

CS286 Solving Big Data Problems – Exam #2 Study Guide

By: Zayd Hammoudeh

Lecture #14 – Introduction to Data Science

Features

| | | |
|--|---|--|
| <p>Model = SampleData + Algorithm(s)</p> <p>Feature – Individual, measurable property of an observation. In supervised learning, it is used for predicting the target.</p> <p>Additional Names:</p> <ul style="list-style-type: none"> • Predictor • Attribute • Variable <p>Vectorizing Data – Given an observation, assign features to dimensions.</p> | <p>Attribute Types</p> <ul style="list-style-type: none"> • Continuous (Numeric) – Data across a continuum. • Binary – Two possible values. • Categorical – Finite set of possible values. • Text – Vocabulary or n-gram. <ul style="list-style-type: none"> ◦ Bi-gram – 2 consecutive words. | <p>Qualities of Good Features/Feature Sets</p> <ul style="list-style-type: none"> • Correlation with the target (i.e. class value) • Features are independent • Number of features is minimized • Number of features is much less than the number of data points. • Features do not leak the target meaning they do not provide information that would not be available under normal circumstances. |
|--|---|--|

| Principal Component Analysis Procedure – Dimensionality Reduction | | |
|--|--|--|
| <p>Set Dimensionality $\{m, n\}$</p> <ul style="list-style-type: none"> • m – Number of observations • n – Number of Features <p>Ideally: $m \gg n$</p> | <p>Procedure</p> <ol style="list-style-type: none"> 1. Calculate covariance matrix of the original input matrix. 2. Perform eigenvalue decomposition of the covariance matrix. 3. Look for orthogonal dimensions with the greatest variance (i.e. principal components) 4. Project input matrix into those dimensions. <p>Relative covariance preserved across projection.</p> | <p>λ – Representation of the total variance of the principal components.</p> $\lambda = [4.2, 0.24, 0.08, 0.02] = \frac{4.2 + 0.24}{4.2 + 0.24 + 0.08 + 0.02} = 98\%$ <p>Number of principal components and eigenvectors equivalent to the number of dimensions in the original input matrix.</p> <p>Since the data is mapped to new dimensions, not easy to map principal components to the original data.</p> |

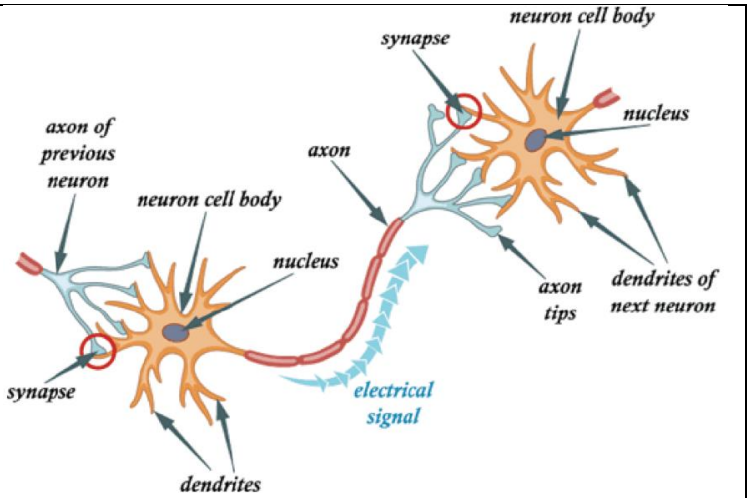
| Feature Scale | | |
|--|--|--|
| <p>Potential Issues</p> <ul style="list-style-type: none"> • Relative Magnitude – Some features may have different magnitude than others causing those features to dominate. • Absolute Magnitude – When features are too large or too small, data may be lost to overflow or underflow respectively. | <p>Potential Solutions</p> <ul style="list-style-type: none"> • Standardization – Transform data to measure how many standard deviations it is from the mean. $x' = \frac{x - \bar{x}}{\sigma_x}$ • Normalization – Transform data such that the range is between 0 and 1. $x' = \frac{x - x_{min}}{x_{max} - x_{min}}$ <p>Both methods preserve covariance.</p> | <p>Target Leaks – Sample data has “extra” information/features not found outside the sample. May cause the model to be unrealistically good.</p> <p>Bias</p> <ul style="list-style-type: none"> • Sampling Bias – Non-representative sample • Observer Bias – Corrupted/invalid measurement of features • Funding Bias – “Fudging”/manipulating of the data. <p>Simpson’s Paradox – Trend that appears in different groups disappears (or changes) when groups are combined. Example: The average of averages may not equal the average all points.</p> |

Algorithms and Evaluation

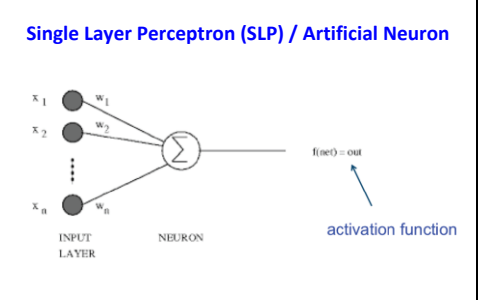
| | | | |
|--|--|---|--|
| <p>Applications of Clustering</p> <ul style="list-style-type: none"> • Image Segmentation • Dendrogram (including Phylogentic Dendrogram) • Market Segmentation | <p>K-Means Clustering Algorithm</p> <ol style="list-style-type: none"> 1. Initialize the K centroid 2. Assign points to the nearest centroid. 3. Update the centroids. 4. Check Stop Conditions | <p>Methods for Selecting the Initial Clusters</p> <ol style="list-style-type: none"> 1. Random 2. El Agha et. al. Method 3. All at the same point (e.g. origin) <p>Stop Condition Options</p> <ol style="list-style-type: none"> 1. Min change of mean from last iteration. 2. Sufficiently small intra-cluster distance 3. Sufficiently large intercluster distance | <p>Intracluster Distance – Distance between points within the same cluster. Measures the cluster’s similarity.</p> <p>Intercluster Distance – Distance between the centroids. Measures data separation between clusters.</p> |
|--|--|---|--|

| | | |
|---|--|--|
| <p>Naïve Bayes</p> $P(C F_1, \dots, F_n) = \frac{P(C) \cdot P(F_1 C) \cdot \dots \cdot P(F_n C)}{P(F_1, \dots, F_n)}$ <ul style="list-style-type: none"> • Assumes all features are conditionally independent. • C – A class value • F_i – The i^{th} feature. • The predicted class value is the one with the highest Naïve Bayes probability. | <p>Sensitivity/True Positive Rate</p> $TPR = \frac{TP}{TP + FN}$ <p>Specificity/True Negative Rate</p> $TNR = \frac{TN}{TN + FP}$ <p>Accuracy</p> $Acc = \frac{TP + TN}{TP + FN + FP + TN}$ | <ul style="list-style-type: none"> • For some domains (e.g. SPAM and cancer detection), you may use a weight function as different types of misclassification may have different costs. |
|---|--|--|

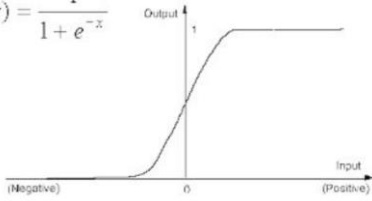
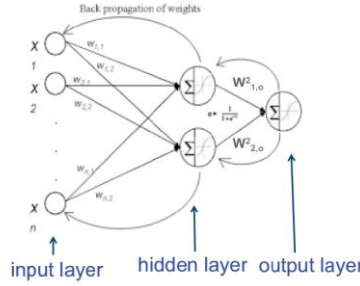
| | | | |
|--|--|--|---|
| <p>K-Nearest Neighbors (KNN)</p> <p>Example Distance Metrics</p> <ul style="list-style-type: none"> • Euclidean • Manhattan • Cosine <p>Voting Related Issues:</p> <ul style="list-style-type: none"> • May be problematic if class distribution is skewed. • Can use distance as a weight. • May want to pick an odd number for k to avoid ties. | <p>Types Recommendation Filtering</p> <ul style="list-style-type: none"> ◦ Item Based ◦ User Based <p>Precision – Proportion of top-search results that are relevant.</p> $Precision = \frac{TP}{TP + FP}$ <p>Recall – Proportion of relevant results that are top-scoring.</p> $Recall = \frac{TP}{TP + FN}$ | <p>Underfitting – Model is too simple for the data set.</p> <p>Overfitting – Model is overly trained to match the training set.</p> <p>Occam's Razor – Most plausible explanation is the simplest one.</p> <p>Pessimistic Estimation – Impose a penalty (λ) in cost function based on model.</p> <ul style="list-style-type: none"> • Smaller λ – Overfitting • Larger λ – Underfitting | <p>Sunrise Problem</p> <p>$P(\text{sun will rise tomorrow}) = ?$</p> <p>Using existing observations the probability is one but there will be a day when it does not rise.</p> <p>Smoothing: Use probability distributions for parameter estimates for sparse data.</p> $\hat{\theta}_i = \frac{x_i + \alpha}{N + \alpha \cdot d}$ <ul style="list-style-type: none"> • Smaller α – Overfitting • Larger α – Underfitting |
|--|--|--|---|

| | | |
|--|---|--|
| <p>Cross Validation</p> <ul style="list-style-type: none"> • Leave k out. • Train on the $N - k$ points. Use the k points to estimate the out of sample error. • Repeat this process until all points have appeared in the training set exactly once. • Example: Leave-1-Out. Generate and train N models. Using this approach, the error is: $Error_{LeaveOneOut} = \frac{1}{N} \cdot \sum_{i=1}^N e_i$ <ul style="list-style-type: none"> • Cross validation can be repeated with different model parameters. Select the model with the lowest cross validation error. | <p>K-Fold Cross Validation – Divide the training set into k equally sized pieces. Example: $k = 10$ means each testing set is size $\frac{N}{10}$.</p> <p>Downside of Leave-One-Out – Computationally expensive.</p> |  |
|--|---|--|

Lecture #15 – Introduction to Machine Learning

| | | | |
|--|---|--|---|
| <p>Machine Learning – A subfield of computer science that deals with the construction and study of systems that can learn from data rather than follow explicitly programmed instructions</p> | <p>Neuron – A brain cell.</p> <p>Dendrite – Appendages that stick out from a neuron.</p> <p>Axon – Long chain that connects two neurons' dendrites. Carries the electrical signal between two cells.</p> <p>Synapses – Small cleft or separation between two neurons.</p> | <ul style="list-style-type: none"> • When neuron A repeatedly participates in firing neuron B, the strength of the action from A to B increases. Cells that fire together wire together. <ul style="list-style-type: none"> ◦ If they fire separately, the weight decreases. • Synaptic Strength – Can be modeled as a set of weights. • Hebbian Learning – Changing of the weights between the neurons in the brain. | <p>Single Layer Perceptron (SLP) / Artificial Neuron</p>  |
|--|---|--|---|

| | | |
|---|---|---|
| <p>Activation Function – Defines the relation between the summed input to the perceptron and the perceptron's output.</p> <p>Perceptron's can be taught to recognize a binary pattern (e.g. 0 or 1)</p> | <p>Rosenblatt Perceptron Learning Algorithm</p> <p>Step #1: Initialize weights and perceptron thresholds (e.g. randomly).</p> <p>Step #2: For each input (\vec{x}) and expected/desire output (d), do the following:</p> <ul style="list-style-type: none"> ◦ Calculate the net signal via: $\vec{x}^T \vec{w}$ ◦ Calculate the perceptron's output: $y_j(t) = \text{sign}(\vec{x}^T \vec{w})$ ◦ Compute the error (one of the set $\{-1, 0, 1\}$): $e_j(t) = d_j(t) - y_j(t)$ ◦ Update the weights $w_i(t+1) = w_i(t) + \alpha \cdot e_j(t) \cdot x_{ij}$ <p>Step #3: If stop condition met, then terminate, else repeat step #2.</p> | <p>Sigmoid Activation Function</p> $f(x) = \frac{1}{1 + e^{-x}}$ <p>Benefits of the Sigmoid Function:</p> <ul style="list-style-type: none"> • Good approximation of the step function • Permits a continuous output • Solves noise saturation for large signals. • Solves noise attenuation for small signals. • Has a simple first order derivative. |
|---|---|---|

| | | | |
|---|---|--|--|
| <p>Sigmoid Function</p> $f(x) = \frac{1}{1 + e^{-x}}$  <p>Output</p> <p>Input (Negative) 0 (Positive)</p> | <p>Deficiency of the Perceptron – Cannot solve all functions. Example: XOR.</p> <ul style="list-style-type: none"> Discovered by Minsky and Papert. <p>Perceptron can only solve functions that are linearly separable.</p> | <p>Multilayer Perceptron</p>  <p>Back propagation of weights</p> <p>input layer hidden layer output layer</p> | <p>Multilayer Perceptron: Developed by Werbos</p> <p>Two Stages: Repeat until convergence</p> <ul style="list-style-type: none"> Forward Pass – Use the inputs to calculate the outputs and deltas for each neuron. Reverse Pass – Adjust the weights by fractions (α) of the deltas. Done via back propagation layer to layer. <p>Multilayer perceptrons are trained to recognize patterns. Example: Handwriting recognition.</p> <p>Can compute XOR.</p> |
|---|---|--|--|

| Algorithm Name | Year Published | Author |
|------------------------|----------------|--------------------|
| Decision Tree | 1986 | Quinlan |
| Support Vector Machine | 1995 | Vapnik and Cortes |
| Boosting | 1998 | Freund and Shapire |
| Random Forest | 2001 | Breiman |
| Deep Learning | 2005 | Hinton et. al |

machine learning

supervised

- Classification**
 - KNN
 - Decision tree
- Regression**
 - linear
 - logistic

Lots of other algorithms that sit in the middle here somewhere

unsupervised

- Clustering**
 - K-means
 - DBSCAN
- Dimensionality reduction**
 - PCA
 - SVD

| | | |
|--|---|--|
| <p>Single Layer Perceptrons can learn linearly separable functions only.</p> <p>Multilayer perceptrons can learn linearly and non-linearly separable functions.</p> | <p>Supervised learning has input-output pairs during the training phase (i.e. labeled data).</p> <p>Unsupervised learning uses only inputs, not outputs (i.e. unlabeled data).</p> | <p>Classification Problems – Have a finite set of outputs.</p> <p>Regression Problems – Do not have a finite set of outputs.</p> |
|--|---|--|

Lecture #17 – Using Mahout for Naïve Bayes

| | | | |
|---|---|---|--|
| <p>Classification – A supervised machine learning algorithm. Distinguishes objects of one class from another. Emulates human decision making.</p> <p>Example of Classification:</p> <ul style="list-style-type: none">• Credit card fraud detection• Credit approval <p>Example of Regression:</p> <ul style="list-style-type: none">• Credit card profitability. | <p>Feature Extractor – Given an input object, builds/extracts the set of features.</p> <p>Machine Learning Algorithm – Generates the classifier model.</p> <p>Classifier Model – Given an input, predicts a class label/value.</p> | <p>Feature / Predictor Attribute Types</p> <ul style="list-style-type: none">• Continuous – Decimal or floating point.<ul style="list-style-type: none">◦ Example: 0.1• Categorical – Predefined, finite set of values.<ul style="list-style-type: none">◦ Example: { true, false }• Word-like – Large set of defined values.<ul style="list-style-type: none">◦ Example: English dictionary• Text-Like – Sequence of word-like objects.<ul style="list-style-type: none">◦ Example: Email message subject. | |
| <p>Bayes Theorem</p> $P(A B) = \frac{P(A) \cdot P(B A)}{P(B)}$ <ul style="list-style-type: none">• Can be used to prove the Monty Hall problem. | <p>Assumptions of Naïve Bayes</p> <ol style="list-style-type: none">1. Input data is labeled.2. Predictors are in <i>N</i>-dimensional space.3. Target is a set of categorical values4. Class C is dependent on the set of input attributes.5. All features are conditionally independent. | <p>Naive Bayes may not be Bayesian since it simplifies the problem by the independence assumption.</p> <p>Bag of Words – Transform a text input into an unordered set of words.</p> | <p>Calculating Conditional Probability</p> $P(C A) = \frac{N_{A,C} + 1}{N_C + k}$ <ul style="list-style-type: none">• <i>N_C</i> – Number of elements of class <i>C</i>• <i>N_{A,C}</i> – Number of elements of class <i>C</i> with attribute <i>A</i>• <i>k</i> – Number of classes. |
| <p>Sequence of Mahout Commands</p> <ul style="list-style-type: none">• mahout seqdirectory• mahout seq2sparse• mahout split• mahout trainnb• mahout testnb | | | |

Lecture #13a – HBase

| | | | |
|--|--|--|--|
| <p>Relational Database – A datastore that whose type and structure is defined before storage. Uses SQL.</p> <p>Deficiencies of Relational Databases</p> <ul style="list-style-type: none">Do not scale horizontally.Sharding is difficult to manage (join and transactions do not scale across shards) | <p>HBase – A distributed database where puts and gets are accessed via a key.</p> <p>Table Splits – Occur automatically as the table grows.</p> <p>Horizontal Partitioning – Putting rows into different tables. HBase uses key ranges/regions to define the horizontal partitions (shards).</p> | <p>Row Key – Used to store and access data.</p> <p>Column Family – Used to group and store similar data (i.e. subcolumns – column qualifiers). Attributes of column families include:</p> <ul style="list-style-type: none">Number of versionsTime to Live (TTL)CompressionKeep in memory or preserve to disk. <p>Column family are separated into different files</p> | |
| <p>Pros of HBase:</p> <ul style="list-style-type: none">Scalable to handle data volume and velocity.Fast reads and writes by key. <p>Cons of HBase:</p> <ul style="list-style-type: none">Does not support joinsGood schema design is needed to achieve the best performance. | <p>Data is stored in a key-value model.</p> <p>Information Required to Access a Single Cell Value in HBase</p> <ul style="list-style-type: none">Row KeyColumn FamilyColumn Qualifier (i.e. subcolumn name)Timestamp/Version | <p>Version in HBase</p> <ul style="list-style-type: none">Each put and delete adds a new version/cell.Stores last 3 versions by default.Version – Stored as a long which is the current time in milliseconds (if no specific version is specified) | <p>Table Physical View</p> <ul style="list-style-type: none">Stored as a sorted map.Ordered by row key and column qualifier in ascending order.Ordered by version/timestamp in descending order. |
| <p>Basic Table Operations</p> <ul style="list-style-type: none">Create table and defining of column families is done before data is imported. <p>Basic CRUD Operations</p> <ul style="list-style-type: none">Put – Insert data into rows (both create and update)Get – Access data from one rowSCAN – Access data from a range of rowsDelete – Delete a row or a range of rows/columns. | <p>Region – Another name for a key range. All key range is contiguous.</p> <p>Region Server – Serve data for reads from or writes to a particular region/key range.</p> <p>Write Ahead Log (WAL) – Disk commit log used for recovery.</p> <ul style="list-style-type: none">Primary Role: Durability (log on disk)Updates appended sequentially <p>Block Cache – Read cache. Uses the Least Recently Used (LRU) paradigm for block eviction.</p> | <p>Memstore – Write cache.</p> <ul style="list-style-type: none">In memory.Sorted set of key-value pairs. Updates quickly stored since in memory.One memstore for each column family. <p>HFile – Sorted key-value pairs on disk.</p> <ul style="list-style-type: none">Ideally one per column family. <p>HBase Region Flush – All contents of memstore flushed to an HFile on disk.</p> <ul style="list-style-type: none">Since Memstore is sorted, HFile is also sorted. | <p>Minor Compaction – Merge multiple, small HFiles into fewer larger ones.</p> <p>Major Compaction – Merge all HFiles into one large HFile with all records marked for deletion removed.</p> |
| <p>Data Model for Fast Writes and Reads</p> <ul style="list-style-type: none">Layout on physical disk is predictable minimizing disk seek.Fast access since get and put is by row key.Since tables are sorted, scan of a key range is fast. | <p>Creation of a Key Region – Table starts as a single range. When a region becomes too large, it splits into two child regions.</p> <p>Region Server – Initiates a region split. The child/daughter regions are opened on the same server.</p> | <p>HBase Use Case #1 – High Velocity and Volume Writes.</p> <ul style="list-style-type: none">Examples: Stock ticker, sensors, log files, system metrics. Real time monitoring. <p>HBase Use Case #2 – Information exchange with high velocity and volume read and write.</p> <ul style="list-style-type: none">Examples: Email, chat, Facebook. | <p>HBase Use Case #3 – High velocity and volume read.</p> <ul style="list-style-type: none">Examples: Content serving, web application back end, search index, online pre-computed view, online catalog. |

Lecture #13b – Recommendation

| | | | |
|---|--|--|---|
| <p>Input to Recommender Systems – Interactions between users and items.</p> <p>Output of Recommender Systems – Suggestions of additional interactions.</p> <p>Common Examples: Movies, music, restaurant choices, sale items at stores</p> | <p>Basis of Recommendation – <i>Behavior of a crowd helps us understand what individuals will do.</i></p> <p>Popular Items – Co-occur with everything making them not very useful for recommendation.</p> <p>Anomalous Co-Occurrence – Far more useful for recommendation. They are the source indicators of preference.</p> | <p>History Matrix – Constructed from the log files to show the history of the users for the items.</p> <ul style="list-style-type: none"> Users by items <p>Co-occurrence Matrix – Quantifies how often two items appear together.</p> <ul style="list-style-type: none"> Items by items <p>Indicator Matrix – Represent anomalous (i.e. interesting) co-occurrences.</p> <ul style="list-style-type: none"> Items by items | <p>Log Likelihood Ratio (LLR) – Can be helpful to judge with co-occurrences can be used with confidence as indicators of preference.</p> |
|---|--|--|---|

Cooccurrence Analysis

Which one is the anomalous co-occurrence?

| | A | not A |
|-------|------|---------|
| B | 13 | 1000 |
| not B | 1000 | 100,000 |

| | A | not A |
|-------|---|--------|
| B | 1 | 0 |
| not B | 0 | 10,000 |

| | A | not A |
|-------|----|---------|
| B | 10 | 0 |
| not B | 0 | 100,000 |

Bottom right is anomalous co-occurrence.
Consistently appear together or separate.

User-Based Filtering – Recommend items by finding similar users. Harder to scale.

Item-Based Filtering – Calculate similarity between items and make recommendations. This can be done offline.

Dithering – Random re-ordering of recommendation results. This makes for more off-line computation, but generally makes results better.

- Helps test relevancy of new items that would otherwise go unviewed.

- A** – Matrix of User search queries
 - Matrix is *Users* by *Queries*
- A^TA** – Gives a matrix of *Queries* by *Queries*
 - Query Co-occurrence Matrix**
 - This enables query recommendation.
- B** – Matrix of User video views
 - Video View Co-occurrence Matrix**
 - Matrix is *Users* by *Videos*
- B^TB** – Gives a matrix of *Videos* by *Videos*
 - This enables recommendations in the form “you may like these videos.”
- A^TB** – Gives query and video view co-occurrence. Just like a search engine.

Lecture #16a – Introduction to Recommendation

| | | | |
|---|--|--|---|
| <p>Recommendation – Class of machine learning that seeks to predict a user’s preference for or rating of an item.</p> <p>Two Main Approaches for Recommendation</p> <ul style="list-style-type: none"> Collaborative Filtering Content-based Filtering | <p>User Based Filtering – Recommendation are based on similarity to other users.</p> <p>Item-based Filtering – Recommendations are made based on similarity to other items.</p> | <p>Similarity Metrics</p> <ul style="list-style-type: none"> Pearson Correlation – Ratio of co-variance to product of standard deviations. <ul style="list-style-type: none"> -1 – Inversely Proportional 0 – No correlation 1 – Directly proportional Euclidean Distance – Coordinates indicate item preference. Smaller distance means more similarity. | <p>Tanimoto Coefficient – Ratio of intersection to union.</p> <ul style="list-style-type: none"> Between 0 and 1. Bigger is more similar. |
| <p>Collaborative-Based Filtering – Deductive. Needs user-item history to start to learn item associations.</p> <ul style="list-style-type: none"> Example: last.fm <p>Content Based Filtering – Inductive. Relies on domain specific knowledge to construct association between users and items.</p> <ul style="list-style-type: none"> Example: Pandora | <p>Challenges of Collaborative Filtering</p> <ul style="list-style-type: none"> Cold Start – No user history means no associations on day 1. Scale – Huge number of products of items and users means lots of computation. Sparsity – Most users express very little behavior with very few items and no behavior with the vast majority of items. | <p>Neighborhood – A group of similar users.</p> <p>Types of Neighborhoods:</p> <ul style="list-style-type: none"> Fixed Size – Cardinality (size) of the neighborhood is fixed in advance. Threshold-based – Cardinality is dependent based on a threshold of similarity that is fixed in advance. | <p>Precision – Proportion of top-scoring results that are relevant.</p> $precision = \frac{TP}{TP + FP} = \frac{a}{a + c}$ <p>Recall – Proportion of relevant results that are top-scoring.</p> $recall = \frac{a}{a + b} = \frac{TP}{TP + FN}$ |
| <p>Methods for Improving Recommender Performance</p> <ul style="list-style-type: none"> Change Model Type <ul style="list-style-type: none"> User Based Item Based Content Based Change Distance Metric <ul style="list-style-type: none"> Euclidean Tanimoto LogLikelihood Ratio (LLR) Model Parameters <ul style="list-style-type: none"> K in KNN | <p>Apriori Algorithm – Used to determine what people buy or use together.</p> <ul style="list-style-type: none"> Can be used for recommendation. <p>KNN Algorithm – Can be used for both classification and recommendation.</p> <p>Two Main Algorithms for Recommendation Engine</p> <ul style="list-style-type: none"> Collaborative Filtering Matrix Factorization | <p>Fundamental Problem of Collaborative Filtering – Finding a distance metric.</p> <p>Base Idea of Collaborative Filtering – Like KNN. Make decisions based off similar things (e.g. items, users, etc.). (Use collaboration from others to make recommendations)</p> | <p>Alternative Least Squares (ALS) – Inspired by matrix factorization but has almost nothing to do with it.</p> <ul style="list-style-type: none"> More similar to neural networks. Can be random or unstable, but it works fine for some cases. |