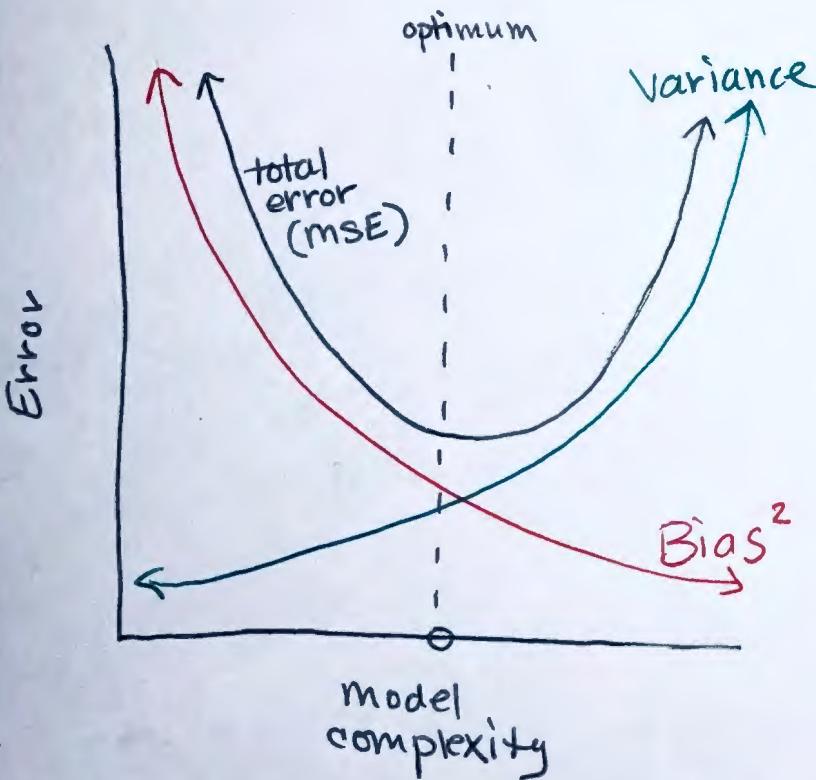


# Bias - Variance Trade off



**Bias** = error from missing relevant relations btwn features + output  
e.g. underfitting

- High in: linear algorithms  
resampling (model selection)  
simpler models

• The Fix: increase model complexity, reduce regularization

**Variance** = error from sensitivity to small fluctuations (noise) in training data

e.g. overfitting

- High in: decision trees  
SVM  
KNN

• The Fix: bagging, bootstrapping, regularization, more traindata

# Regularization

**Def:**

a process of adding additional info in order to solve a problem  
→ adding constraints or penalty to model

Why use it?

- reduce overfitting + increase generalization

How it works

$$\text{minimize}(\text{Loss(Data|model)} + \lambda \text{complexity(model)})$$

$$\text{e.g. L}_2 \text{ Regularization term} = \sum (\text{weight}_n)^2$$

weights  
of  
features

What methods?

- Ridge + LASSO regression

# Correlation

Def.

direction + strength of linear relationship between two variables

value lies between -1 and 1

Methods

Pearson's R

$$r = \frac{\text{covariance}(x, y)}{\text{stdv}(x) \cdot \text{stdv}(y)}$$

Spearman's Rank

-assesses monotonic relationships

(one var increases as the other var does)

Point-Biserial

-assesses relationship between continuous + binary vars

! Note: Spurious correlation are things with a mathematical relationship, but no real-life relationship

# XGBoost

##

##



accuracy , efficiency , feasibility

Def.

a linear model + tree learning algo that does parallel computations on a single machine ; a fast implementation of gradient boosted trees

## Ridge Regression

L2 Regularization method that shrinks the size of the coefficients

Encourages :

- Coefficients that are closer to true population parameters
  - equal shrinkage

## Combats:

- overfitting
  - multicollinearity
  - sparseness  
(avoids it)
  - high variance

### Other Uses:

## Cons:

- sensitive to outliers

R code :

```
> library(glmnet)
```

```
> model = cv.glmnet(x, y, alpha = 0)
```

minimize →

$$\sum_{i=1}^n \left( y_i - \underbrace{\sum_j x_{ij} \beta_j}_{\text{sum of squared error}} \right)^2 + \lambda \underbrace{\sum_{j=1}^p \beta_j^2}_{\text{penalty}}$$

$y_i$  = Response data point     $x_{ij}$  = covariate matrix     $\beta$  = coefficients

# LASSO Regression

a.k.a. Least Absolute Shrinkage & Selection Operator

L1 Regularization method that shrinks data points towards a central value, like the mean

Encourages:

- simple, sparse model
- fewer parameters by elimination
- shrinking  $\beta$  to exactly 0

minimize  $\rightarrow$

Combats:

- overfitting
- multicollinearity
- high variance

Other Uses:

- automate feature selection

R code:

```
> library(glmnet)  
> model = cv.glmnet(  
    x, y, alpha=1)
```

$$\sum_{i=1}^n \left( \hat{y}_i - \sum_j x_{ij} \beta_j \right)^2 + \lambda \sum_{j=1}^p |\beta_j|$$

sum of squared errors      penalty

tuning parameter  
of strength  
of penalty

$y_i$  = response var  
data point     $x_{ij}$  = covariate matrix     $\beta$  = coefficients

# Bootstrapping

Def:

a resampling method that uses sampling with replacement to simulate multiple samples/experiments

## Uses

- estimate C.I.
- estimate standard errors

## Benefits

- non-parametric (e.g. distribution-free)
- recreates any # of resamples

## R Code

```
> library(boot)  
> boot(data, statistic,  
      R = # of resamples)
```

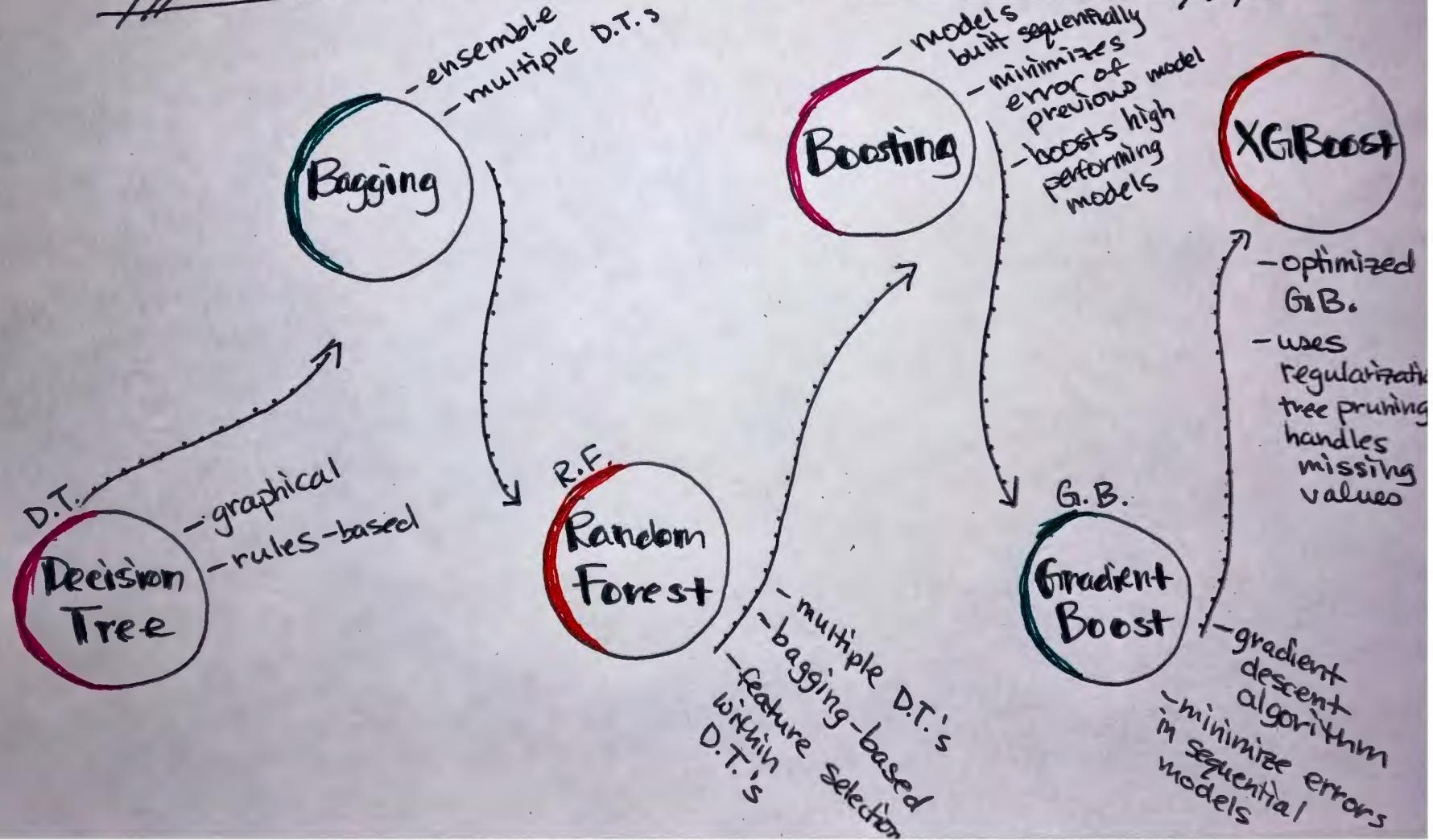
# Cross - Validation

Def : a resampling method used to evaluate a model's performance (generalization)

Validation Set	Leave One Out CV	K-fold CV
<p>training      validation</p> <ul style="list-style-type: none"><li>• simplest approach</li><li>• can overestimate accuracy</li></ul>	<ul style="list-style-type: none"><li>• highly unbiased</li><li>• high variance</li><li>• computationally expensive</li></ul>	<ul style="list-style-type: none"><li>• typical K=5 or 10</li><li>• middle of the road bias + variance</li><li>• not as computationally expensive</li></ul>

# Tree-Based Algorithms

\*non-linear\*



# Type I + II Errors

"false positive"

## Type I Error

- mistaken rejection of an actually true null hypothesis

a.k.a saying 2 groups are different when they are not

Tolerance for this error set by your **alpha  $\alpha$  level**

confidence level  
 $= 1 - \alpha$

e.g.  $\alpha = 5\%$ , 5% chance of a Type I error

"false negative"

## Type II Error

- mistaken acceptance of an actually false null hypothesis

a.k.a saying 2 groups are not different when they are

Tolerance for this error set by your **power level**

e.g. power = 80%, 20% chance of Type II error

# Linear Regression

\*supervised\*

Def: a method of fitting a linear trend line to data points in order to predict future values

Goal minimize the sum of squares of the difference between the observed data + predictions ↗ variance

## Assumptions

linear relationship ; features are not correlated ; errors are uncorrelated + homoscedastic + normal

## Validation

$$y_i = \beta_0 + \beta_1 x_i \dots + \epsilon$$

$R^2$  = goodness of model fit ; explanatory power of features

adj.  $R^2$  = same as  $R^2$  but penalizes for more features

AIC = used for model selection ; estimator of prediction error

Residual plots + Durbin-Watson = check if residual assumptions are met

# Naive Bayes

\*supervised\*

111

111

Def:

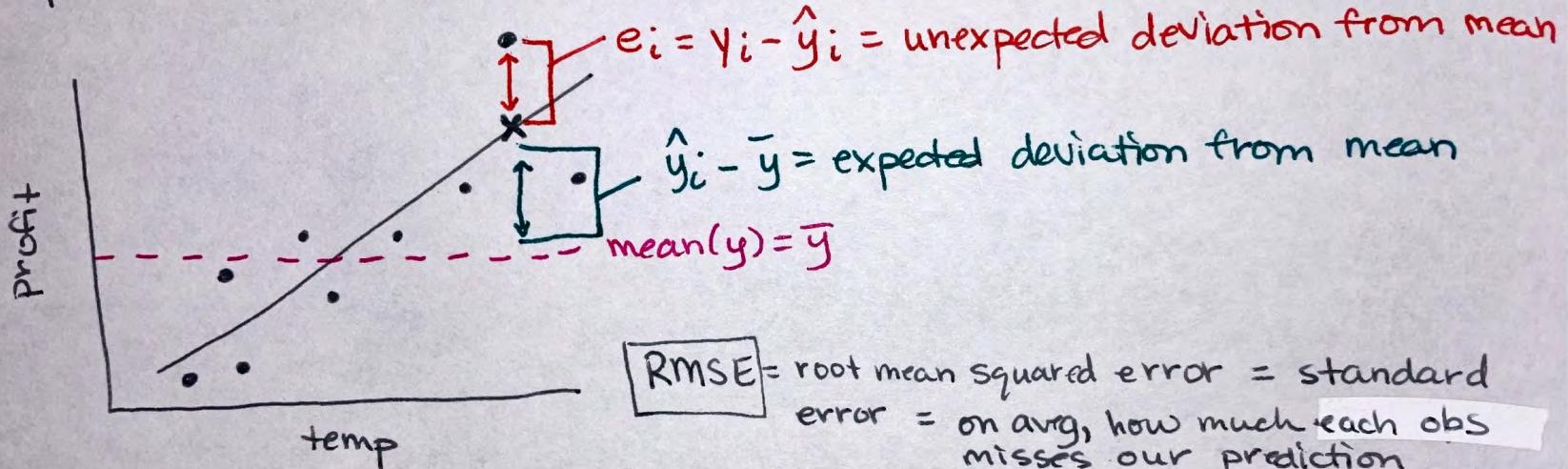
probabilistic classifier based on Bayes' Theorem with strong independence assumptions between features

- not iterative, so very scalable
- uses MLE
- used often for text classification & recommendation systems

$$P(A | x_1, \dots, x_n) = P(x_1 | A) \cdot P(x_2 | A) \cdot P(x_3 | A) \dots \cdot P(x_i | A) \cdot P(A)$$

# Error

pertains to how well a model explains the true population



**SSR** = sum squared regression =  $\sum (\hat{y}_i - \bar{y})^2$ , also called ESS

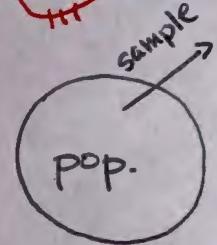
**SSE** = sum squared error =  $\sum e_i^2$ , also called RSS

**SST** = SSR + SSE = sum squared total, also called TSS

**MS** = mean square = sum squared errors / degrees freedom = unbiased estimate of error variance  
↑ use this to calculate F-statistic to accept/reject Null that all regression coefficients equal zero

# C.I. - Proportion

One Sample



$$n \rightarrow \hat{p} = \text{Sample proportion}$$

## Steps:

- 1) choose confidence level, e.g. 95%
- 2) obtain critical value for that confidence level, e.g. z-score
- 3) calculate C.I.

$$\hat{p} \pm \text{z-score} \left( \sqrt{\frac{\hat{p}(1-\hat{p})}{n}} \right)$$

*if you know the true population stdv., use that*

Two samples

$$P_1 \rightarrow n \rightarrow \hat{p}_1$$

$$P_2 \rightarrow n \rightarrow \hat{p}_2$$

$$\text{confidence interval for } (P_1 - P_2) = (\hat{p}_1 - \hat{p}_2) \pm z \left( \sqrt{\frac{\hat{p}_1(1-\hat{p}_1)}{n_1} + \frac{\hat{p}_2(1-\hat{p}_2)}{n_2}} \right)$$

standard error  $\rightarrow \sigma_{\hat{p}_1 - \hat{p}_2}$

\*unsupervised\*

# K-Means

hard clustering

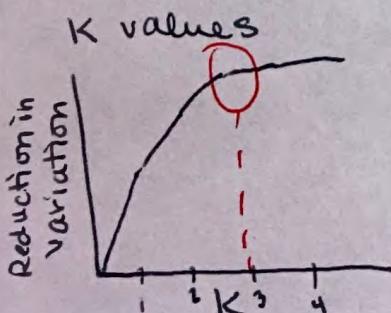
**Def:** unsupervised clustering algorithm that groups data points based on the nearest mean (cluster center)

**Goal:** minimizes within-cluster variance (squared Euclidean distance)

**Steps:**

① Define # of clusters

- use elbow method to plot total variation within cluster for diff. K values

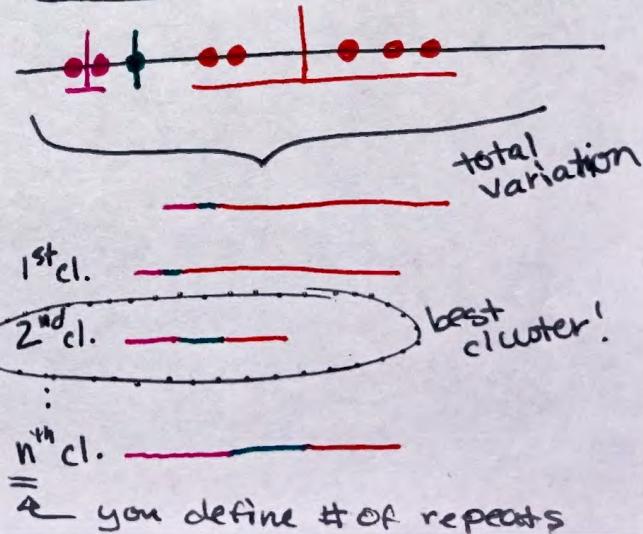


② Randomly select K datapoints to be initial clusters



- assign rest of points to cluster based on shortest distance
- after assigned, calculate mean of each cluster & re-assign based on point distance to cluster mean

③ Repeat clustering method (#2) until best total variation within cluster is reached



unsupervised\*

# Mixture Models & EM Algo

soft clustering

Def: unsupervised clustering method based on probability distributions  
(Gaussian or multinomial)

Goal: discover unknown mean & variance of distributions

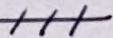
Steps:

- ① Start by placing the  $K$  distributions somewhere on fine
- find the  $(\mu, \sigma^2)$  for each one
  - calculate probability that data point belongs to each distro
  - softly assigns point to group based on likelihood to be in that group

- ② After first probabilities found, adjust distro  $(\mu, \sigma^2)$  to fit points & iterate process
- maximize the log-likelihood function
- ↑  
the most probable assignments

# Hierarchical Clustering

\* unsupervised\*



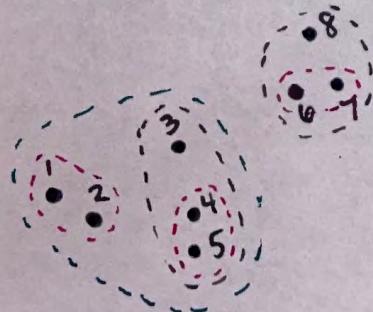
Def: an algorithm that builds a hierarchy of clusters based on a defined distance metric

e.g. Euclidean distance, squared Euclidean distance, Manhattan distance, maximum distance, mahalanobis distance

## Agglomerative

"bottom-up"

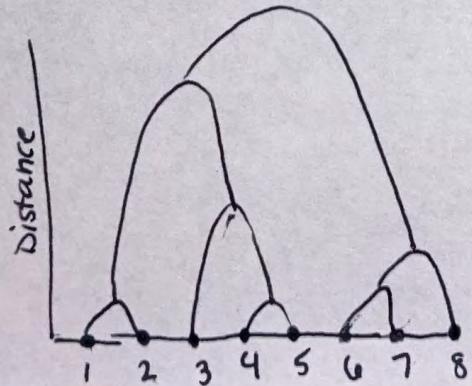
- each obs. starts as its own cluster



## Divisive

"top-down"

- all obs. start in 1 cluster



# DBSCAN

\*unsupervised\*

##

##

Def:

clustering algorithm that groups any shape

① a point is randomly selected as starting point + based on the "minimum points" and the declared epsilon (distance)

② each point in the initial cluster will broadcast out their perimeter, looking for new point members

border points  
part of cluster, but not within  
epsilon

③ border points broadcast out their perimeter to gain new members

④ method is repeated for any remaining, unassigned points

# Non-Negative Matrix Factorization

Def:

a dimension reduction algorithm which decomposes a matrix  $V$  into 2 lower dimension matrices

$$\begin{matrix} W \\ \boxed{\quad \quad \quad} \end{matrix} \times \begin{matrix} H \\ \boxed{\quad \quad \quad} \end{matrix} = \begin{matrix} V \\ \boxed{\quad \quad \quad} \end{matrix}$$

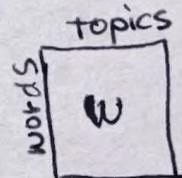
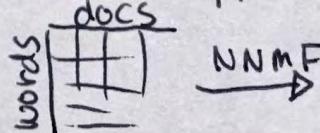
each column  
is a "basis" element,  
or a component that  
comes up many times  
in the data points

tells you how  
to reconstruct  
the original data  
point, given the  
basis elements

each element in matrix  $\geq 0$



when applied to text & documents



$W$  matrix columns represent topics

# Multicollinearity

Def:

the state of independent variables being correlated to each other in regression ; not a good thing ; violation of assumptions

Why is it bad?

- muddies interpretation of marginal effects of each variable , since "holding other variables constant" does not apply
- difficult to determine which vars actually had an effect on  $y$
- inflates variance of affected variables  $\rightarrow$  increases standard errors & p-values

Detection

- correlation of independent variables
- VIF

The Fix

- Do nothing!
  - OK if model is only being used for prediction (doesn't affect model fit)
- remove correlated vars
- combine vars or use PCA

Applicable to

- linear regression
- generalized linear model
- logistic regression
- Cox regression

# Margin of Error

confidence interval }  $\mu \pm z\text{-score} \cdot \text{se}$   
margin of error

mean

$$ME = z\text{-score} * \sqrt{\frac{s^2}{n}}$$

Proportion

$$ME = z\text{-score} * \sqrt{\frac{p(1-p)}{n}}$$

# Precision + Recall

///

Def:

classification metrics

///

Accuracy

tells us  $\frac{\text{total correct}}{\text{total obs.}}$ , but can have issues with imbalanced classes

## Precision

$$= \frac{\# \text{ true positives}}{(\# \text{ true positives} + \# \text{ false positives})}$$

e.g.  $\frac{\text{correct results}}{\text{all returned results}}$

Want high for →

When there's a high cost of false positives  
(pregnancy test, recommendation)

## Recall

$$= \frac{\# \text{ true positives}}{(\# \text{ false negatives} + \# \text{ true positives})}$$

e.g.  $\frac{\text{correct results}}{\text{results that could have been returned}}$

Want high for →

when there's high cost of false negatives  
(identifying fraud or cancer)

# Expected Value & Variance

111

Discrete

$$E(x) = \sum x_i \cdot p(x_i)$$

- weighted average
- for Binomial =  $n \cdot p$

$$\sigma^2 = E[(x-\mu)^2]$$

$$= \sum ((x_i - \mu)^2 \cdot p(x_i))$$

- for Binomial =  $n \cdot p \cdot (1-p)$

Continuous

$$E(x) = \frac{\sum x_i}{n}$$

$$\sigma^2 = \frac{\sum (x_i - \mu)^2}{n-1}$$

# Correlated Topic Modeling

def: a probabilistic approach to infer latent topics of a document ; allows latent topics to be correlated

application: derive themes or topics within book (or corpus of books)

assumptions:

each document is comprised of multiple topics, made of multiple words  
order of word doesn't matter

word distribution per topic & topic distro per document is logistic normal

# Word Network Topic Modeling

Def: method to identify topics within short text using network graphs

## Method

- 1) create word network of document text + apply an edge weighting of occurrence frequency of word pairs
- 2) apply LDA to the graph's adjacency matrix
- 3) LDA gives you topics for each word + will need to be mapped back to documents

# tf-idf

term frequency - inverse document frequency

**Def:** measure of how important a word is to a document in a collection (corpus)

$$\rightarrow \text{idf}(\text{term}) = \ln \left( \frac{n \text{ documents}}{n \text{ documents containing term}} \right)$$

$$\rightarrow \text{tf}(\text{term}) = n \text{ occurrences of term}$$

$$\text{tf-idf}(\text{term}) = \text{tf} \times \text{idf}$$

when low....

e.g. tf-idf = 0

extremely common word

when high....

e.g. tf-idf = 0.009

important words

# P-Value → probability of observed outcome

##

##

Def: the probability of obtaining test results at least as extreme as the results actually observed  
→ used in hypothesis testing

## .....Calculation.....

T, test statistic, provides a single number & follows a distribution determined by the function used to define the test statistic & distro of input data (observations)

e.g. Z-tests → t-statistic ; normal distro w/ known variance

t-tests → t-statistic ; normal distro w/ unknown var.

F-tests → F-statistic; for variance

Chi-squared test → chi-squared - statistic ; discrete data

To calculate a p-value, you need ① null hypothesis, ② test statistic (1 or 2 tailed), ③ data

# the Sign Test

\*non-parametric\*

III

III

**Def:** statistical method to test for consistent differences between pairs of obs.

e.g. weight pre-treatment (x)  
weight post-treatment (y)

## \*When to use\*

when determining if ...  $x > y$

$x = y$

$x < y$

median(x) >  
some #

## Assumptions

- data randomly selected
- Samples must be dependent or paired

! will not work on independent samples

# Assumption of Normality

Inferential statistics rely on the assumption that the values of interest will exhibit a bell-curve distribution if many random samples are taken. → parametric

\* means across samples assumed to be normal, not  
that the sample obs. or population obs. are normal  
*also called sampling distro*

Why do we assume a normal distro? → Central Limit Theorem

Methods that rely on assumption

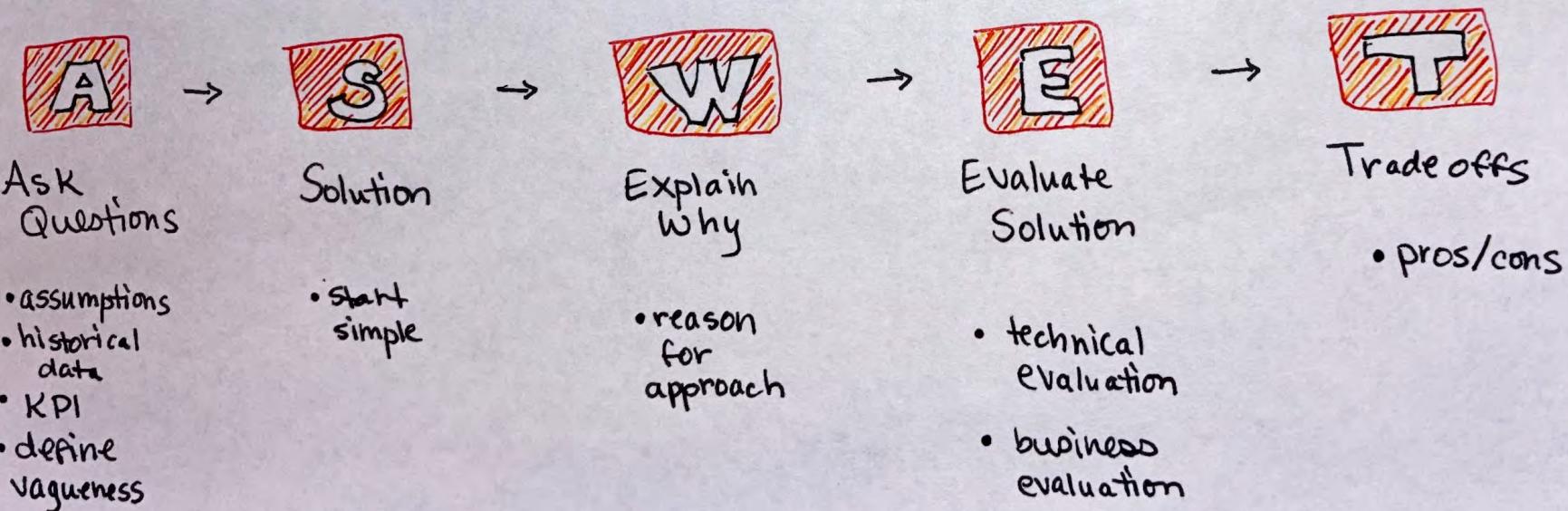
t-tests, ANOVA, regression

Requirements to use this assumption

N sample size needs to be large enough

# Data Science in Business

- 1 What business question are you solving?
- 2 What methods will you use?
- 3 How will you evaluate performance of methods?
- 4 How do you generalize it to similar business problems?



# Endogeneity

Def. when an explanatory variable is correlated with the error term because of measurement error or confounding in omitted or unobserved variables  
relevant for: time series, causal analysis

How to Detect it



How to Correct It



- Instrumental Variables

# Standard Error

+++

++

Def:

the measure of uncertainty in the sample mean

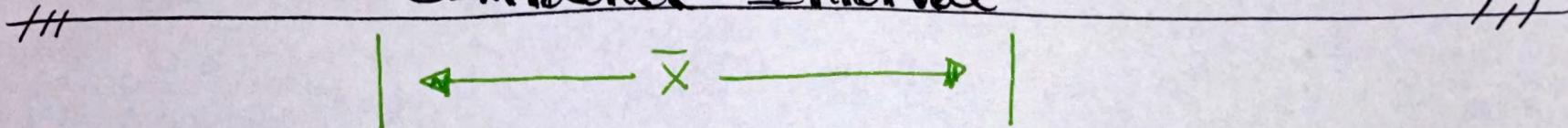
e.g. how precise is our estimate of the mean?

\*decreases as sample size grows

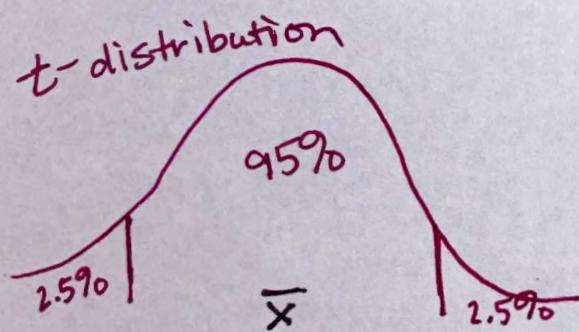
sample  
sd  $\sigma$

sample mean	sample proportion	diff. of means	diff. of proportions
$SE(\bar{x}) = \frac{s}{\sqrt{n}}$	$SE(p) = \sqrt{\frac{p(1-p)}{n}}$	$SE(\bar{x}_1 - \bar{x}_2) = \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}$	$SE(p_1 - p_2) = \sqrt{\frac{p_1(1-p_1)}{n_1} + \frac{p_2(1-p_2)}{n_2}}$

# Confidence Interval



| Def: an interval in which the population mean lies within, with a certain % of confidence



$$\bar{X} \pm SE(\bar{X}) * t\text{-statistic}$$

0.975, n-1  
both part of normal distro

Z score

- use when  $n > 30$  or population is known

vs. T-score

- use when  $n < 30$

# 7 Common Distributions

Name	$f(x)$	$E(x)$	$\text{Var}(x)$	Use Case
Bernoulli	$P \quad x=\text{success}$ $1-p \quad x=\text{failure}$	$p$	$p(1-p)$	single coin toss
Binomial	$\binom{n}{x} p^x (1-p)^{n-x}$	$n \cdot p$	$n \cdot p \cdot (1-p)$	$n$ Bernoulli trials
Geometric	$p(1-p)^x$	$\frac{1-p}{p}$	$\frac{1-p}{p^2}$	# of failures before first success
Uniform (continuous)	$\frac{1}{b-a}, x \in [a, b]$	$\frac{a+b}{2}$	$\frac{(a-b)^2}{12}$	random # generator
Normal	$\frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}$	$\mu$	$\sigma^2$	
Poisson	$\frac{\lambda^x}{x!} e^{-\lambda}$	$\lambda$	$\lambda$	# of phone calls in one day *related to binomial
Exponential	$\lambda e^{-\lambda x}$	$\frac{1}{\lambda}$	$(\frac{1}{\lambda})^2$	how long it takes between calls *related to poisson

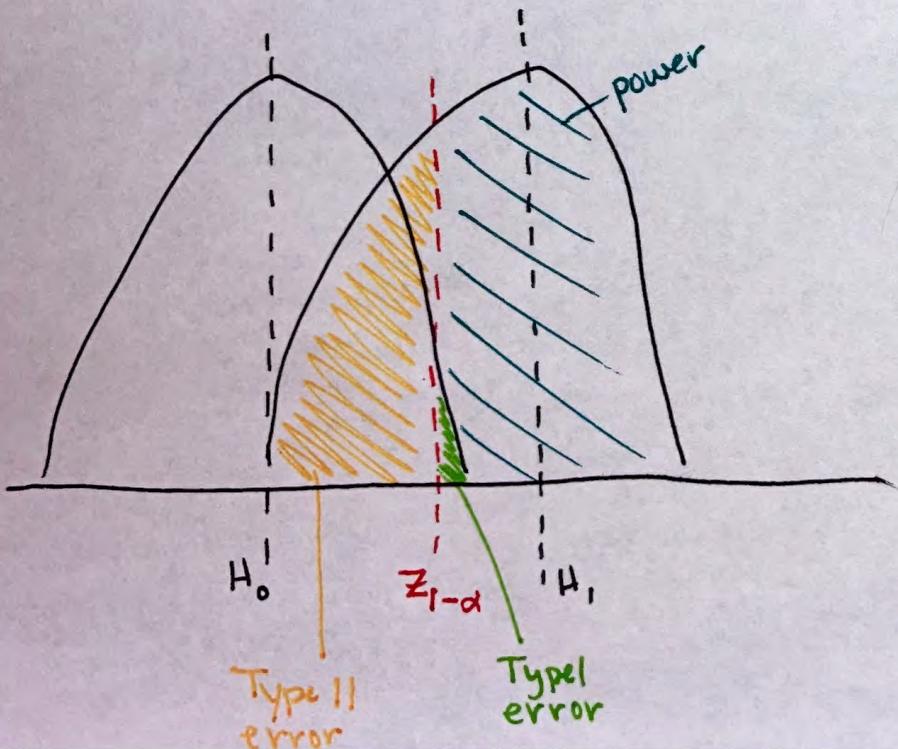
# Impacts of Sample Size

SE = standard error (dispersion of sample mean distribution)

power =  $1 - \beta$  (probability of not getting Type II error)

confidence level =  $1 - \alpha$  (probability of not getting Type I error)

effect = how big a difference is practically significant (minimum detectable effect, mde)



- as sample size  $\uparrow$ , SE  $\downarrow$  because sample mean gets closer to population mean (less dispersion)
- as sample size  $\uparrow$ , power  $\uparrow$  bc the range of acceptance decreases (smaller "area" of TII errors)
- as sample size  $\uparrow$ , confidence  $\uparrow$  bc the area of the distribution included within the acceptable range increases
- as sample size  $\uparrow$ , mde  $\downarrow$

# Sample Size

111

111

Def: total # of set of participants that adequately represent the population

## Assumptions:

- every individual has equal chance of being put in sample
- participants don't affect each other

$$n = \frac{2(z\text{-statistic}_{\alpha} + z\text{-statistic}_{1-\beta})^2 \cdot s^2}{\text{estimated effect size}}$$

$n$       =  
required  
sample  
size

# Types of Hypothesis Tests

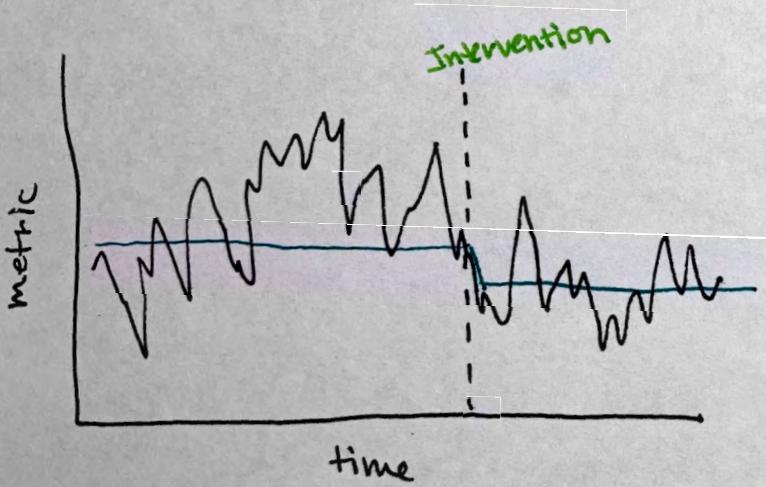
Name	When to Use	Test Statistic	Underlying Distro
t-test	<ul style="list-style-type: none"> <li>unknown mean + variance</li> <li>small samples</li> <li>finding means, difference</li> </ul>	$t = \frac{\hat{\theta} - \theta_0}{SE(\hat{\theta})}$	t-distribution <small>*at infinity this becomes Normal</small>
Z-test	<ul style="list-style-type: none"> <li>unknown mean, known variance</li> <li>large samples</li> <li>finding means, proportions, differences</li> </ul>	$z = \frac{\hat{\theta} - \theta_0}{SE(\hat{\theta})}$	normal
chi-squared	<ul style="list-style-type: none"> <li>tests independence of a pair of variables</li> <li>tests if there's a relationship (no direction)</li> <li>categorical data (contingency table)</li> </ul>	$\chi^2 = \sum \frac{(obs. value_i - expected value_i)^2}{expected value_i}$	chi-squared
Fisher's exact test	<ul style="list-style-type: none"> <li>tests independence</li> <li>alt. to chi-squared when n is small</li> <li>categorical data (contingency table)</li> </ul>	$P = \frac{(a+b)!(c+d)!(a+c)!(b+d)!}{a!b!c!d!n!}$	hypergeometric
Mann-Whitney U-test	<ul style="list-style-type: none"> <li>nonparametric</li> <li>lower power</li> <li>data from both groups need to be symmetric &amp; identical</li> <li>specific Null hypoth.</li> </ul>		
ANOVA (F-test)	<ul style="list-style-type: none"> <li>tests mean between 3+ groups</li> <li>tests explained variance vs. unexplained in model</li> </ul>	$F = \frac{\text{between group variability}}{\text{within-group variability}}$	F-distribution
Welch's t-test			

# Gathering Business Requirements

- What customer problem are we trying to solve?
- How will we know we've completed what's required?
- What does success look like?
- How will we measure success?
  - What specific metrics indicate success vs. failure?
  - How will those metrics be measured?
  - Who is responsible for taking those measurements?
  - When will the measurement take place?
  - Who needs to review the measurement + make the call?
- What decision will be made based on the outcome?
  - Who needs to approve that decision?
  - What is the business question or purpose?

# Intervention Analysis

Def: a time series analysis that estimates the effect of an external intervention on a time series, assuming the same ARIMA structure holds before + after intervention



The effects of the intervention are evaluated by changes in the level + slope of the time series + statistical significance of the metric

# Time Series

##

##

Def:

analysis of data that changes over time ; using a variable's own past values to forecast its trend

## Types of Models

Auto regressive (AR) = current data is built upon from data at previous time ; how far back do we need to look back? ; e.g. stock prices

Moving Average (MA) = current data is a result of previous unexpected events ; how many days of unexpected events do we look at?

Auto regressive Moving Average (ARMA) = combination of AR + MA models

## Assumptions:

- constant mean + constant variance for whole series → Stationary
  - use Lag Differencing to transform non-stationary data
- errors can be written as a linear function of past observations → Invertible