

1. Introductory Classical Statistics: t-test, Wilcoxon, Correlation

$$\text{test statistic} = \frac{\text{Observed data - what we expect if } H_0 \text{ is true}}{\text{Avg Variation}}$$

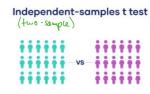
t-test - a test used to compare means of 2 groups to see if treatment group has stat sig difference

- Assumes: 1. Independent events 2. Normal Dist. 3. Similar Variance amongst both groups
- * If Assumptions fail, use WILCOXON Signed-Rank test

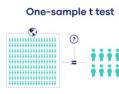
- Types



Investigate whether there's a difference within a group between two points in time (within-subjects).



Investigate whether there's a difference between two groups (between-subjects).



Investigate whether there's a difference between a group and a standard value or whether a subgroup belongs to a population.

One-tailed: $H_0: \mu_1 < \mu_2$ Pop₁ is less than Pop₂

Two-tailed: $H_0: \mu_1 \neq \mu_2$ Pops are diff

$$t = \frac{\bar{\mu}_1 - \bar{\mu}_2}{\sqrt{s^2(\frac{1}{n_1} + \frac{1}{n_2})}}$$

$\bar{\mu}_i$ - mean

s - Standard error of both groups

n_i - # of observations

Calculated difference in samples
[represented in units of stand. error]/
[relative to the data's variation]

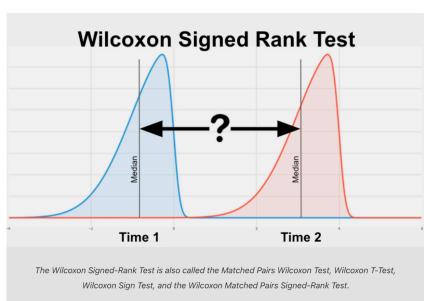
- Code
 - `t.test (variable ~ category var, data)`
 - `t.test (length ~ Species, flower.data)`

Wilcoxon

see above

- Compares Paired-Samples to determine if there's a stat. sig. diff.
- for nonparametric data: isn't quantitative but CAN be ranked (e.g. customer satisfaction)
- Assumes: 1. data is Paired-Sample 2. continuous data (since time) 3. Randomly sampled data

The Wilcoxon Signed-Rank Test is a statistical test used to determine if 2 measurements from a single group are significantly different from each other on your variable of interest. Your variable of interest should be continuous and your group randomly sampled to meet the assumptions of this test.



- types:
 - 1. Rank-sum: tests H_0 : Samples have same distribution i.e. no diff in samples

- 2. Signed-Rank: nonparametric version of t-test tests H_0 : Sample 1 is less/greater than Sample 2 Assumes skewed dist.

- Code: `wilcox.test (sample1.vector, sample2.vector, alternative = "greater")`

Correlation Tests

- Measures strength + direction of association b/w two vars
- Types: Pearson, Kendall, Spearman

$0 < r < 1$
weak \rightarrow strong

2. Resampling Methods: Randomization tests, Bootstrap, Cross-Validation

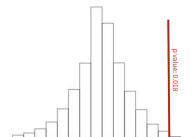
- the process of creating new samples based on 1 observed sample

Randomization tests

- aka Permutation tests

- nonparametric w/ very few assumptions \Rightarrow helpful when much is unknown about Pop.
- Randomly sample (without Replacement) OG data set into samples of different permutations, then perform test-stat to see if diff. exist (P-value)
- All samples are the same dim as OG

Permutation Test Distribution



Bootstrap

- Randomly select observations from OG data set (with Replacement)
 - St. your new sample dim = OG dim. Call this Sample 1 (S_1)
 - Find the M_i : **mean**(S_i) doesn't have to be mean, can be median or some other stat.
 - Repeat this for many S_i and M_i .
 - Add M_i 's to histogram to estimate what our data would look like if we repeated OG experiment/sample many times
- Basically it's a way to simulate what would happen if we repeated experiment 1,000's of times

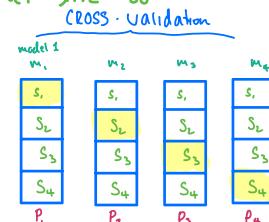
- Code:

```
library(boot)
meanfun <- function(vector, i) {
  mean(vector[i]) # any stat
```

`boot(data, meanfun, R=1000)`

Cross-Validation

- Automatically splits OG set into differently sized subsets for each train + test set like so:



CV chooses one fold (indicated with highlight) as test set (validation set) from train on the remaining folds. Performance P_2 is used then all P_i 's are compared and the model with best P_i is used.

- types: Validation-set, Leave one out CV, k-fold, Repeated CV

- Code:

```
library(caret)
train_control <- trainControl(method = "LOOCV")
model <- train(y ~ ., data,
               method = "lm",
               trControl = train_control)
```

k-fold

```
method: "cv"
number = 10
```

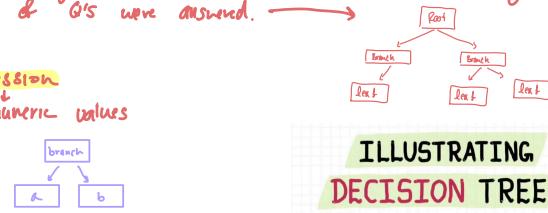
Repeated k-fold

```
method: "repeatedcv"
number = 10
repeats = 3
```

3. Decision Trees, Random Forest, Boosted TREES

Decision Trees

- nonparametric supervised learning method used for classification + regression based on how a previous set of Q's were answered.
- types - classification into categories + Regression predicts numeric values
- leaves are PURE if a or b is empty
- Quantify impurity to train model: Gini Impurity, Entropy, Information Gain
- Output of leaf is the category with most votes
- Code: Sklearn
- Main Disadvantage: tends to overfit training data

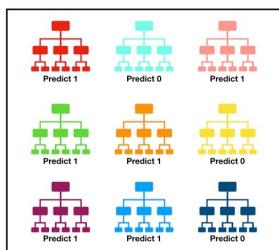


ILLUSTRATING DECISION TREES

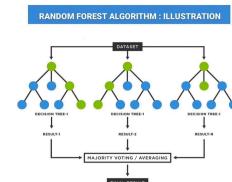


solution to this

RANDOM FOREST



- a set of many decision trees: each outputting a class prediction \Rightarrow the class that is predicted the most is our model's prediction
- called: Bagging (bootstrap aggregating)



- Fundamental idea: a large # of uncorrelated models (trees) operating as a committee will outperform any individual model
 \rightarrow the truth will out!

- bagging: trees are independent and are computed in parallel.
Thus, Forest is an amalgamation of a bunch of Ind tree models

- Main Drawback: if one tree has a discrepancy then whole forest has discrepancy. They are all interdependent "Parallel Learning"

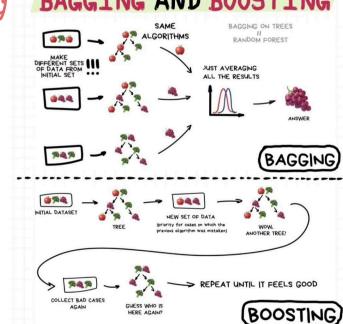
solution to this

BOOSTED TREES

- Boosting: trees are created sequentially one after another.
Each tree is created by previous tree, taking into errors of its predecessor into account

- Code: extreme gradient boosting (XGBoost)
Adaptive boosting (AdaBoost)

DIFFERENCE BETWEEN BAGGING AND BOOSTING



4. Generalized Additive Models (GAM)

- An alternative to linear models when data is non-linear that's easy to interpret and more flexible/complex than linear
Recall: Linear models

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n, \text{ where}$$

β_n - weights

x_n - variables

- GAMs Replace β_n 's with smooth non-linear functions called Splines (S_n)

$$y = S_0 x_0 + S_1 x_1 + \dots + S_n x_n, \text{ where } S_n = \sum_{k=1}^K \beta_k \cdot b_k(x)$$

β_k - weights

$b_k(x)$ - "basis expansion" x^0, x^1, x^2, \dots

Splines - hard to define + understand. Essentially a polynomial func. with small range

- The final model is just the sum of all the splines

- Wigginess - how wiggly our model's line of best fit is. The more splines the more wiggly (overfitting)
So there's strategy in deciding # of splines used

multiply S_k 's by λ penalty term then use CV to optimize

*Start w/ high # of splines then CV with λ to optimize

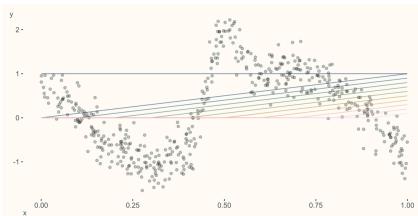
- Any Distribution goes: Normal, Poisson, Binomial

- Code: library(mgcv)

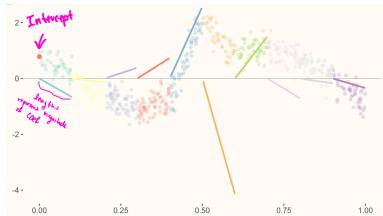
```
model1 <- gam(y ~ x_0 + s(x_1) + s(x_2), data)
```

```
model2 <- gam(y ~ x_0 + s(x_1), data)
```

Set Spline degree = 1 (linear)

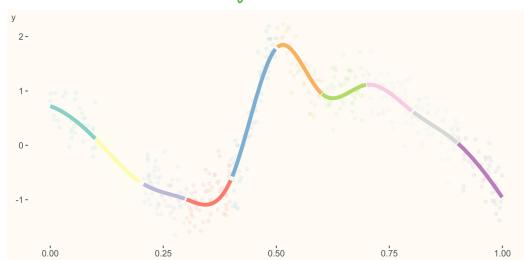


Multiply b_k 's by their corresponding Regression coeff

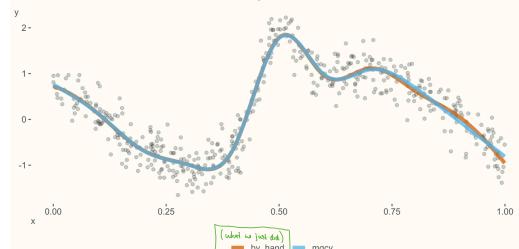


Sum of basis functions

Set Spline deg = 3 (cubic)

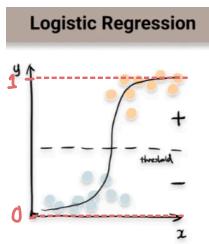


Compare to mgcv



5. Logistic Regression + Classification Models

Logistic Regression



- Sigmoidal Function:

maps any Real x to value b/w 0-1
 $x: \mathbb{R} \rightarrow (0, 1)$

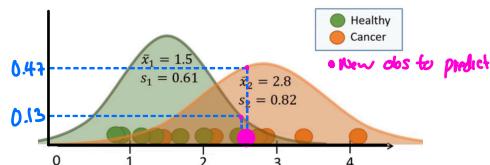
$$f(x) = \frac{1}{1 + e^{-x}}$$

- helpful bc $0 < y < 1$, y - predicted output

ex) X - emails
 y - prob that x is spam $\Rightarrow f(x) = \begin{cases} \text{Spam}, & \text{if } y \geq 0.5 \\ \text{Not Spam}, & \text{if } y < 0.5 \end{cases}$

Code: `logReg <- glm(y ~ ., data, family = "binomial")`

Naive Bayes



1. Find mean (\bar{x}) and std dev (s) of each class

2. Plug into Gaussian function to plot Normal Dist's

$$f(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

3. Predict new observations' classes based on height of each curve (height Reps. prob. that the obs. is in that class) See: Pink

height_cancer = 0.47 \Rightarrow more likely so this class height_healthy = 0.13

- Based on Bayes Thm:

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

A - outcome (classification)

B - Predictors: $B = x_1, x_2, \dots, x_n$

Posterior: $\frac{\text{Prior} \cdot \text{likelihood}}{\text{evidence}}$

"Naive" bc predictors are Independent or "naive" towards one another.
 $\text{key sum} = \text{sum key}$
 $\bar{x}_1 + \bar{x}_2 = \bar{x}_1 + \bar{x}_2$

"Gaussian" because this is a normal distribution
 $P(\text{class} | \text{data}) = \frac{p(\text{data} | \text{class}) \times p(\text{class})}{p(\text{data})}$
 This is our belief

We don't calculate this in naive bayes classifiers

- All predictors have equal effect on outcome. Continuous vars

$$P(y | (x_1, x_2, x_3, \dots, x_n)) = \frac{P(x_1|y)P(y) \cdot P(x_2|y)P(y) \cdots P(x_n|y)P(y)}{P(x_1) \cdot P(x_2) \cdots P(x_n)} = \prod_{k=1}^n \frac{P(x_k|y)P(y)}{P(x_k)}$$

Code: `library(caret)`

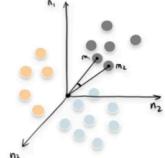
nb.fit \leftarrow naive Bayes(classVar ~ ., training set)

nb.pred \leftarrow predict(nb.fit, testset)

confusion matrix...

K-Nearest Neighbors

K Nearest Neighbour



- nonparametric supervised classifier that uses proximity of training data to make predictions about the grouping of an individual point

- K-lables describe the # of training observations within the neighborhood (for each class) ex:

4 Purple
 3 Blue
 2 Red
 K neighborhood

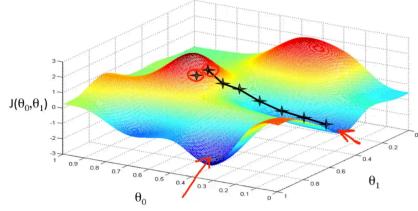
- Regression: returns mean of k-lables
- Classification: "mode"
- Pros: Simple, versatile
- Cons: costly as predictors inc

Code: `library(class)`

```
knn.fit <- knn(train = trainingSet
                 test = testSet
                 cl = classSet
                 k = neighborhood )
```

confusion matrix...

Stochastic Gradient Descent



Code: `sgd(y ~ ., data, model = "lm")`

Lin Reg is used here but can be other models

- Procedure to find minimum point of a cost function. Stochastic = Random
- Blindfolded on a MTN summit \rightarrow feel around with your feet to get down.
- Iterative Step algorithm. At each step, we choose a Random point from calc. Step direction
- Learning Rate - parameter that controls model's flexibility
as learning rate INC steps INC \Rightarrow too large and your model may jump over min. point

Step size: gradient \cdot learning rate
new params = old params - step size
calc. gradient with partial derivatives

Normal GD considers all observations at each step. Stochastic Randomly picks 1.

Normal GD Con: Costly \Rightarrow a observations, b predictors

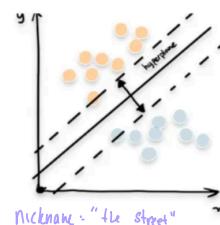
a^b Computations per step

1000 steps is typical $\Rightarrow 1000 \cdot a^b$

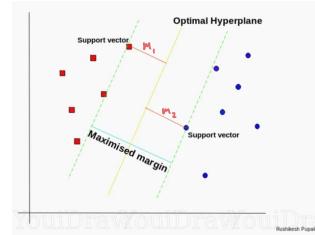
\star mini-batch GD is a solution to this OS it takes a subset of a + Stochastic GD since it chooses a randomly

Support Vector Machines

Support Vector Machine



Algorithm creates a line/Hyperplane to separate data into classes



Support Vector - a point within each class that is closer to estimated Hyperplane

Margin: $m_1 + m_2$

when margin is maximized optimal Hyperplane is found
above underr possible street

Aha moment: models are **ALWAYS** linearly separable when observed from the perspective of its highest dimension

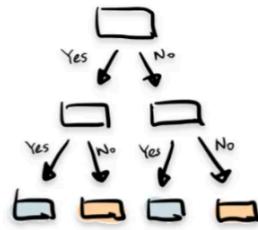
- Mapping to higher dimensions is the toughest part of this AND computationally expensive.
solution: Kernel trick \Rightarrow uses vector dot products instead of mapping to HDims
- Pro: handles higher dims well, versatile
- Con: sensitive to kernels so if np.sum performs poorly,
non-Probabilistic (not for Regression) - but - can measure effectiveness of classification by looking at distance from new point to hyperplane
- Soft vs. hard margin classification: allowing vs. not-allowing data within margin
 \star if soft then use cost function to penalize (derentivize) data from being inside

Code: `library(e1071)`

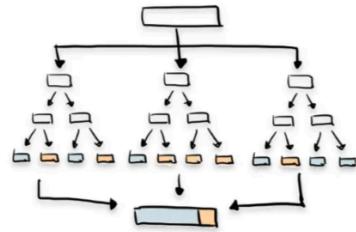
```
svm.fit(y ~ ., data, kernel = "linear", cost = 10, scale = False)
      polynomial      Radial      Sigmoid
      penalized       data within    margin
      scale data      but fitting?
```

Decision Tree + Random Forests

Decision Tree

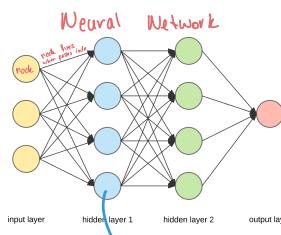


Random Forest



6.

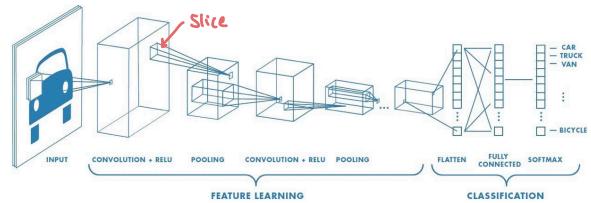
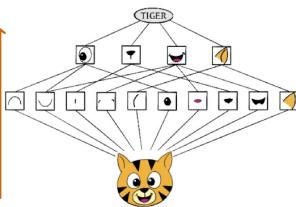
Neural Networks



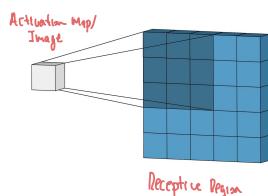
- If network > 3 layers including input/output then considered Deep learning Network
- The more layers the greater ability to recognize complexity
- Each node connection has weight representing how helpful node is in correctly classifying info
- Each layer has a weight threshold that must be satisfied for info to be passed on to next layer
- Initially weights + thresholds are randomly set ✗
↳ then with Forward + Backward Propagation they're updated + optimized
- Corrective Feedback Loop: gives more weight to nodes that guess correctly, less to mistaken nodes

Each node has an Activation Function to Compute weight. Pop functions: Soft Plus, Rectified Linear, Sigmoid
When initializing a NN we must decide: # of hidden layers, # of nodes within each layer, Activ. func.

Convolutional NN
Gives computers Sight

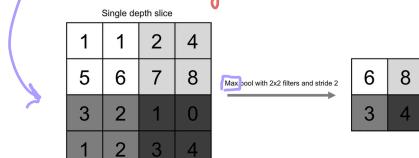


- Best for Image Analysis or any data with grid-like topology
- Picks out patterns + makes sense of them. Layers are arranged s.t. easier patterns (lines, curves, etc) are picked out 1st then more complex (faces, objects)
- 3 Layers
 - ① Convolution Layer: Main computational layer
 - performs dot prod of 2 matrices
 - ↳ One is set of learnable parameters (kernels)
 - Other is the $n \times n$ pixel set from input
 - Response Map: matrix of dot products



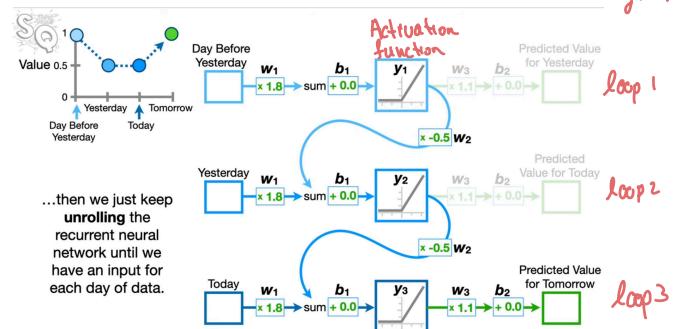
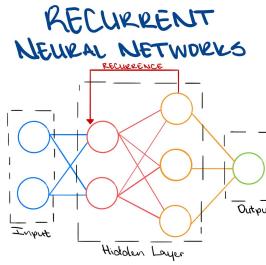
- kernel Slides across height + width of Receptive Region, computing the dot prod each time and storing it in a smaller matrix image called the Activation Map
- Reduces memory ↳ "sparse interaction"

- ② Pooling Layer: Computes + Stores a Summary Statistic for each Activation Map
Ex: mean, L2 norm, maximum of elements within $n \times n$ map
"Pooling function/operation"
- Another layer to Reduce dimensions of input



- ③ Fully Connected Layer: brings everything back together
All neurons are connected/mapped to all those immediately preceding + succeeding

Recurrent NN



Ex: Stocks for 2 companies that started at diff dates

- Works well if inputs have diff dimensions
- Uses Feedback loops
- Sequential data ex. time
- Each loop is for each previous time stamp / sequence in input data
- Not used often bc as loops inc model gets hard to train
"The Vanishing/Exploding Gradient Problem"

* costly!!

Long Short-term memory

7. UnSupervised Learning: Clustering + Principal Component Analysis

PCA

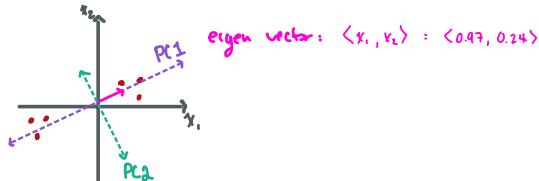
Clustering method to Represent multi-dimensional data ($p > 3$) in 2 dimensions
AND can determine which variable is most valuable

Procedure:

1. Center data about origin + Project higher dim data on 2D plane
2. Fit line or best fit by maximizing sum of squares of data's projections

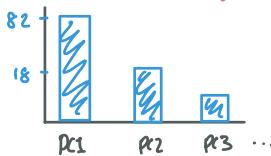
PC1 is linear combo of this line: $\vec{u} = 0.97x_1 + 0.24x_2$

↳ Unit vector is eigenvector
↳ mean (sum of sq.) is eigenvalue for PC1
↳ sum of sq. is single value for PC1



3. PC2 is line \perp PC1
4. Rotate Everything st. PC1 lies horizontally
5. Use projected training obs. to compute eigenvalues for each PC which reps the % of variation each PC accounts for
 - Eigenvalues are measures of variation
6. Repeat 1-5 until you go thru all vars / dimensions
 - * Thus, # of PC's = # of vars OR # of obs (whichever is smaller)

6.5 Scree Plots: bar graph representing % of variation caused by PC1 + PC2



• PC1 var accounted for $>$ PC2 $>$ PC3 $>$...

Code: `pca ← prcomp(data.matrix, scale=TRUE)`

↳ Returns 3 things:
PC's Sdev Rotation
Standard dev Loading Scores

Clustering

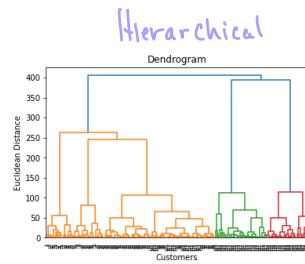
k-means

Procedure:

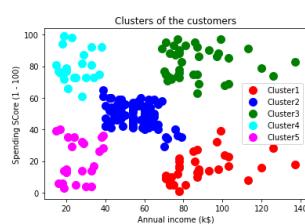
1. Choose $\#$ of clusters /centers you want $\xrightarrow{\text{CV helps optimize this}}$
2. Randomly select C obs. to establish initial clusters
3. Assign all other obs. to their closest C cluster
4. Calc mean of each C
5. Assess quality of clustering by adding up Var within each C : "Total var"
6. Keep track of total vars then Repeat 1-5 a few more times
7. Compute total var values from 1-6 iterations and choose best one

Code: `kmeans(X = data.vector, centers = 3)`

Hierarchical



k-means



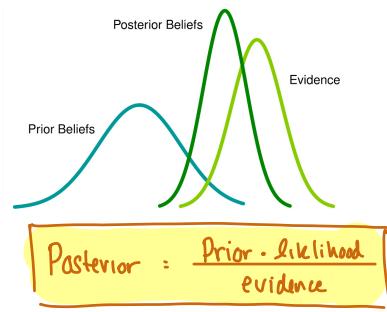
	K-means	Hierarchical
Time Complexity	$O(n*k*t)$	$O(n^3)$
Space Complexity	$O(n(d + k))$	$O(n^2)$
Hyperparameters Tuning	Must specify the number of clusters (k) and retrain model for each k	Can dynamically update k value without retraining model
Data Structure	Better performance when dealing with convex clusters	Generate better result when dealing with non-convex clusters
Variations	Many variations (e.g., K-median, K-medoid...) with different distance matrices	Two approaches: Agglomerative approach and Divisive approach
Optimization	K-means++ introduces smarter initialization of centroids, making convergence faster	Top-down approach reduces time complexity to $O(n^2)$
Result Robustness	Result may be different on different runs	Same parameters generate the same result every time

8.

Bayesian Analysis

Bayesian Statistics

- A theory in the field of statistics based on Bayesian Interpretation of Probability where prob expresses a degree of belief in an event based on prior knowledge about the event (i.e. results of prev. experiment)
- This is different from other Interpretations of Probability such as The Frequentist Interpretation that view prob as the limit of the relative frequency of an event after many trials.
- All parameters assumed to be Random Quantities
- A param is summarized by an entire Dist. of values instead of just one value as is in Frequentist Analysis
- Posterior dist. is the goal / "heart" of Bayesia



Gaussian Naive Bayes

- for continuous data + data w/ Normal Dst.
- Good at handling continuous values

$$f(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

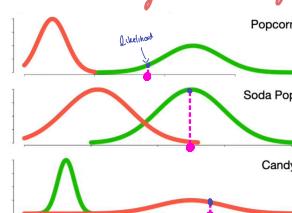
Multinomial Naive Bayes

- text classification: cares about # of occurrences of the word
- Good at handling Discrete values. Can classify non-numeric data
- Main advantage: significantly Reduces Complexity
↳ classification using small training sets without needing to be retrained



	Popcorn	Soda Pop	Candy
Popcorn (grams)	24.3	750.7	0.2
Soda Pop (ml)	28.2	533.2	50.5
Candy (grams)	etc.	etc.	etc.

	Popcorn	Soda Pop	Candy
Popcorn (grams)	2.1	120.5	90.7
Soda Pop (ml)	4.8	110.9	102.3
Candy (grams)	etc.	etc.	etc.



"new obs."

$$P(y|x_1, x_2, x_3) = \frac{P(x_1|y)P(x_2|y)P(x_3|y)}{P(x_1)P(x_2)P(x_3)}$$

$$y = \text{Loves Troll 2}$$

$$x_1 = \text{Popcorn} \quad x_2 = \text{Soda Pop} \quad x_3 = \text{Candy}$$

Gaussian Naive Bayes is named after the Gaussian distributions that represent the data in the Training Dataset.

Bernoulli Naive Bayes

- text classification: cares whether the word happened or not
- Good at handling binary/boolean values

9. Bayesian Ext. Markov Chain Monte Carlo (mcmc)

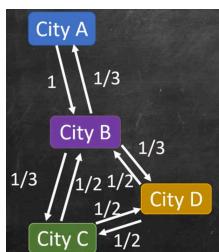
- A method to approximate Posterior Dist. of a parameter if not able to do so Analytically
- Applies Law of Large #s (Central Limit Theorem) to Dependent Results

Markov Chain a sequence of events where the prob of the future only depend on the present

- Markov Property**: given the present, the future is Independent of the past Memoryless-ness
- The posterior probability relies solely on present probability and Randomness of data generation

The next Step is a Random Choice based on Current Step

- Non-Markov has memory

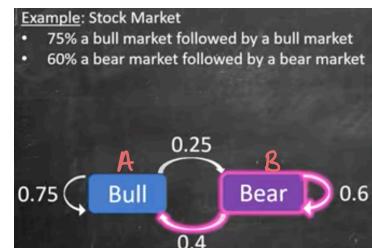


- Traveling example
- Transition Diagram**: same idea as tree diagram
- Transition Matrix**: $P = \begin{bmatrix} A & B \\ C & D \end{bmatrix} = \begin{bmatrix} 0.75 & 0.4 \\ 0.25 & 0.6 \end{bmatrix}$

$$S_n = P^n S_0 \quad \text{where } D = \text{transition matrix}$$

$S_0 = \text{initial state}$

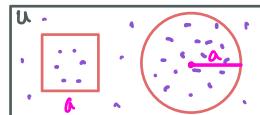
$S_n = n^{\text{th}} \text{ state}$



- A representation of how a var walks around a graph Nicknamed "Random Walk"
- How a Random var changes from one state to the next

Monte Carlo

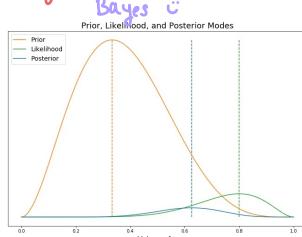
- Simulations evolving Randomly
- Fundamental idea that allows us to say: Random Samples can represent Populations



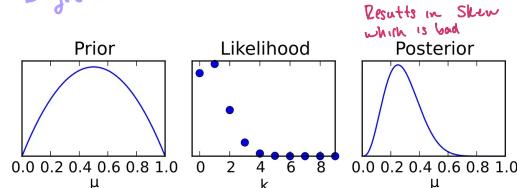
Randomly dropping marbles in U . Proportions of marbles in Circle vs. Square Match area Ratio: $\frac{\pi a^2}{a^2} = \pi$

MCMC

- Allows us to leverage computers to do Bayesian Stats
- Bayesian Stats requires prior Dist. + likelihood Dist to be Normal



Bayes is so enter mcmc



MCMC is used when prior + likelihood aren't Normally Dist.

Procedure

1. Monte Carlo : Generate Random Sample

ex: generate a series θ_t of random $t \in$ from normal dist. $\theta_t \sim N(\mu, \sigma)$

2. Markov Chain : Determine Transition diagram/matrix

$\theta_t \sim N(\theta_{t-1}, \sigma)$ so plugging in the prior term θ_{t-1} in as the new term's ^{Dist.} mean

3. Acceptance - Rejection Sampling : Decide to keep/Discard new obs. generated from 1+2

• Algorithms

1. Metropolis-Hastings

- We are at point x .
- We make a guess for the next step. We will call this x^*
- We then compute the ratio of the probability of x^*/x . This is calculated using the product of the likelihood and prior distributions.
- If the ratio of $p(x^*)/p(x)$ (also called the acceptance probability) is greater than 1 we accept x^* as the new position.
- Even if the acceptance probability is less than 1, we don't automatically reject x^* . We flip a coin by selecting a random number, from a Uniform(0,1) distribution. If the number is smaller than the acceptance probability we accept x^* if it is higher we reject x^* and start the process over again.

2. No U-turn

Summary Procedure

- We randomly generate numbers: This is the Monte Carlo part
- We allow the numbers we generated to influence the next generated number: This is the Markov chain
- We then decide if the new numbers generated are "moving in the right direction": The Acceptance-rejection algorithm
- We then check for convergence: We determine when our data has converged to a reasonable distribution. The randomly generated values after the convergence point become our posterior distribution