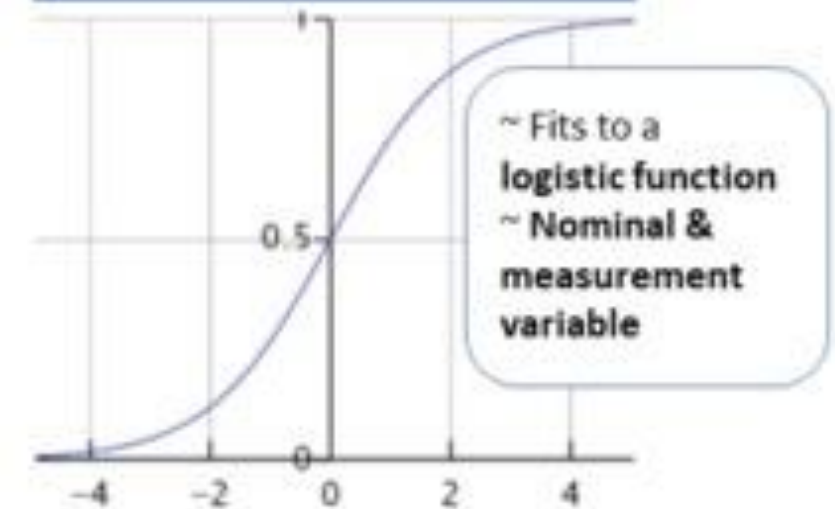
Linear Regression



Polynomial Regression

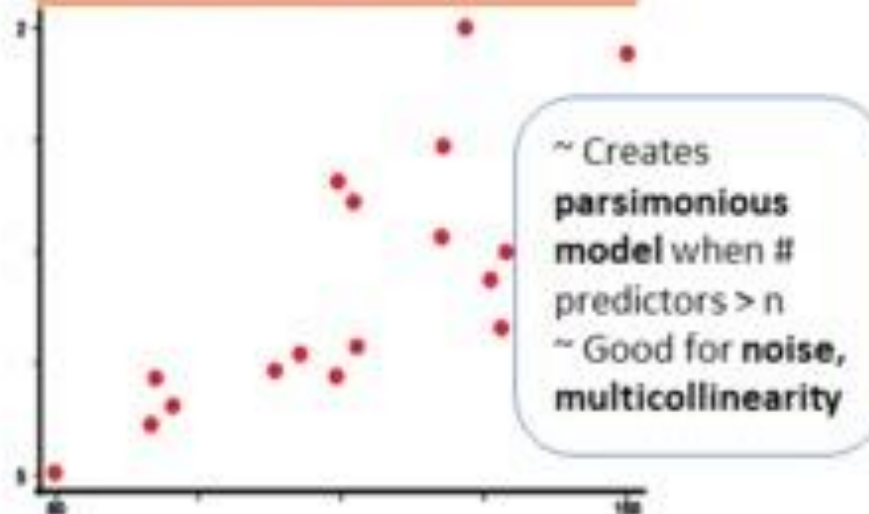


Logistic Regression

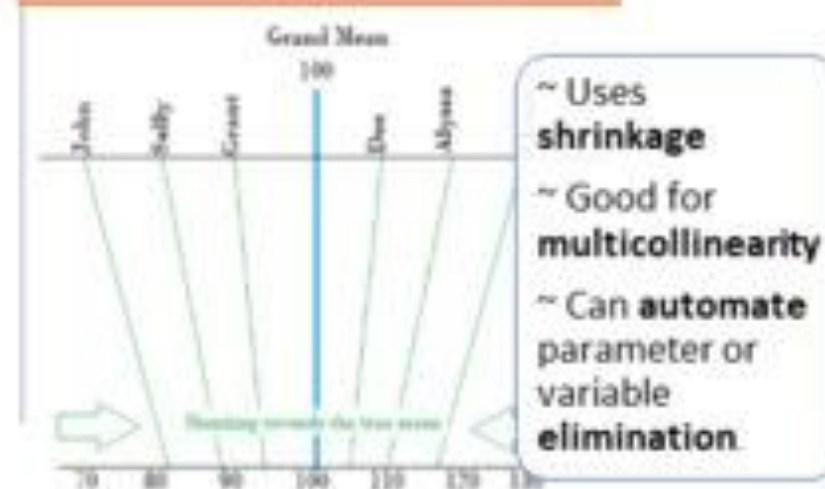


Special Situations

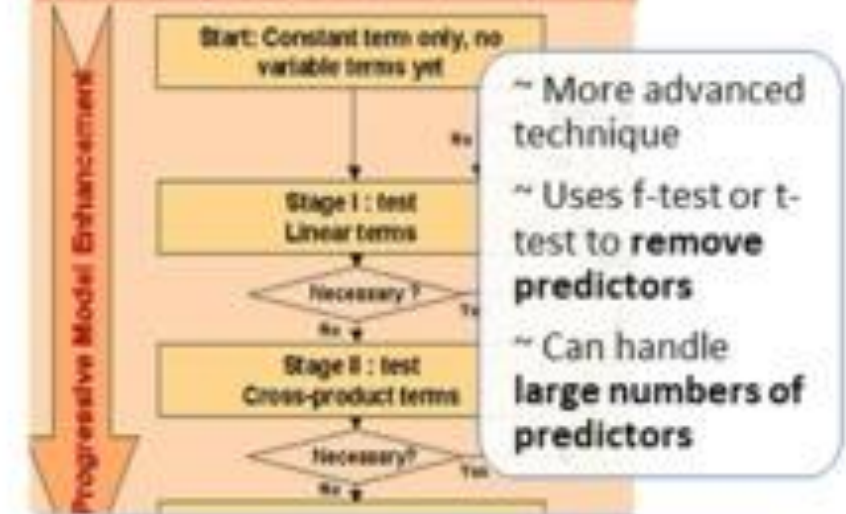
Ridge Regression



Lasso Regression

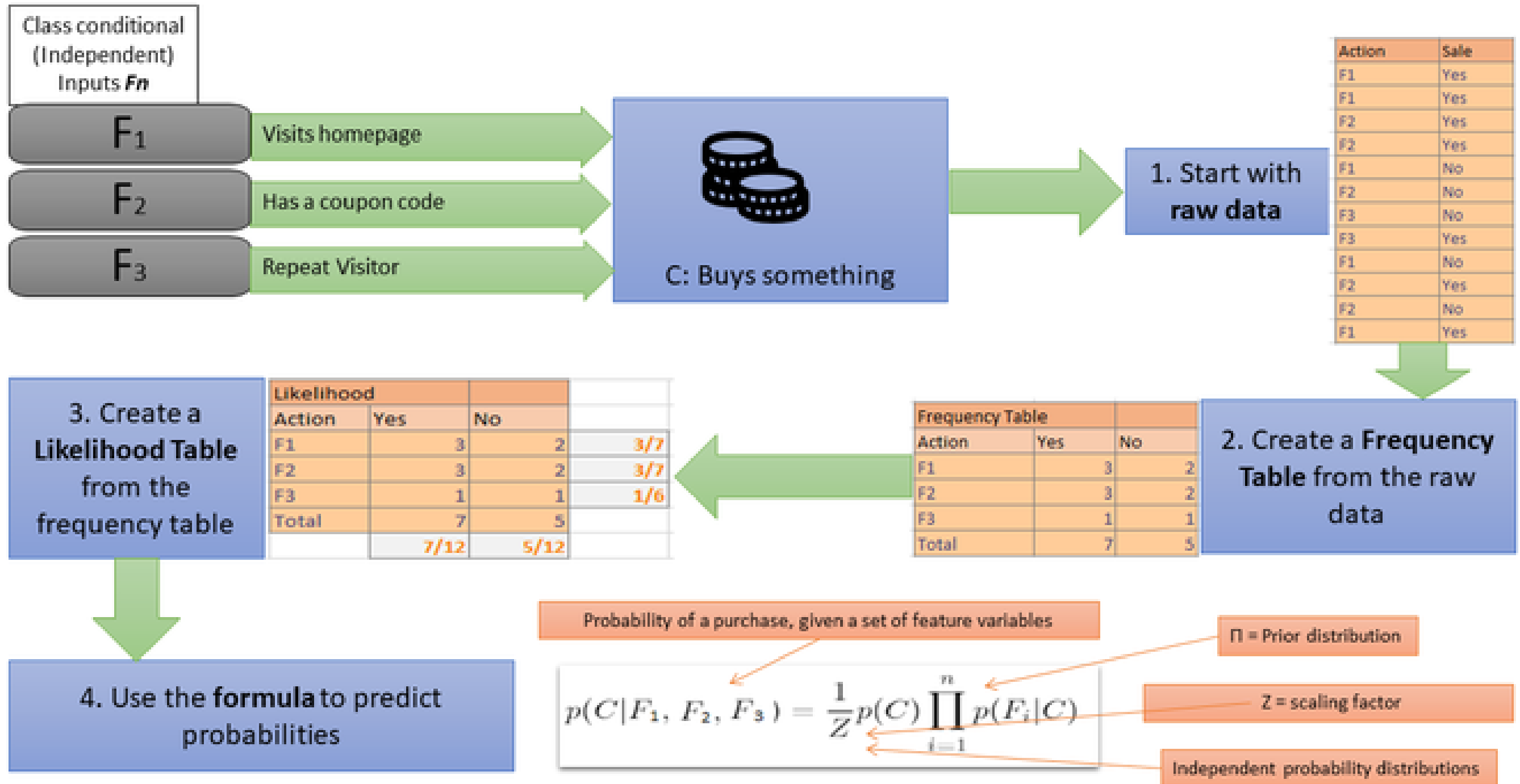


Stepwise Regression

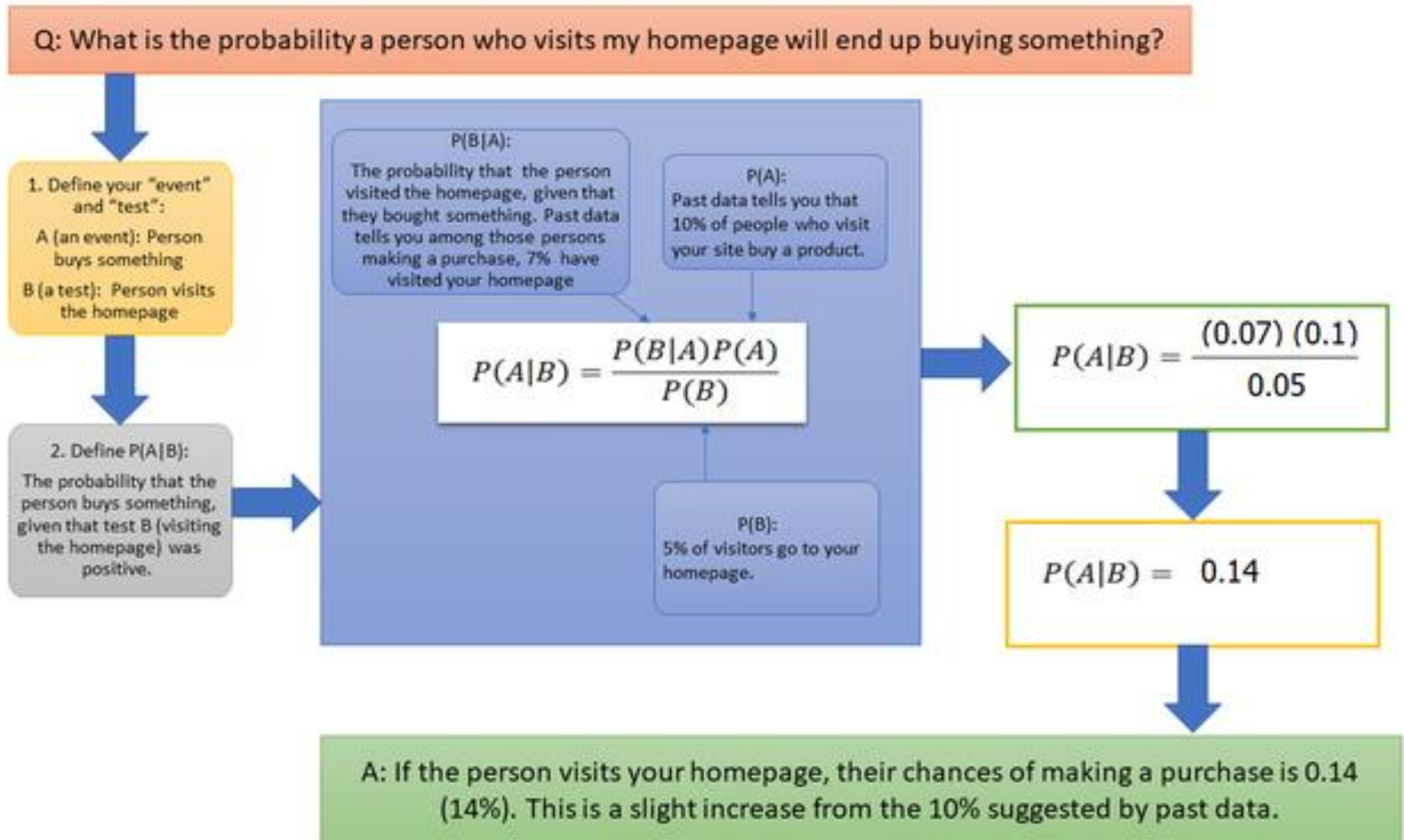


Naïve Bayes in One Picture

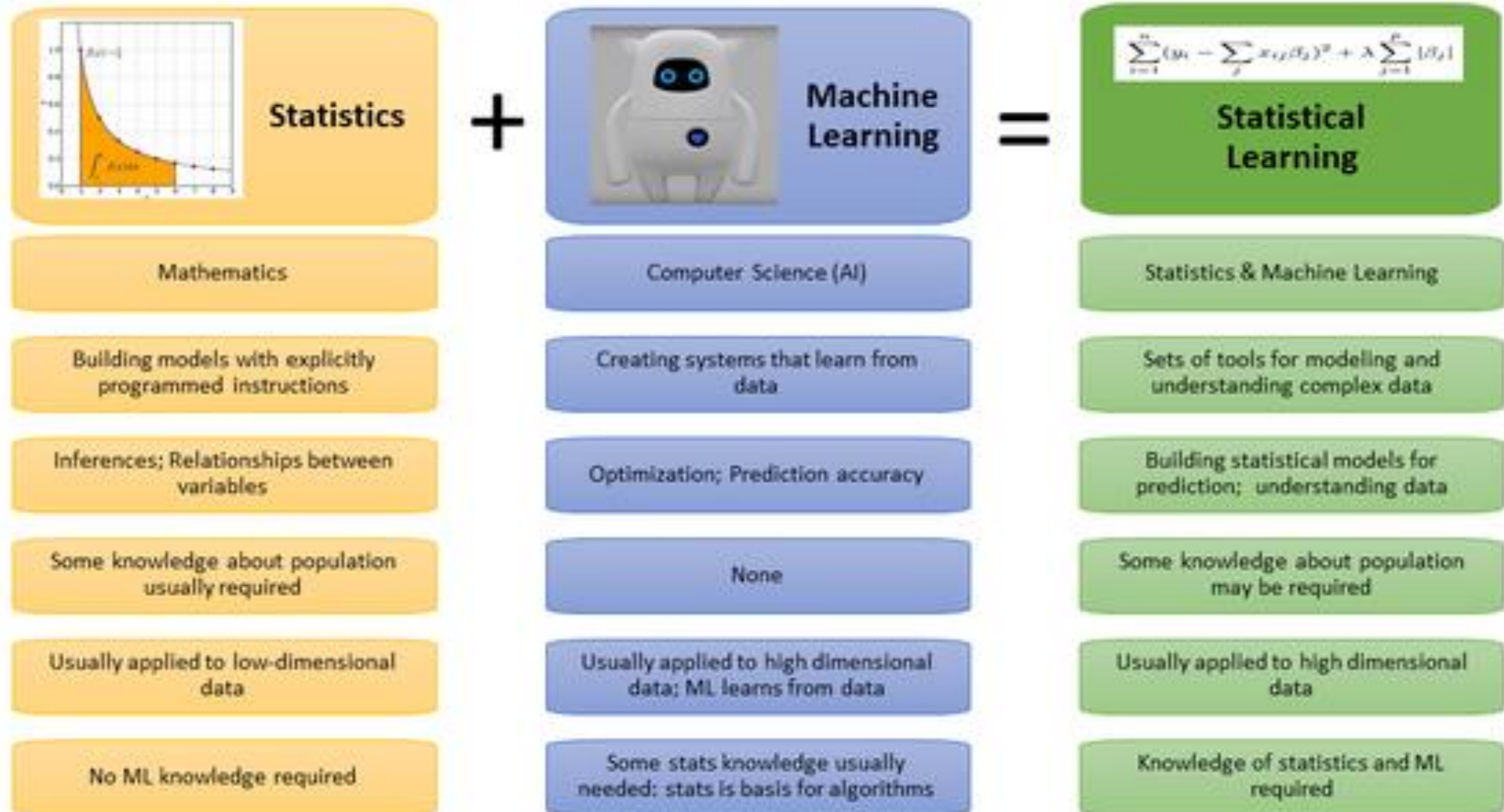
Q. If a person visits the homepage, has a coupon code, and is a repeat visitor, what is the probability they will buy something?



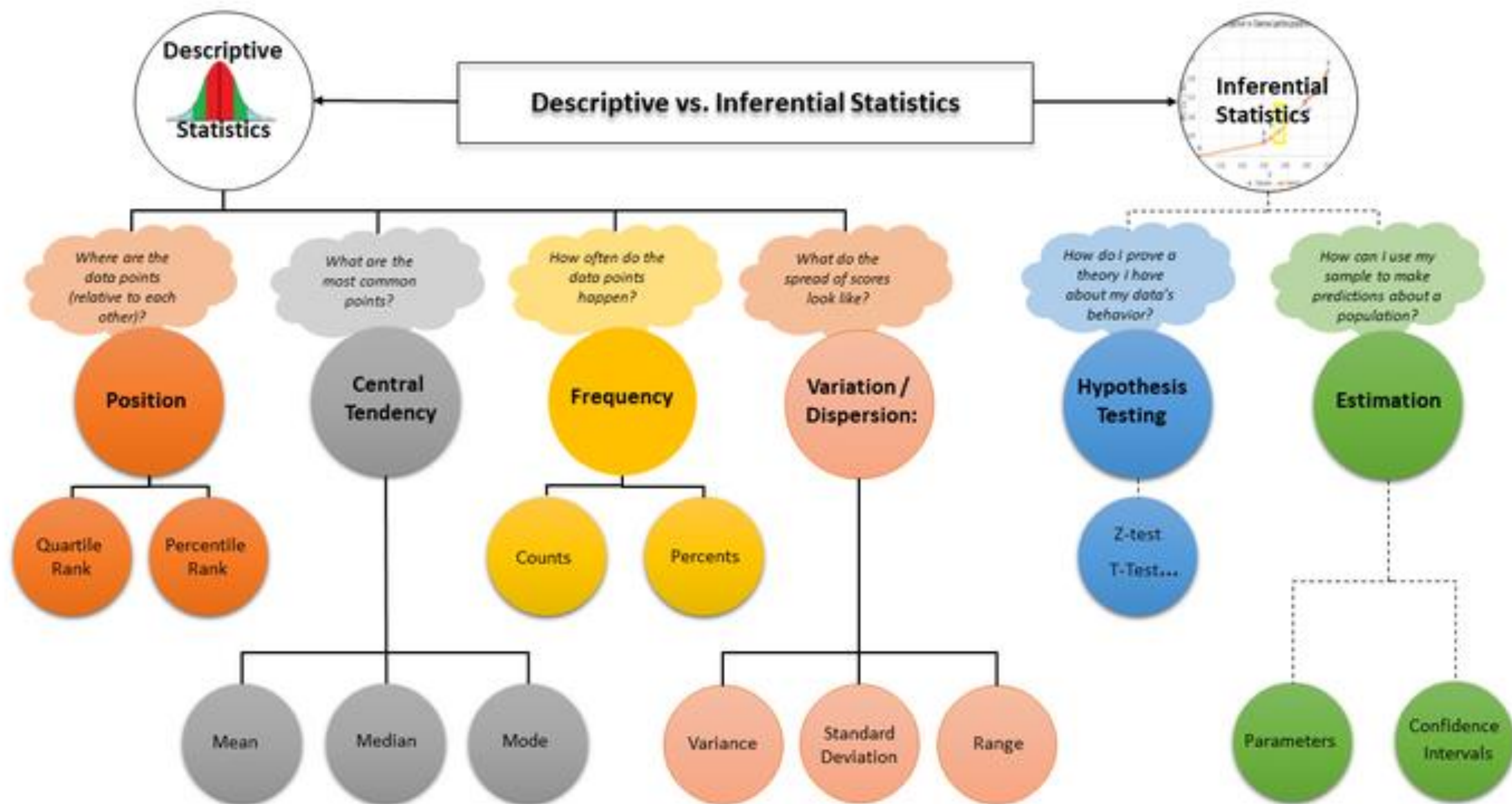
Bayes Theorem in One Picture



Statistics and Machine learning In One Picture



Musio image: Akawikipic [CC BY-SA 4.0 (<https://creativecommons.org/licenses/by-sa/4.0/>)]

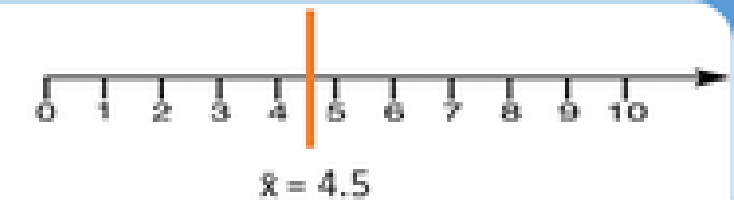


Standard Deviation vs. Standard Error

Standard Deviation with One Set of Measurements

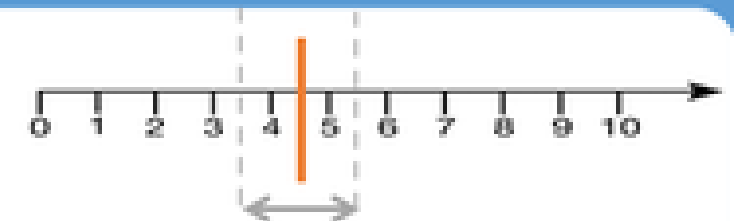
Step 1: Find the Mean
(What is the Average Value?)

Mean (\bar{x}) of **Sample 1**
{4.5, 4.5, 3, 6, 4, 5}
= 4.5



Step 2: Find the Standard Deviation
(What does the spread of data look like?)

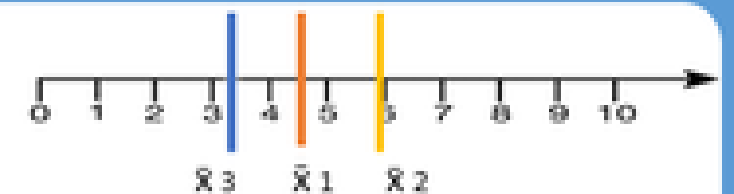
Standard Deviation of **Sample 1**
= 1
(This tells you the bulk of the scores are one standard deviation either side of the mean)



Standard Error with Multiple Sets of Measurements

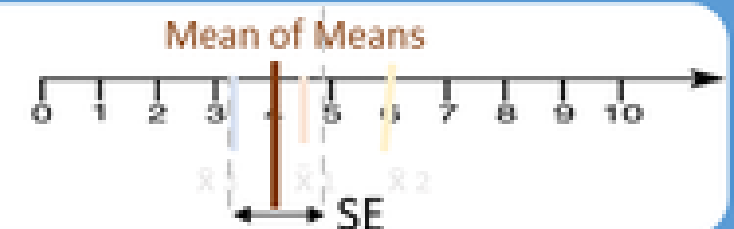
Step 1:
Find the Means for each separate sample

Sample 1: Mean = 4.5
Sample 2: Mean = 6
Sample 3: Mean = 3.23



Step 2: Find the Standard Error
(aka the Standard Deviation for the average mean)

First, find the "Mean of the Means":
MoM = $(4.5 + 6 + 3.23)/3 = 4.577$
Then find the standard error = 0.8



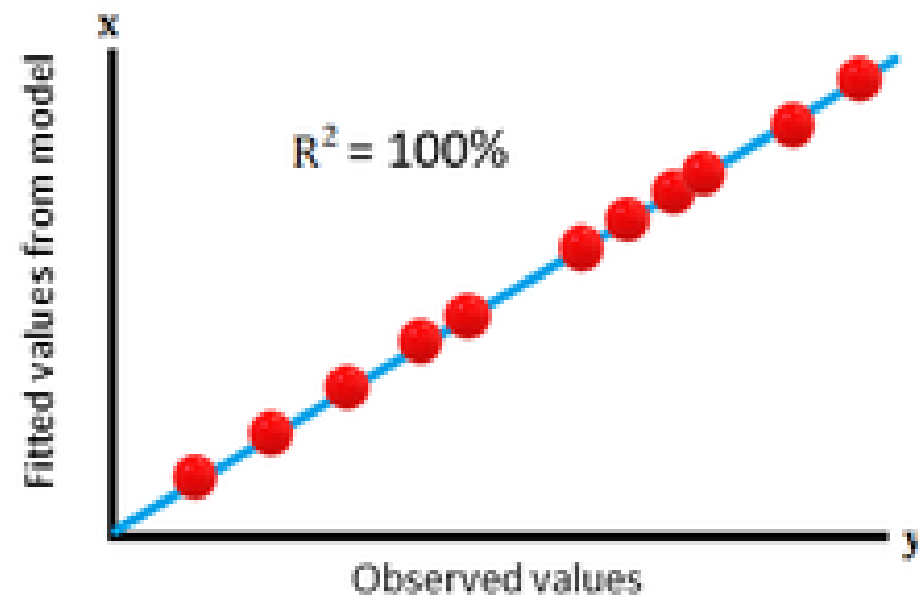
Correlation Coefficients in One Picture

	Pearson's	Spearman's Rank	Kendall's Tau	Gamma	Yule's Q
Relationships evaluated	Linear	Monotonic	Monotonic	Monotonic	Monotonic
Parametric or Non Parametric	Parametric	Non Parametric	Non Parametric	Non Parametric	Non Parametric
Input variables	Continuous	Continuous, ordinal, or ratio	Ordinal	Ordinal	Ordinal (2x2)
Advantages	Easy to use and understand Widely available	Normal distribution not required Widely available	Insensitive to errors More accurate p-value for smaller samples than Spearman's	Good for tied ranks Not affected much by outliers	Simplified version of Gamma for 2x2 table

R-Square metrics in One Picture

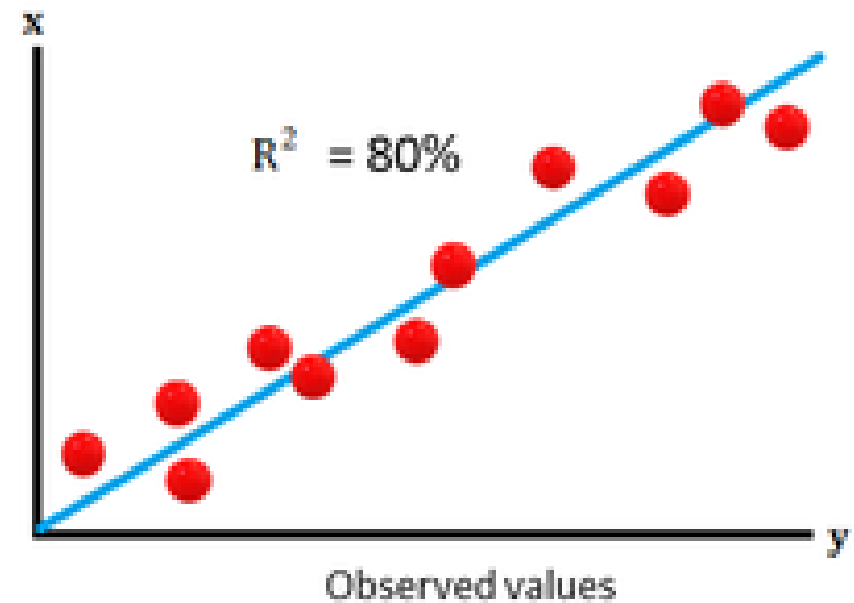
Comparison of R-Squared for Different Linear Models (Same Data Set)

Great



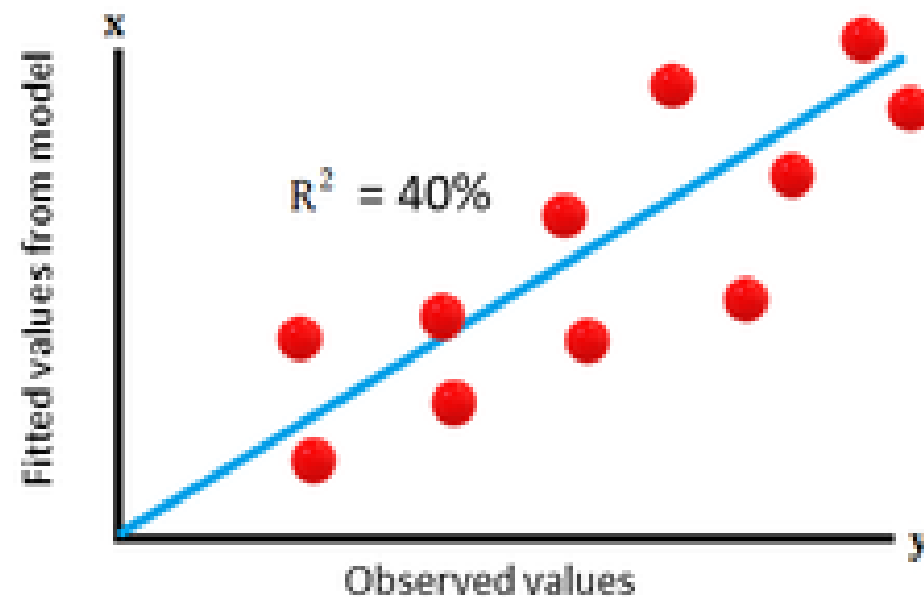
Fitted = observed: Model explains all variance

Good



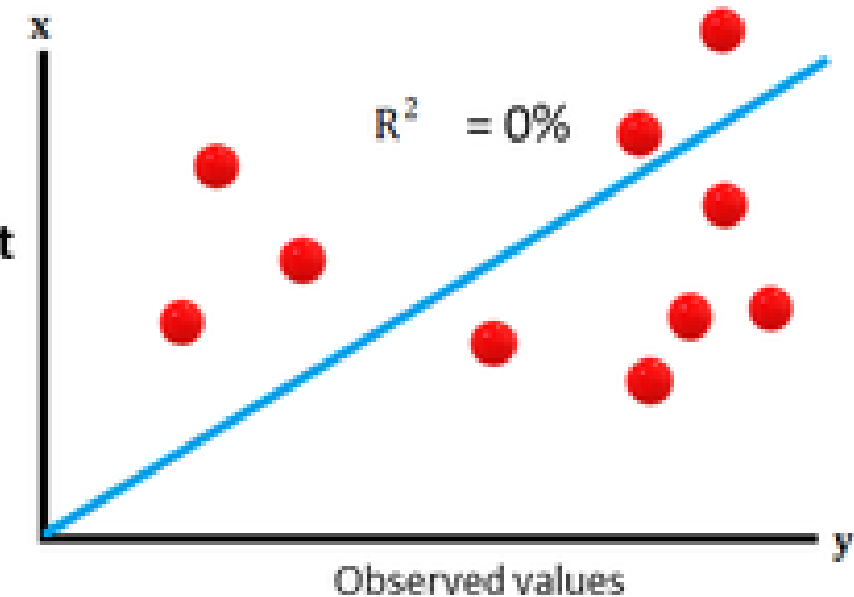
Model explains bulk of variance

OK



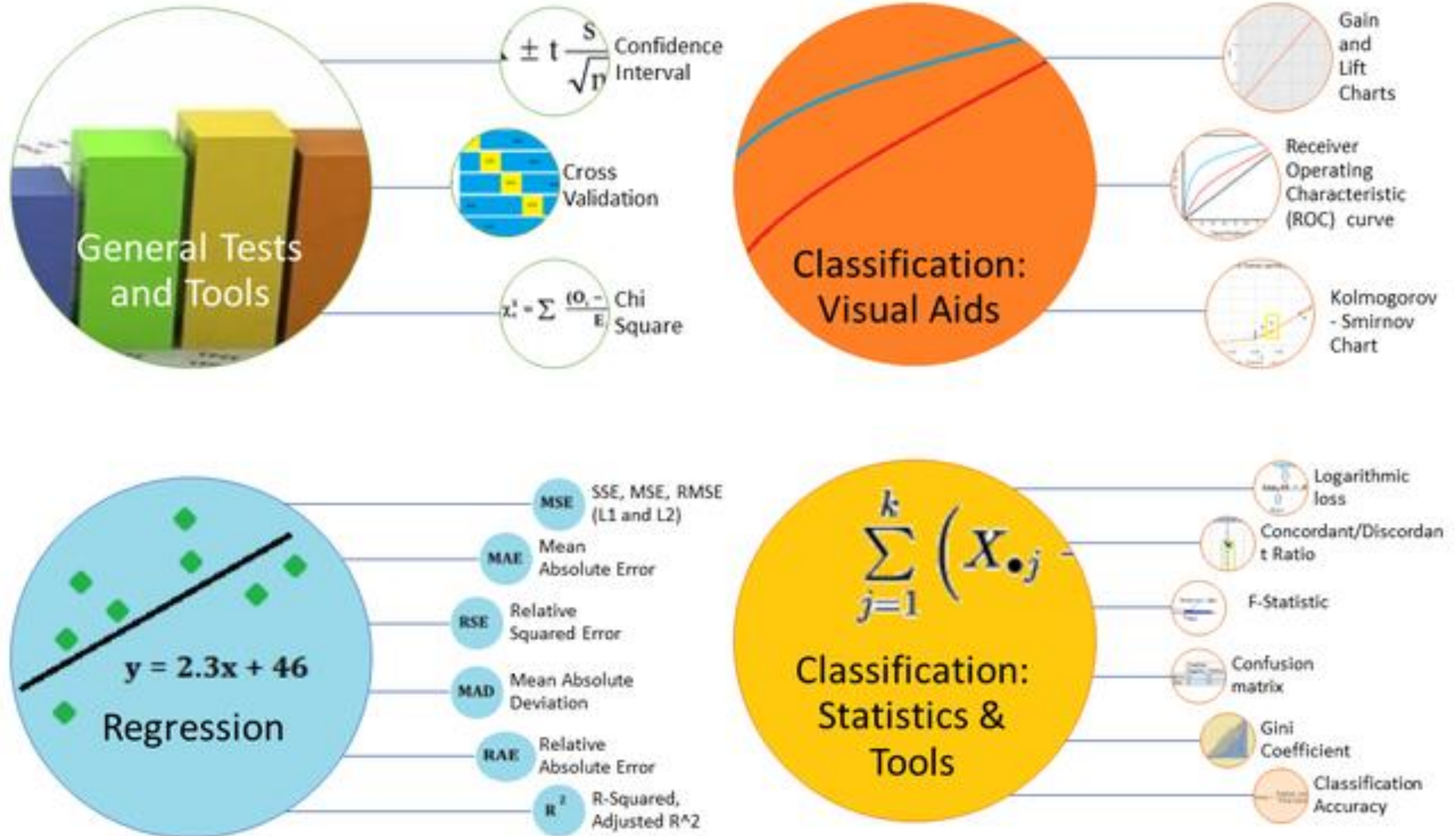
Model explains 40% of variance, so is reasonable.

Inconsistent

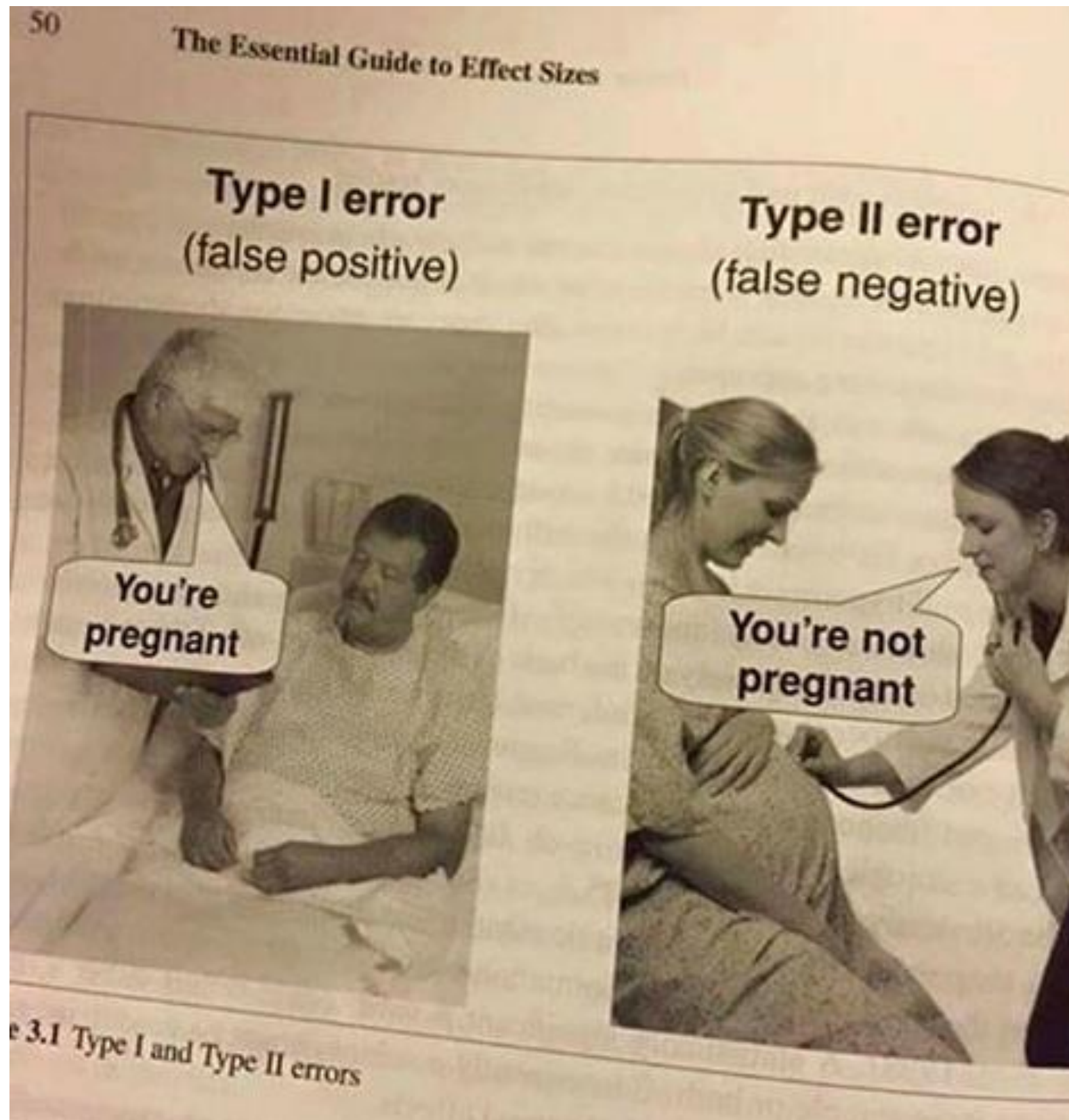


Model fails to explain any variance

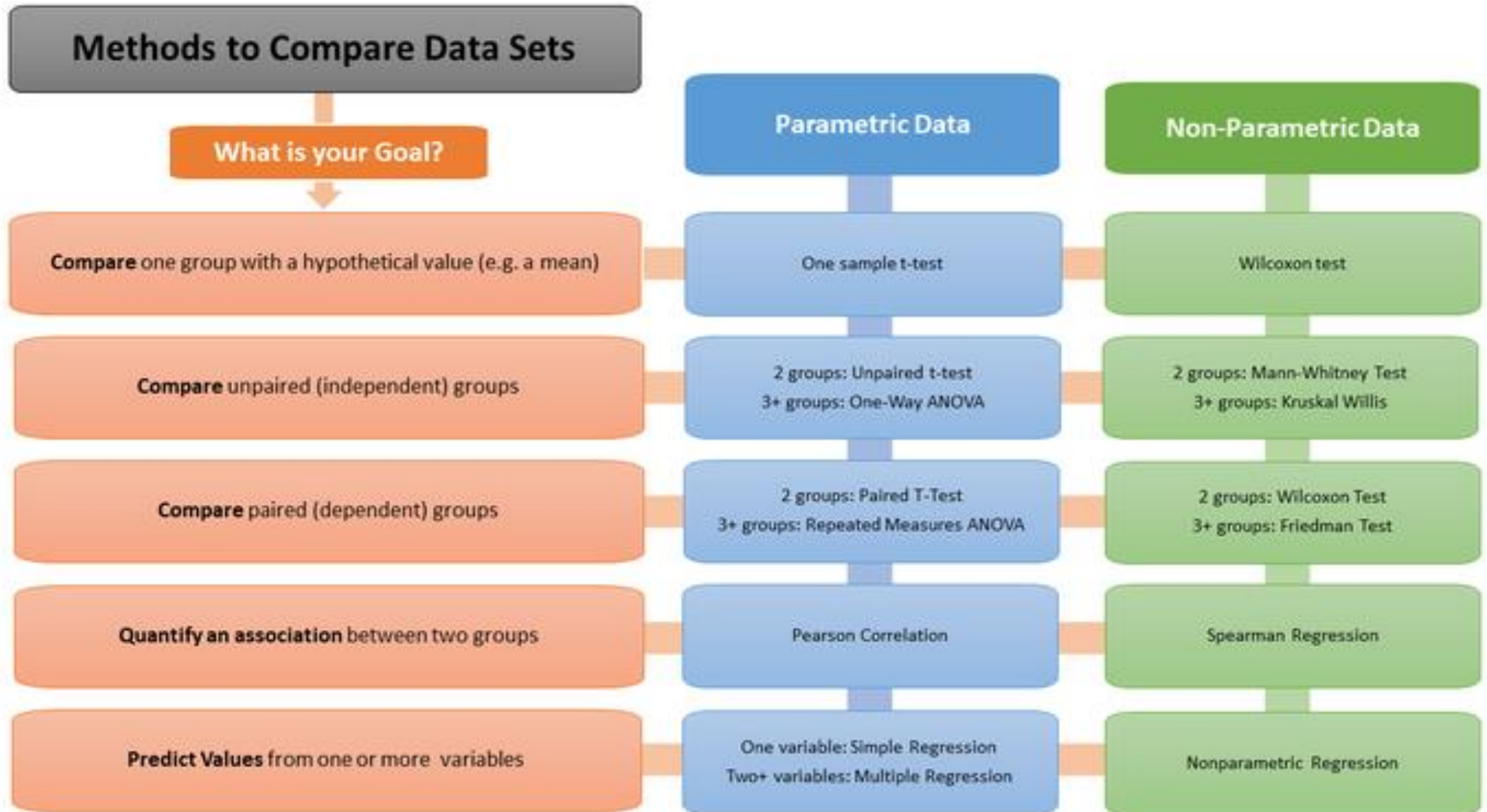
Evaluation metrics in One Picture



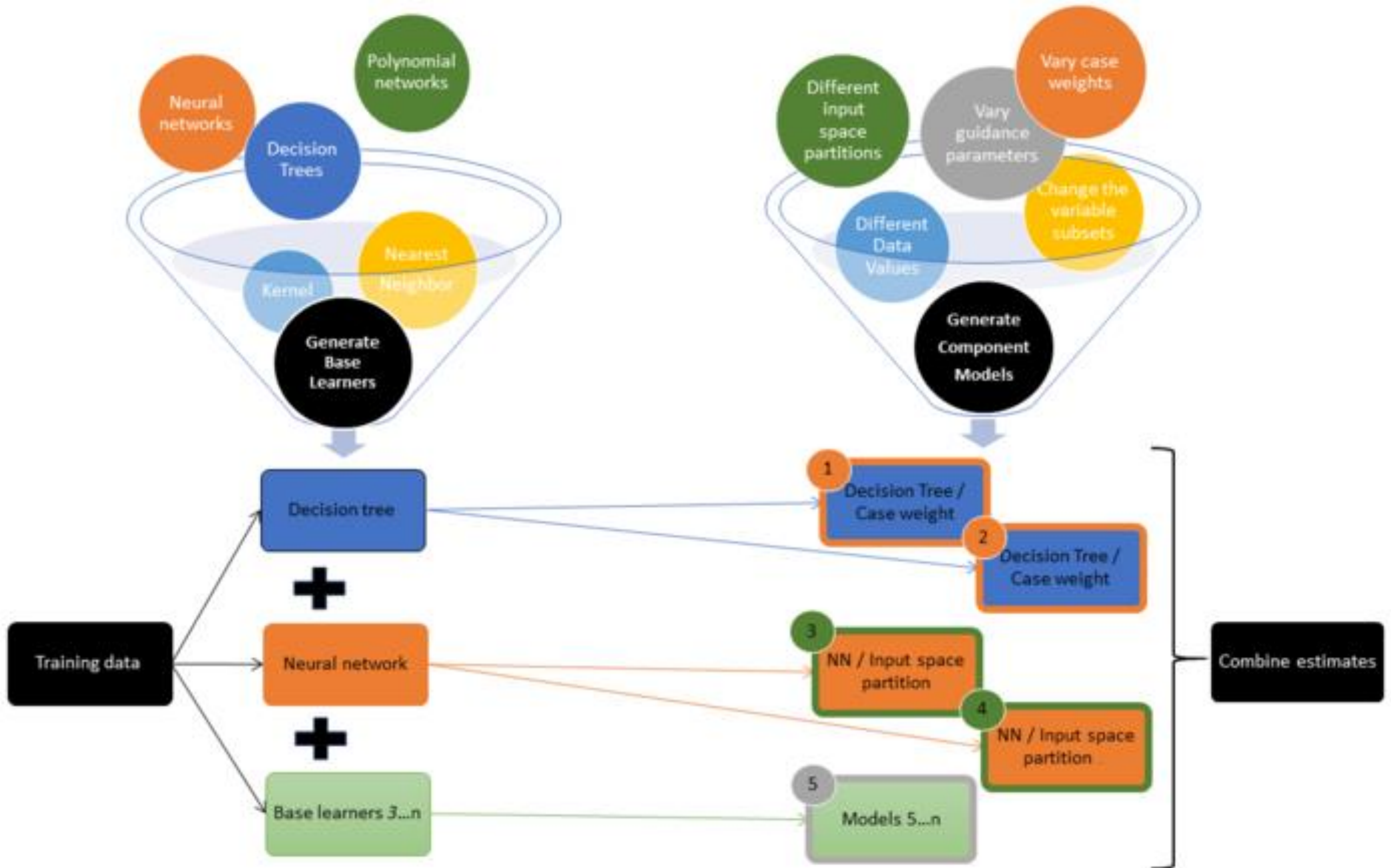
Type-I and Type-II in One Picture



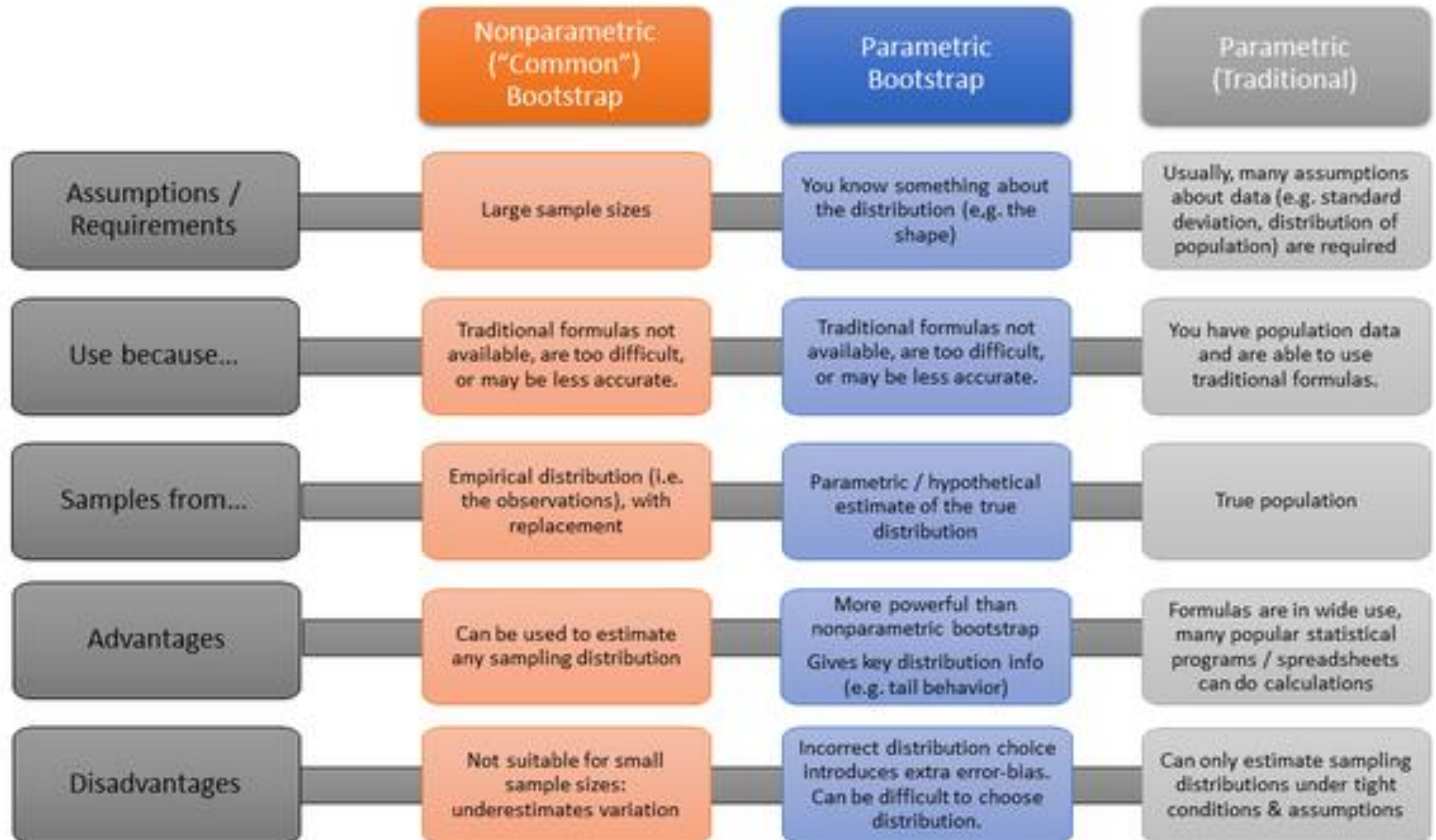
Comparing Dataset in One Picture



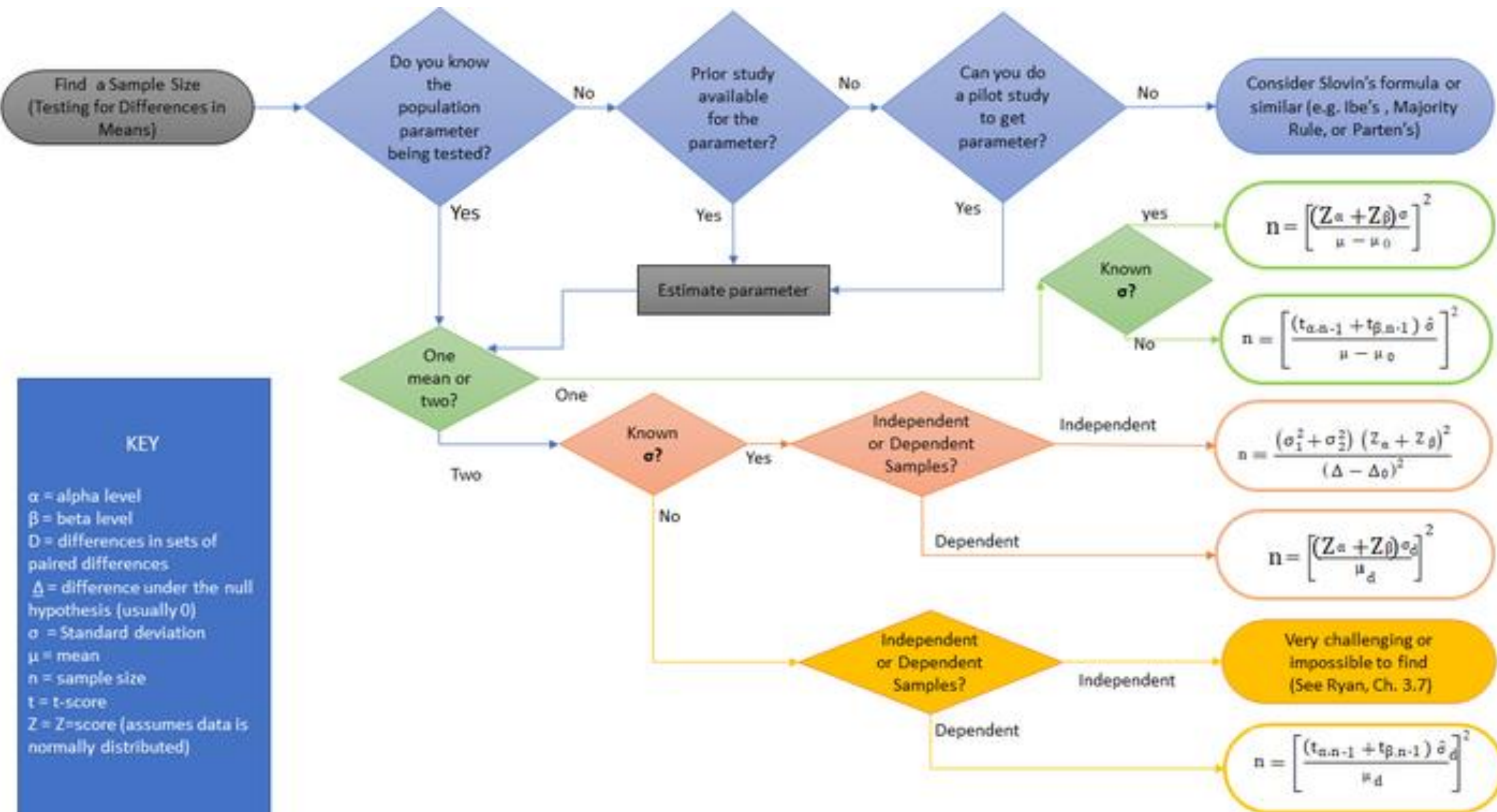
Ensemble Dataset in One Picture



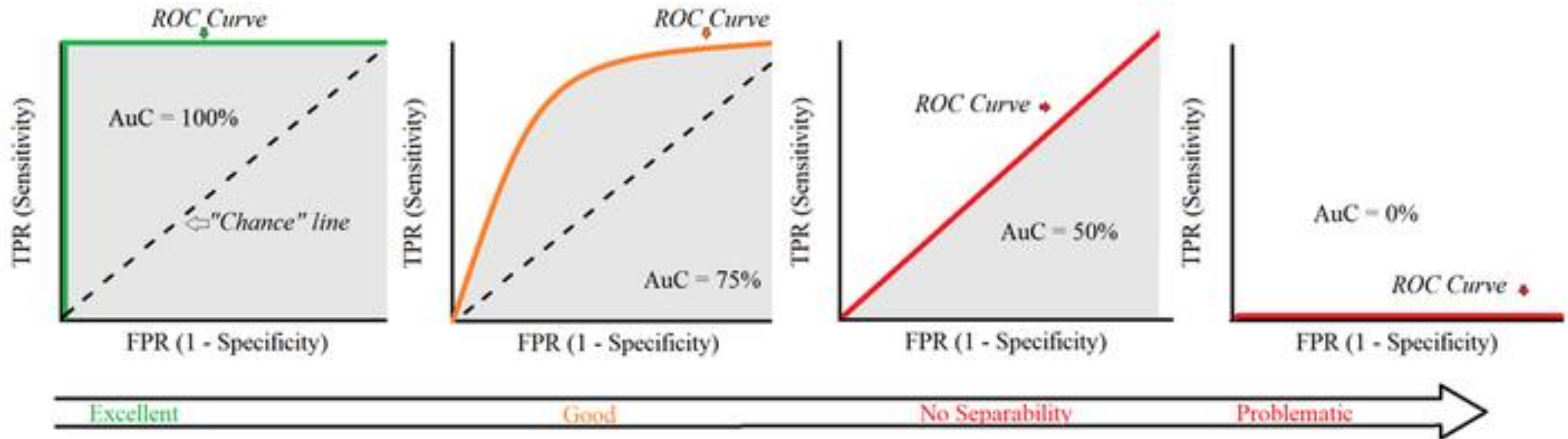
Parametric and Non-Parametric in One Picture



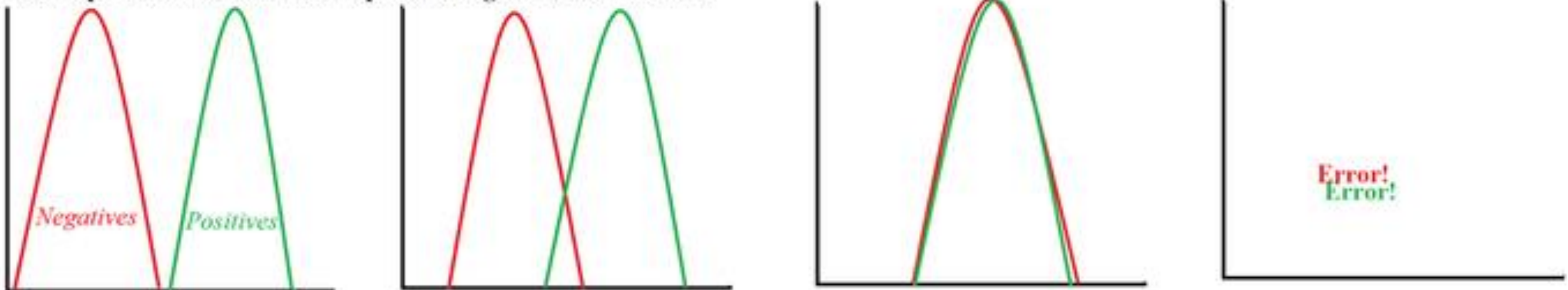
Determining Sample Size in One Picture




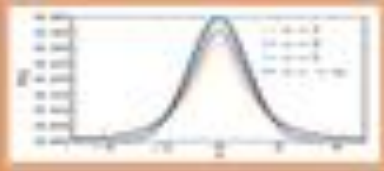
ROC CURVE in One Picture



Overlap = How well the model separates Negatives and Positives



Z-test and T-test in One Picture

	Subtypes	Distribution Used	Population Standard Deviation	Sample Size	Assumptions
Z Test	<ul style="list-style-type: none"> One sample Two independent samples Two proportions 	<ul style="list-style-type: none"> Normal Distribution 	<ul style="list-style-type: none"> Known population standard deviation 	<ul style="list-style-type: none"> Large samples only ($n > 30$) 	<ul style="list-style-type: none"> Interval or ratio data Random selection Independent samples Normality
T Test	<ul style="list-style-type: none"> One sample Paired T-test (dependent means) T-test for independent means 	<ul style="list-style-type: none"> T-Distribution 	<ul style="list-style-type: none"> Unknown population standard deviation 	<ul style="list-style-type: none"> Small samples to large samples 	<ul style="list-style-type: none"> Continuous or ordinal data Random selection Normality Homogeneity of Variance
Notes	<ul style="list-style-type: none"> More than two groups or levels? Use ANOVA instead 	<ul style="list-style-type: none"> The t-distribution is slightly shorter and fatter than the normal 	<ul style="list-style-type: none"> In real life, population standard deviation is usually unknown 	<ul style="list-style-type: none"> Sample size is a rule of thumb For $n > 30$, you can choose either the T or Z 	<ul style="list-style-type: none"> Assumptions can differ for specific tests (e.g. independent vs. dependent means)

Anova in One Picture

	One Way ANOVA	Factorial ANOVA	Two Way ANOVA
Basic Description	Identifies differences between the means of 3+ independent (unrelated) groups	Compares mean differences between groups split on two or more independent variables	Special Case of Factorial ANOVA with 2 factors
Independent Variables (Factors)	1	2 or more (although 3/4 is usually the max due to complexity of interpreting results and higher probability of Type I errors)	2 IVs: Factor A (2 or more levels) crossed with Factor B (2 or more levels)
What is being compared in the test?	Means of three or more IV's groups on a dependent variable (although 2 groups is possible, 3 or more groups is the norm).	Effects of multiple groups of multiple IVs on a dependent variable, and on each other.	Effects of multiple groups of two IVs on a dependent variable, and on each other.
Assumptions	Continuous dependent variable Normality Sample independence Homogeneity of Variance	Continuous dependent variable Normality Sample independence Homogeneity of Variance Categorical independent variables	Continuous dependent variable Normality Sample independence Homogeneity of Variance Categorical independent variables

Predictive Analytics in One Picture

Technique	Purpose	Advantages / Disadvantages	What questions can be answered?
Decision Trees	Predict future classes of data.	Advantages: Easy to implement and understand. Disadvantage: Can be overly simplistic for many problems.	"Which one" (i.e. one discrete variable) or answers to Yes/No questions.
Neural Networks	Cluster and classify by recognizing patterns in data.	Advantage: Outperforms most ML algorithms. Disadvantages: Tough to implement; requires parallel processor.	Pretty much any question (as long as there's sufficient data and some kind of correlation/causation)
Regression (Linear and Logistic)	Estimates relationships between variables, uses those to predict future outcomes.	Advantage: Results are easy to understand. Disadvantage: Limited to linear/logistic data.	How much or how many of a certain target variable?
Time Series	Forecasts continuous variables over time.	Advantage: Easy to understand results. Disadvantage: Limited to time-dependent data.	What is the result going to be at a future data?
Clustering (K-means)	Organization of data into similar groups.	Advantage: Easy to implement. Disadvantage: Can be hard to predict "k", the number of clusters.	What kinds of patterns are in the data?
Factor Analysis	Find explanations for outcomes / correlations.	Advantage: Reduction of data to concise "picture". Disadvantage: Factors can be hard to interpret; information can be lost.	What are the explanations for the themes in the data?

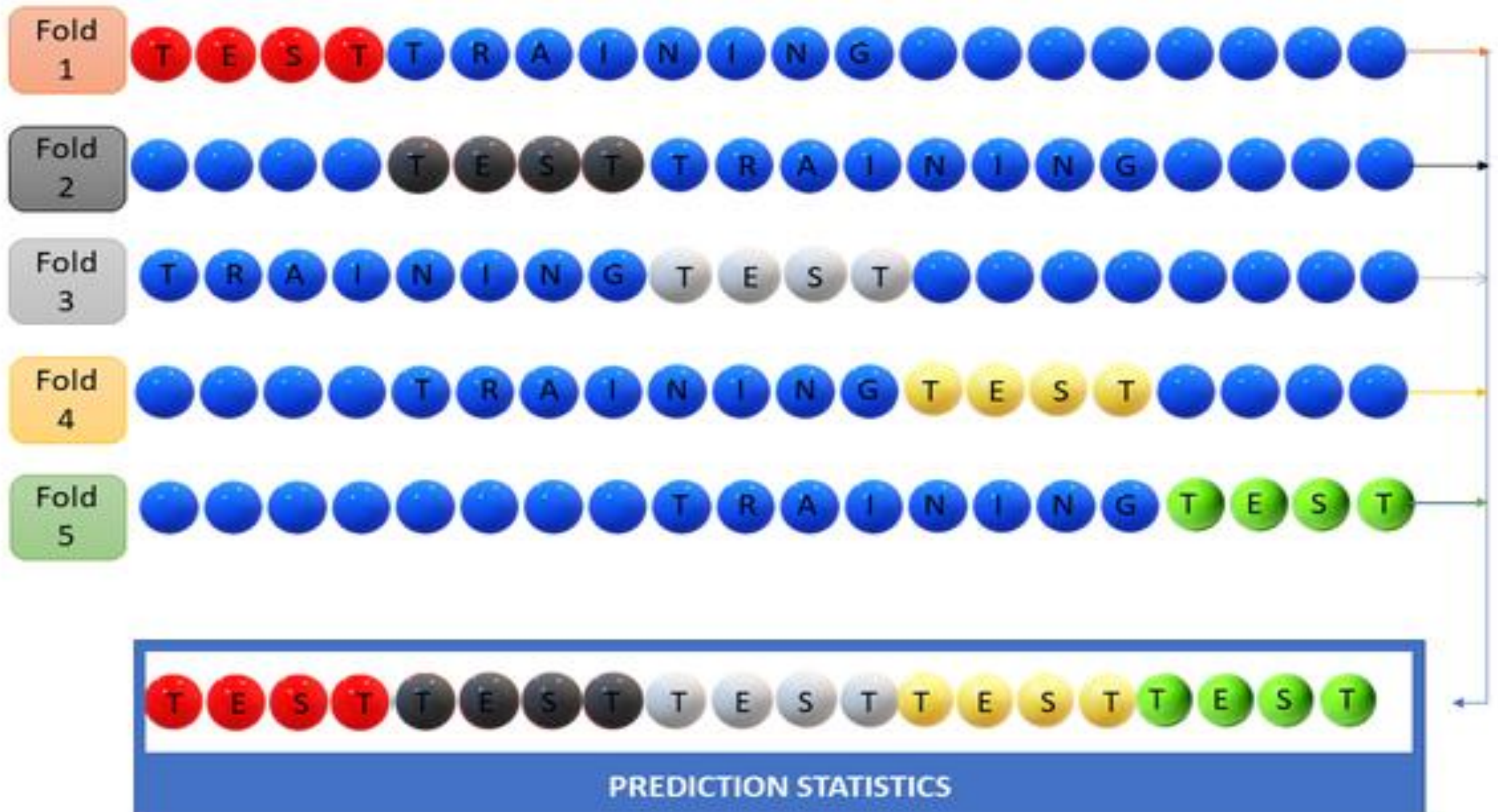
Time Series Method in One Picture

List of Time Series Method

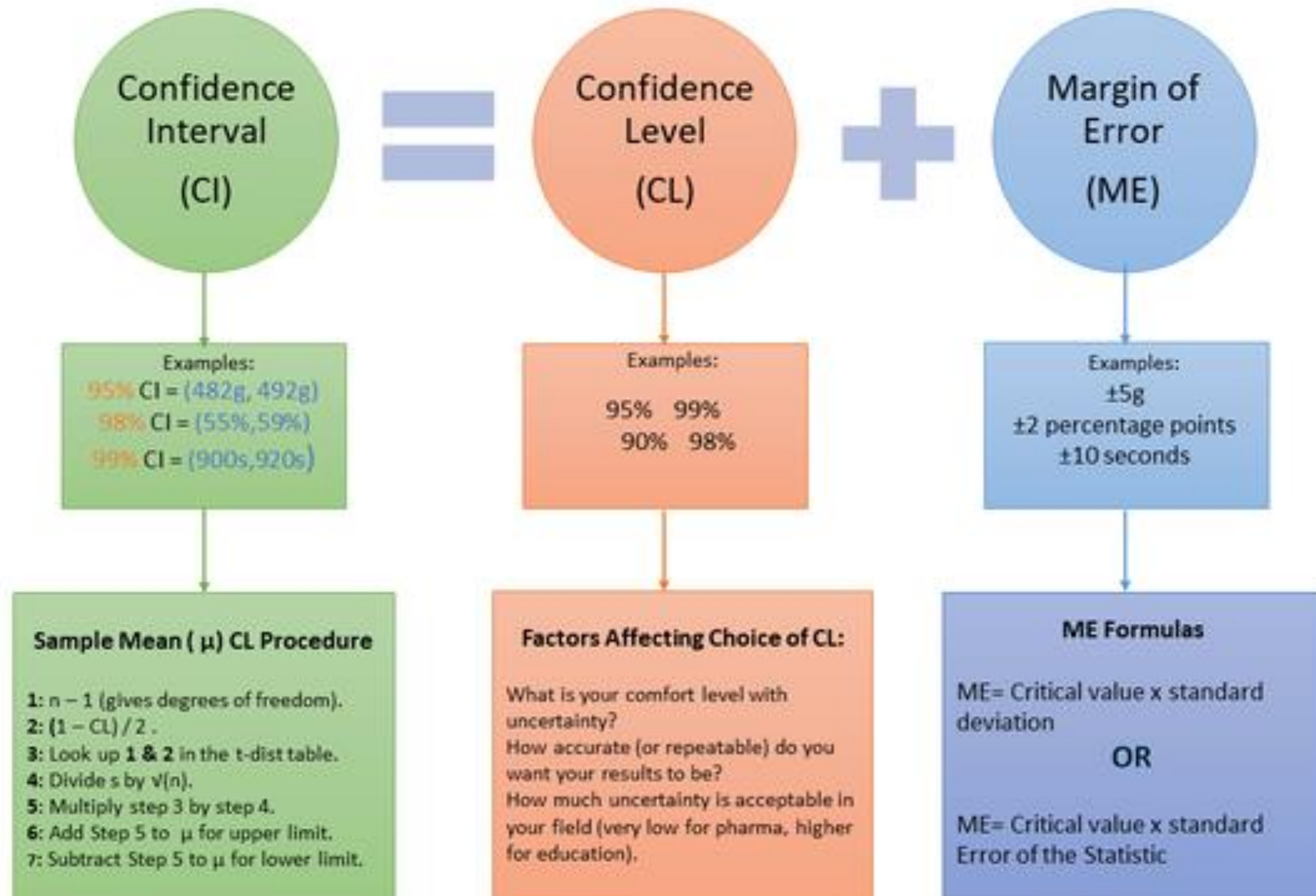


1. Autoregression (AR)
2. Moving Average (MA)
3. Autoregressive Moving Average (ARMA)
4. Autoregressive Integrated Moving Average (ARIMA)
5. Seasonal Autoregressive Integrated Moving-Average (SARIMA)
6. Seasonal Autoregressive Integrated Moving-Average with Exogenous Regressors (SARIMAX)
7. Vector Autoregression (VAR)
8. Vector Autoregression Moving-Average (VARMA)
9. Vector Autoregression Moving-Average with Exogenous Regressors (VARMAX)
10. Simple Exponential Smoothing (SES)
11. Holt Winter's Exponential Smoothing (HWES)
12. Prophet
13. Naive method
14. LSTM (Long Short Term Memory)
15. STAR (Space Time Autoregressive)
16. GSTAR (Generalized Space Time Autoregressive)
17. LSTAR (Logistic Smooth Transition Autoregressive)
18. Transfer Function
19. Intervention Method
20. Recurrent Neural Network
21. Fuzzy Neural Network

Cross Validation in One Picture



Confidence Interval in One Picture



Unsupervised Learning in One Picture

Type

Overview of process

Disadvantages

Advantages

Works well for

K Means



Clustering

Nonhierarchical method that finds mutually exclusive spherical clusters based on distance.

- Requires known "k" clusters; Can be difficult to choose.
- Initial cluster choices and order of data strongly affect results.
- Can be difficult to reproduce results due to random initial "centroid" choice.

- Easy to implement.
- Fast (for small k).
- Clusters can be recalibrated.

Big data; Hyper spherical (e.g. 3D sphere).

Hierarchical



Clustering

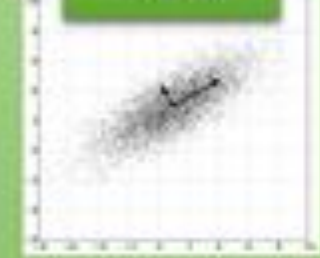
Repeatedly links pairs of clusters until every data object is included.

- Initial seeds, data order have strong effect.
- Merges cannot be undone.
- No statistical / theoretical foundation for results.
- Sensitive to outliers.

- Easy to implement.
- Dendrogram makes for easy visualization of "k".
- Results easily reproducible.

Small to medium data. Performance and execution time increase dramatically for large data sets.

PCA



Dimension Reduction

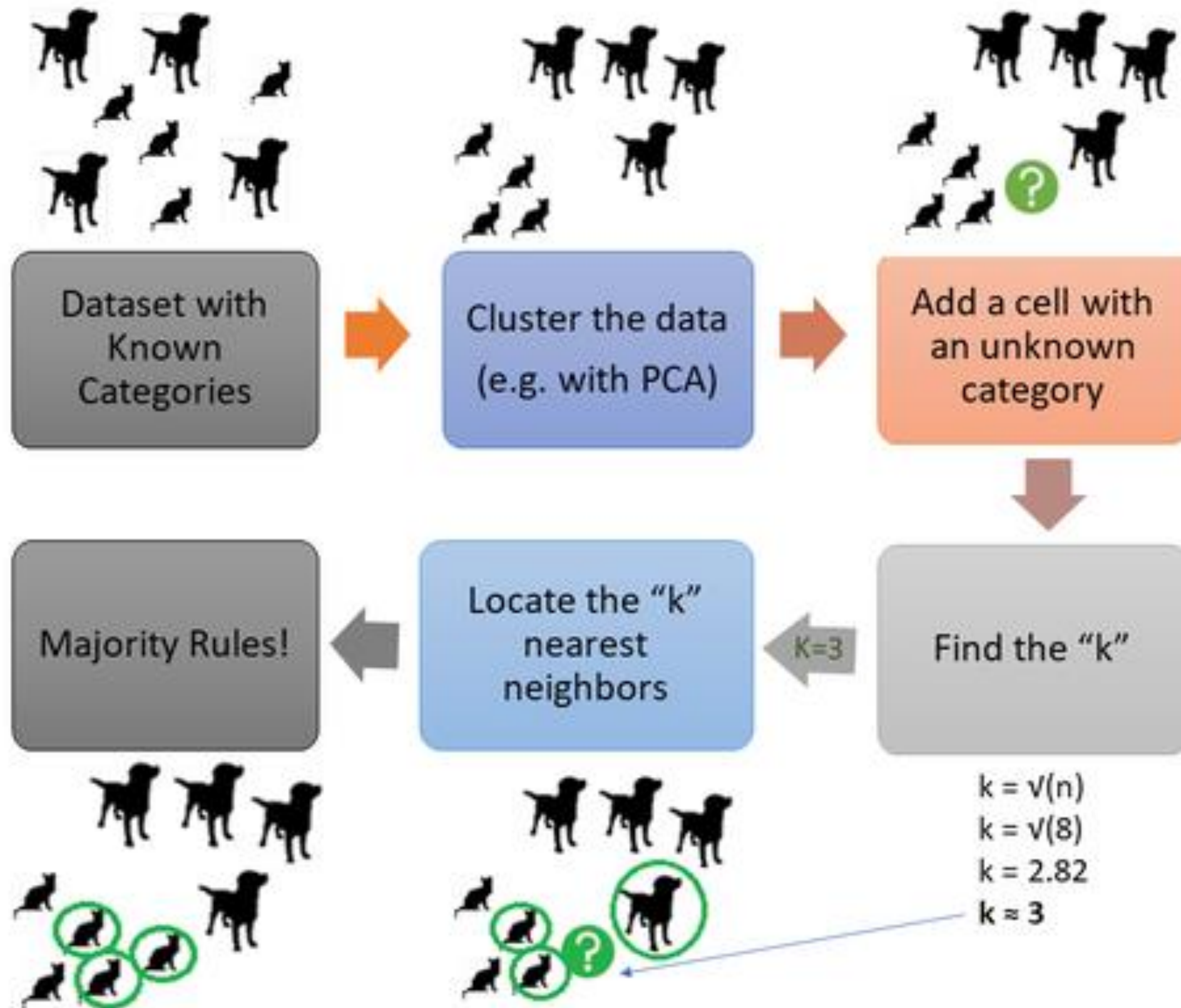
Transforms high dimensional data into low dimensional data using orthogonal transformations.

- Principal components (a linear combination of the original features) can be challenging to interpret and read, compared to the original features.
- Data must be standardized beforehand.

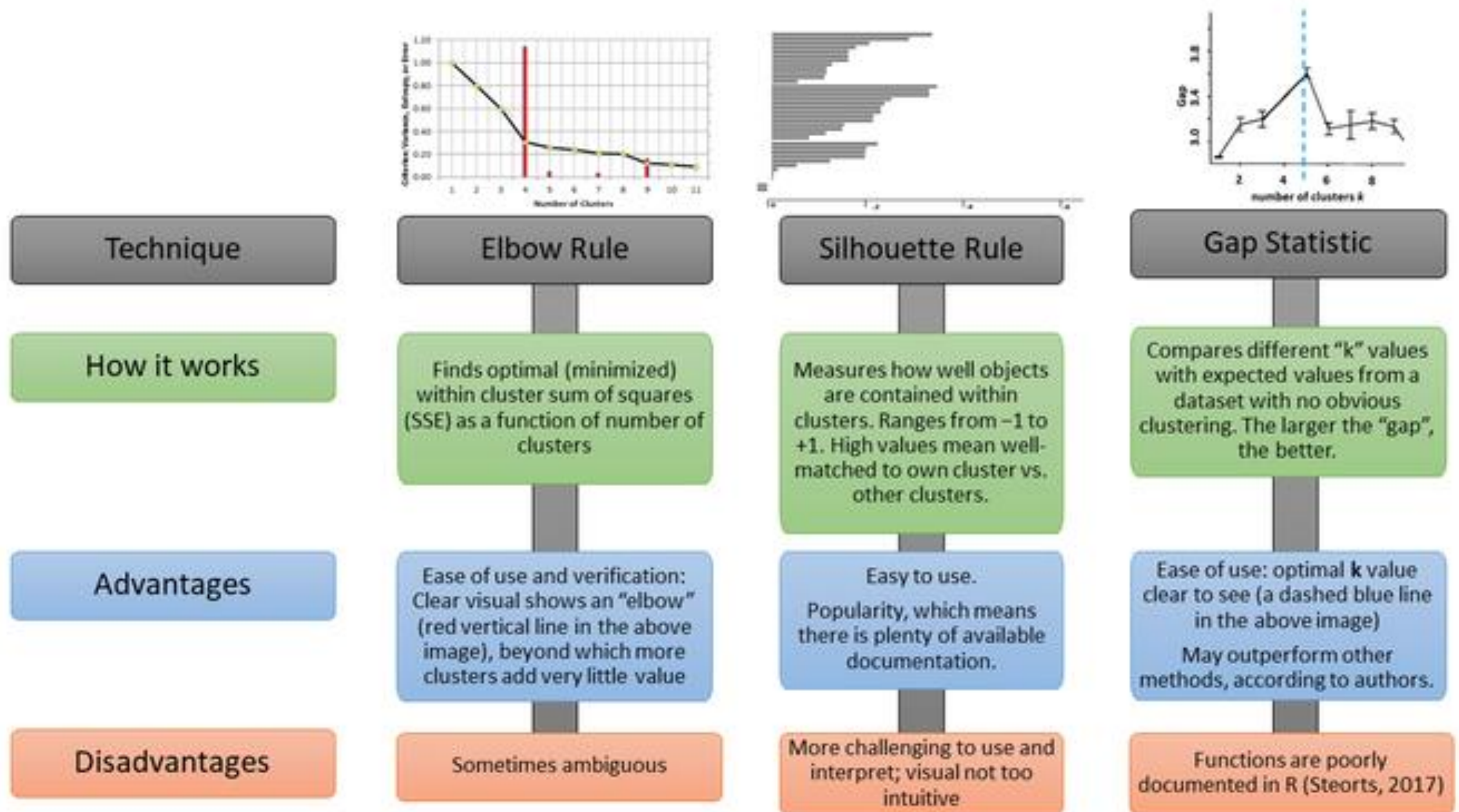
- Good visualization tool.
- Reduces irrelevant or redundant features.
- Reduces overfitting (by reducing features).

Extracting important features (components) from big data.

KNN in One Picture



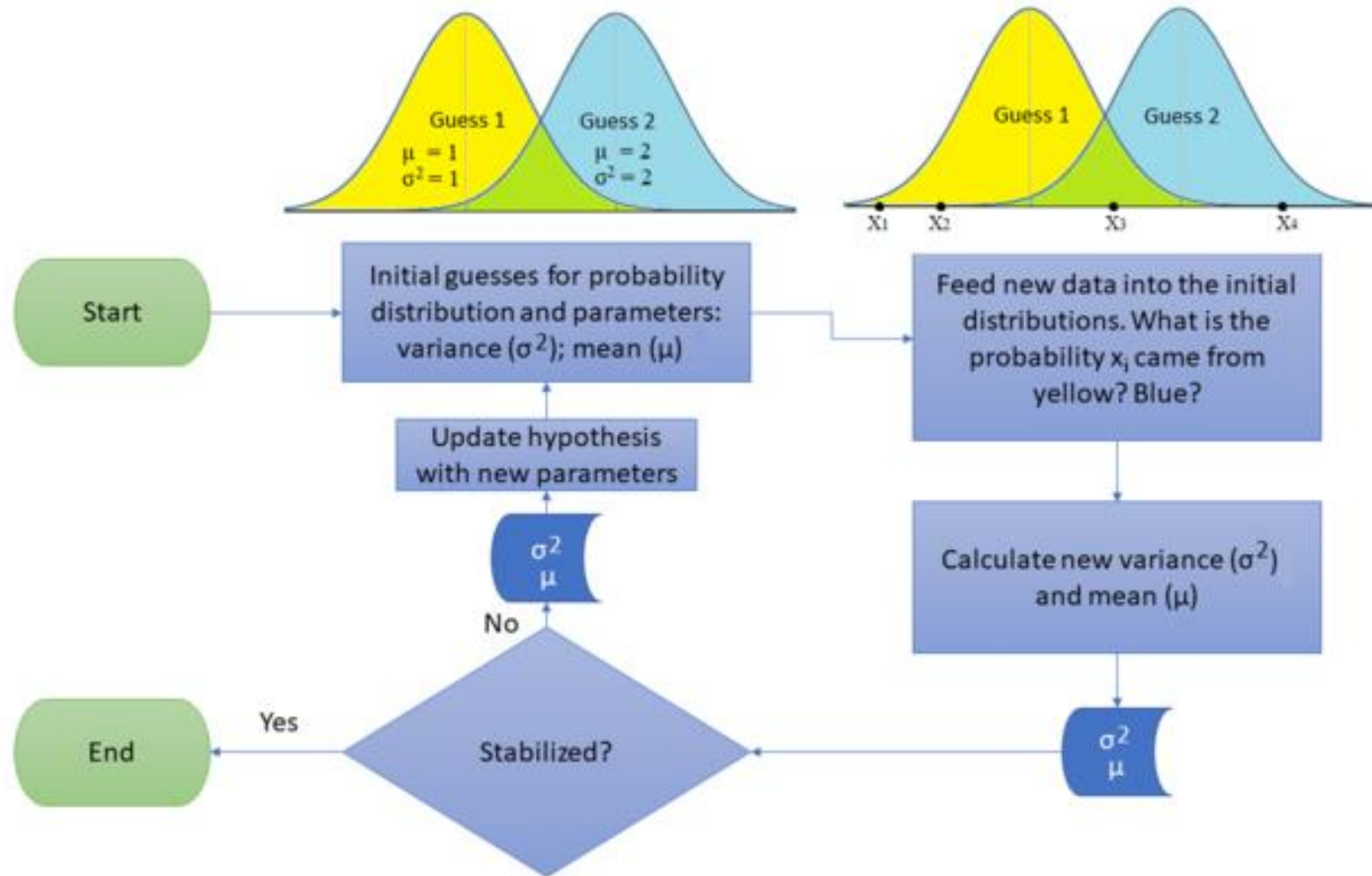
Number of Clusters Selection Method in One Picture



AB Testing in One Picture



EM Algorithm in One Picture

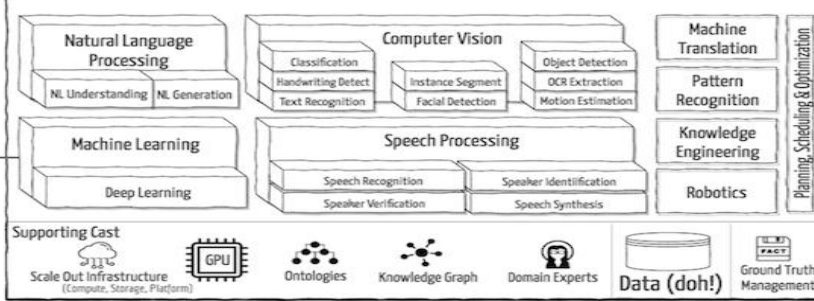


DEMYSTIFYING ARTIFICIAL INTELLIGENCE

by Swami Chandrasekaran (@swamichandra)

Building Blocks

These are the core technologies and essential building blocks that are used in the design - build of AI systems. One or many of these are typically combined together to help realize AI systems.



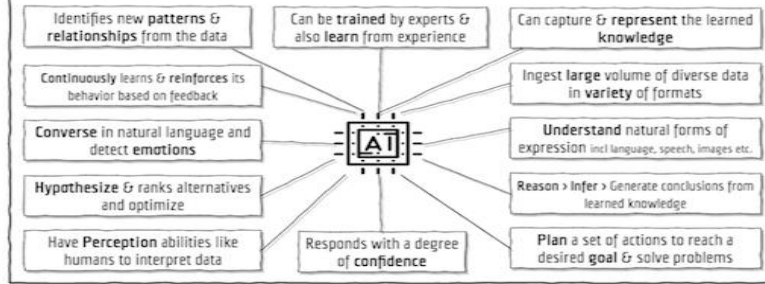
learns & fits to

AI Systems work on various types of data

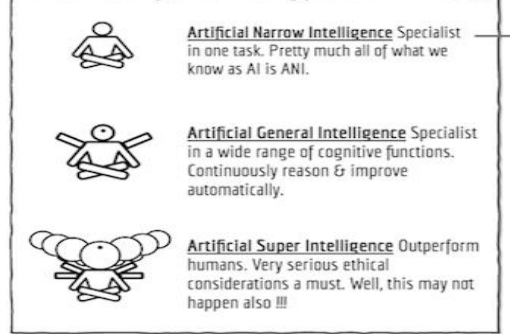


Primary Characteristics of an AI System

AI Systems must exhibit the following characteristics. Call them the DNA, these characteristics are aligned with the key goals of an AI system that include reasoning, knowledge, planning, learning, natural language processing, perception and the ability to move and manipulate objects.



Theoretically, different types of AI



Eliminate Unconscious Bias

Author: Swami Chandrasekaran

Augment Human Intelligence

Explainable & Evidence Driven

"ARTIFICIAL INTELLIGENCE IS THE SCIENCE AND ENGINEERING OF MAKING INTELLIGENT MACHINES" - JOHN MCCARTHY

Used Responsibly & Ethically

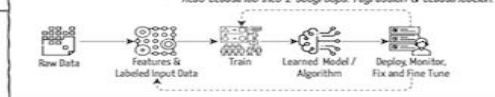
Transparent & Accountable

Respecting Privacy of the individual

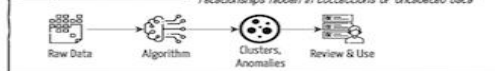
Types of Learning Approaches

3 major approaches through which machines learn without being explicitly programmed

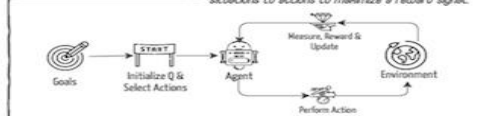
Supervised Learning Learning from set of labeled data provided by an expert. Also classified into 2 subgroups: regression & classification.



Unsupervised Learning Uncovering inherent structures, patterns and relationships hidden in collections of unlabeled data



Reinforcement Learning Learning what to do in an environment, and how to map situations to actions to maximize a reward signal.



Top Algorithms

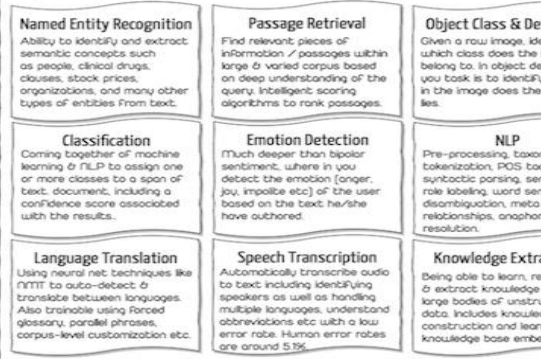
- Linear Regression
- Logistic Regression
- Naive Bayes
- Hidden Markov
- XG Boost
- SVM

- Deep Learning Algorithms**
 - Convolutional Neural Nets (CNN)
 - Deep Belief Networks (DBN)
 - Auto Encoders
 - LSTM
 - RBM
- K-Means Clustering
- Principal Component Analysis
- Anomaly Detection
- Singular Value Decomposition (SVD)

- Markov Decision Process
- Q-Learning
- Temporal Difference Learning

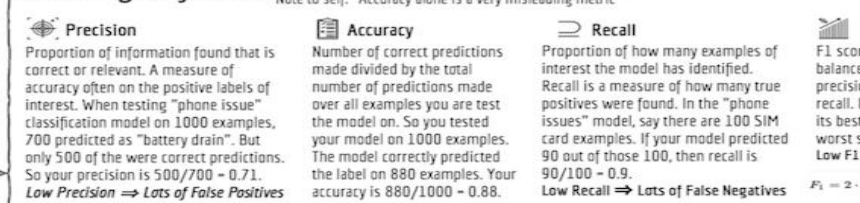
Most common AI Workloads / Tasks

These are some of the most common tasks that AI systems are taught to perform and also learn continuously over time. Each of these tasks MUST provide responses with a degree of confidence or probability associated with them along with underlying evidence.



Measuring AI Systems

The most common measures and other approaches to evaluating the performance of AI systems. Note to self: "Accuracy alone is a very misleading metric"



Other Evaluation Criteria / Approaches	Bias - Variance Tradeoff Homogeneity Score	Akaike Information Criterion(AIC) Completeness Score	Passage relevancy rating Jaccard Similarity Score	Average Steps per Conversation Inter Annotator Agreements	Problem Benchmarks Conversation abandonment rates	Hamming Loss Peer Confrontation	Adjusted Rand Index (ARI) In-Sample Error	Log-Loss Lift
--	--	--	---	---	---	---------------------------------	---	---------------

Receiver operating characteristics (ROC): Technique for visualizing, organizing and selecting classifiers based on their performance. ROC curves show the tradeoff between false positive and true positive rates. We plot "True Negative Rate" on X-axis and "True Positive Rate" on Y-axis.

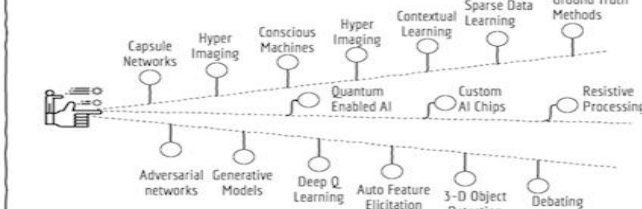
Turing Test: Suppose that we have a person, a machine, and an interrogator. The interrogator is in a room separated from the other person and the machine. The object of the game is for the interrogator to determine which of the other two is the person, and which is the machine.

Chinese Room: Does the machine "literally" understand Chinese? Or is it merely simulating the ability to understand Chinese? Programming a digital computer may make it appear to understand the language but does not produce real understanding.

API's, Libraries & Frameworks



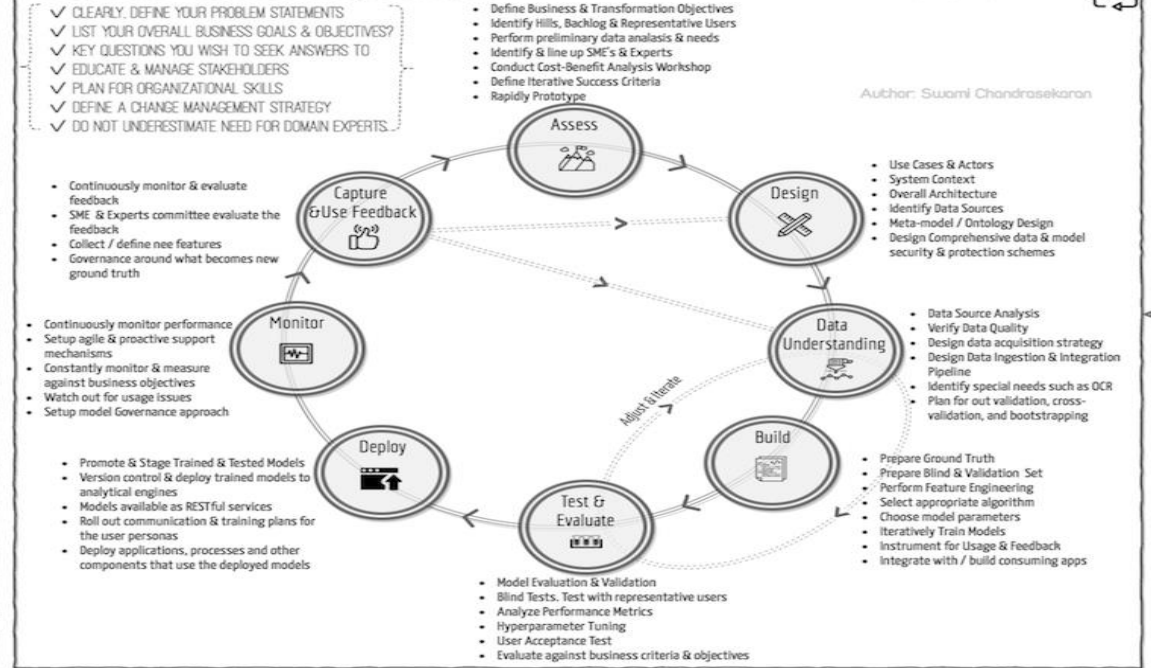
Pragmatically, what's next for AI?



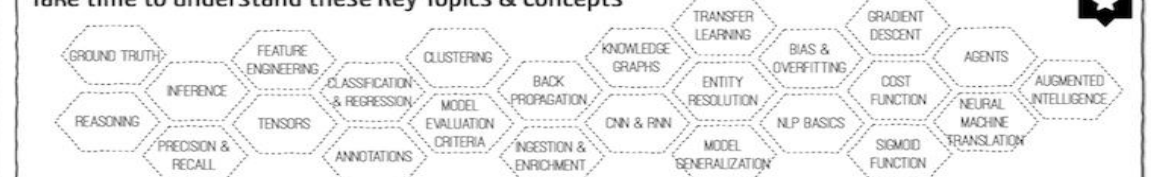
Common Examples of AI Systems in Action

Review Contract Docs	Audio Call Transcription	Conversational Bots	Search	Customer 360	Threat Prevention
Product Recommendations	Sentiment Analysis	Precision Marketing	Equipment Maintenance	Abnormalities from MRI, x-ray etc	Visual Search
Fraud Detection	Customer Churn Prediction	Enforce Compliance	Malware & Spam Detection	Insurance Claims Damage Assessment	Case Law
Social Listening	Employee Attrition Prediction	Lead Generation	Identify Skin Cancer from Lesions	Multi-modal, no touch Commerce	Amazon Alexa
Taxi Dispatch	Home Automation Control	Sort Vegetables	Patient Readmission Prediction	AI enabled RPA	Drug Repurposing

DevOps for AI - How are these systems built?



Take time to understand these Key Topics & Concepts



Number Representation Systems in One Picture

	Logistic Map ($p = 1$)	Logistic Map ($p = 0.5$)	Nested Square Root	Base b ($b > 1$, integer)	Base b ($b > 1$, not an integer)	Continued Fraction
Support domain for $x = x(1), x(n)$ and $g(x)$	$[0, 1]$	$[0, 1]$	$[1, 2]$	$[0, 1]$	$[0, 1]$	$[0, 1]$
Support domain for digits $a(n)$	$\{0, 1\}$	$\{0, 1\}$	$\{0, 1, 2\}$	$\{0, 1, \dots, b-1\}$	$\{0, 1, \dots, \lfloor b \rfloor\}$	$1, 2, 3, \dots$
$g(x)$ [$x(n+1) = g(x(n))$]	$4x(1-x)$	$\sqrt{4x(1-x)}$	$x^2 - \lfloor x^2 \rfloor + 1$	$bx - \lfloor bx \rfloor$	$bx - \lfloor bx \rfloor$	$1/x - \lfloor 1/x \rfloor$
$h(x)$ [$a(n) = h(x(n))$]	$\lfloor 2x \rfloor$	$\lfloor 2x \rfloor$	$\lfloor x^2 \rfloor - 1$	$\lfloor bx \rfloor$	$\lfloor bx \rfloor$	$\lfloor 1/x \rfloor$
$x = f(\{a(n)\})$	Unknown	Unknown	$\sqrt{a(1)} + \sqrt{a(2)} + \dots$	$\sum_{k=1}^{\infty} \frac{a(k)}{b^k}$	$\sum_{k=1}^{\infty} \frac{a(k)}{b^k}$	$\frac{1}{a(1) + \frac{1}{a(2) + \dots}}$
Digits distribution $P(a(n) = k)$	Uniform on $\{0, 1\}$	$\sqrt{2}/2$ if $k = 1$	$r(k+2) - r(k+1)$ $r(k) = \sqrt{5\sqrt{k} - 1}$	Uniform on $\{0, 1, \dots, b-1\}$	Non uniform	$\log_2 \left(\frac{(k+1)^2}{k(k+2)} \right)$
Equilibrium distribution $P(x(n) < y)$	$\frac{1}{\pi} \int_0^y \frac{1}{\sqrt{x(1-x)}} dx$	$1 - \sqrt{1-y}$	$-2 + \sqrt{5y-1}$	Uniform on $[0, 1]$	Non uniform	$\frac{\log(1+y)}{\log 2}$
Correlation between $x(n+1)$ and $x(n)$	0	-1/2	Non zero	1/b	Non zero	Non zero
Correlation between $a(n+1)$ and $a(n)$	0	-1/4	Non zero (but close)	0	Non zero	$E[a(n)]$ is infinite

The equilibrium distribution is the one satisfying $P(X < y) = P(g(X) < y)$. The statistical properties listed here, are valid for almost all seeds x , except for a set of seeds of measure 0 (depending on the system), known as bad seeds. This document was produced by www.datasciencecentral.com.