

CS 194 - Midterm Notes

Lecture 1: Introduction

Data Science Hacking Skill + Math & Statistics + Expertise
Research: Math & Statistics + Expertise
Machine Learning Hacking Skill + Math & Statistics

Database Engineer

- Experience with No-SQL systems: SQL systems, i.e. column-family stores, text and image DBs.
- Data cleaning: reading entities, fuzzy matching, bounds checking etc.
- Schema-on-read (or never) – late commitment semantics of data
- Data validation, using checksums, range checks statistical checks.

Scientific Computing Expert

- We were looking for the relevant terms in scientific computing toward integration of machine learning to address alternative to precise physical modeling learned in lecture one and of the project topics. So skill in modeling with ML is important.
- Another answer was all in using cluster tools (Spark, Hadoop, GraphLab etc.)

Machine Learning Expert

- Appropriateness of particular algorithms for various types of data
- Also how to feature data, BOW, n-grams, skip-grams, dependency parsing.
- Knowledge of ML tools: tuning for performance and accuracy

Judging Machine Learning Model

1. Identify problem
2. Instrument data sources
3. Collect data
4. Prepare data
5. Build model
6. Evaluate model
7. Communicate results

What's hard about Data Science?

1. Overcoming assumptions
2. Making ad hoc explanations of data patterns
3. Overgeneralizing

4. Communication
5. Not verifying enough (model validation, data pipeline integrity, etc.)

Lecture 2: Data Preparation

ETL: Extract, Transform, Load

1. Extract data from the source
2. Load data into a staging area
3. Transform data at source into a staging area

Some examples of sinks are: Python, R, SQLs, NoSQL stores. Schema can be anything, from structured to unstructured data

Schemas

Schema specify the structure and types of a data repository, and specify constraints within or between data fields. They are traditionally schema-on-write, and you can't write data into a table without a schema.

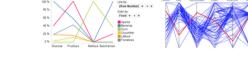
However, new data stores are schema-on-read or schemas.

CREATE SCHEMA Sprockets

```
CREATE TABLE KinFroggs (source int, cost int, partnumber int) GO
INSERT INTO KinFroggs (source, cost, partnumber) VALUES (5, 100, 45312453)
```

XML

- Generalizes HTML and helps to encode data, e.g. DOM
- Data can be represented and stored without schema
- More flexible than schema representation in DB
- Standard query/transformation languages like XSLT and XQuery
- Use Mark Logic Server as testing
- Example of defining an XML Schema: The right.



3. Radar Chart: One discrete variable with an arbitrary number of other variables.

Principal Component Analysis

Allows visualization of high-dimensional continuous data in 2D using principal components, and well from the orthogonal dimensions in the dataset. Great for Dimensionality Reduction.

Conclusion

Deceive of Marian Canal (false conclusions) - expect the unexpected!

Lecture 3: Manipulating Tabular Data

Structured Data

Data model is a collection of concepts for describing data.

Schema is a description of a particular collection of data.

Relational database is part of relations (true for argumental:

- At least two or more rows
- Schema specifies name, relation plus name and type of each column
- Instance: the actual data at a given time.

SQL is a Data Definition Language and Data Manipulation Language.

SQLs

Creating Tables

```
CREATE TABLE Students (sid CHAR(20), name CHAR(20), login CHAR(10), age INTEGER, gpa FLOAT);
```

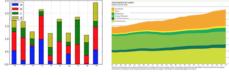
Adding and Deleting Tables

```
INSERT INTO Students (sid, name, login, age, gpa) VALUES ('53688', 'Smith', 'smithj', 18, 3.2);
```

Scatter plots can quickly explore relationships between two variables.

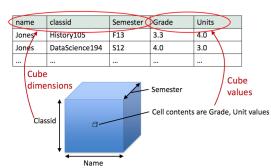
For multiple variables, we can use:

1. Stacked Plots, where stack variable is discrete.



2. Parallel Coordinate Plot: One discrete variable, and an arbitrary number of other variables.

Constructing OLAP cubes



Aggregation

- We can project a field on another axis, and then this is precomputed and maintained automatically in an OLAP cube, so queries are instantaneous. This helps to support real-time querying.
- However, aggregates increase space and cost of updates, but storage overhead is still small.
- Good integration with clients.
- NumPy and Matlab use dense nd-arrays which are OLAP



What's Wrong with RDBMS Tables?

- Indices: Typical RDBMS tables mostly use indices, but we can't afford this for large datasets.
- Transactions: Safe state changes require journals and are slow

- Relations: Checking relations adds further overhead to updates
- Sparse Data Support: Very wasteful, however there is both row-based and column-based

NoSQL Storage Systems

- Cassandra, CouchDB (JSON), BigTable (MapReduce jobs), Hive (relational database built on Hadoop, developed at Facebook)
- Could be dynamic columns, or no schema (like CouchDB)

Pig

• SQL functionality built on top of Hadoop

• Example Problem

An Example Problem

Suppose you have user data in a table and you want to find the most visited pages by user (top 10).

Users = LOAD 'hdfs://.../userdata' AS (name, age);

Pages = LOAD 'hdfs://.../pageviews' AS (user, page, count);

Visits = group User by name;

SortVisits = sort Visits by name desc;

Top10Visits = limit SortVisits 10;

store Top10Visits into 'top10';

In Pig Latin

Users = LOAD 'hdfs://.../userdata' AS (name, age);

Pages = LOAD 'hdfs://.../pageviews' AS (user, page, count);

Visits = group User by name;

SortVisits = sort Visits by name desc;

Top10Visits = limit SortVisits 10;

store Top10Visits into 'top10';

XML Schema

classifications

```
<element name="location">
<complexType>
<sequence>
<element name="latitude" type="xsd:decimal"/>
<element name="longitude" type="xsd:decimal"/>
</sequence>
</complexType>
<!-- location -->
```

An XML schema for this element:

```
<!-- location -->
<xs:schema xmlns:xsd="http://www.w3.org/2001/XMLSchema"
elementFormDefault="qualified">
<xsd:complexType name="location">
<xsd:sequence>
<xsd:element name="latitude" type="xsd:decimal"/>
<xsd:element name="longitude" type="xsd:decimal"/>
</xsd:sequence>
</xsd:complexType>
</xs:schema>
```

JSON

- Schemas data description language.
- Typically used to represent hierarchical data structures. Procedural in the target language.
- More succinct.
- Transformation/injection rely on code.
- Use MongoDB as modeling

Tablear Data

- Table is a collection of rows and columns
- Each row has an index, and each column has a name
- Cells in a row form a column pair
- Cells may or may not have a value
- Schema = minimally column types
- Can be stored as CSV or TSV

Syslog - A Standard for System Messages (Splunking)

• Grab data from many machines

• Check for errors, system failures, network congestion, security attacks

• Monitor resources: network, memory usage, disk use, latency, threads

Event-Driven Parsing: SAX

Finds all the open-close tag events in XML documents and does callbacks to user code.

- User code can respond to subset of events it's interested in
- User code can correctly compute aggregates from the data instead of implementing for each tag

• User code has to implement statement to keep track of where it is in the DOM tree

• User code can implement recovery strategies for XML errors

REST: Simple Object Access Protocol

- Web service is a software system designed to support interoperable machine to machine interaction over a network. There are two kinds of interactions: XML-based (SOAP) or RESTful interactions.
- SOAP RPC messages typically encode arguments that are presented to calling program as parameters and return values. This is longer because it requires a request-response cycle and verb, e.g. get/put, add/delete.

REST: Representation State Transfer

- Stateless Client/Server Protocol. Each message in the protocol contains all the information needed by the receiver to understand or process it. Keeps things simple, with set of well-defined operations (PUT, GET, ...). Resources are typically stored in a structured data format that supports hypermedia links, such as XML, JSON. Design emphasis is on nouns, e.g. user, location.

Data Cleaning

Openface is a spreadsheet-like tool checking data quality and reformats, substitutes, constraint check, etc.

Data Exploration

Histograms can show:

- Densities with long-tailed data. You can sort histogram counts by descending magnitude, and plot count vs bucket number on a log-log plot.

Reductions and Group By

Use GroupBy to select which columns to average, count, sum, average, min, max.

```
SELECT SID, Name, AVG(GPA)
FROM Students
GROUP BY SID
```

Python / Pandas

- Pandas is lightweight and fast. Full SQL expressiveness plus expressiveness of Python, integration with plotting functions like Matplotlib.
- Pandas: Tables must fit into memory. No post-load indexing functionality, indices are built when a table is created. No transactions or journaling. Large and complex joins are slow.

Series

- Keys of dictionary are indexes, built on NumPy's "ndarray".

Dataframe

- Tables with named columns represented as a dictionary
- Each Series object represents a column

Operations

- Map, filter, sort group by aggregate, pivot or reshape
- Relational functions, e.g. union, joins, renames

Online Analytical Processing - OLAP

- Like a n-dimensional spreadsheet (cube)
- Columns become dimensions, where some of the variables are used as qualifiers, whereas those that we want to measure are normally numeric

Three Important Tests

• T-test: compares two groups to see if they have the same distribution

• Chi-squared: test if the observed data matches the expected data.

• You could also have a single sample test for one group of individuals in two conditions.

• This is called a within-subjects design

• Chi-squared and Fisher's exact: compare counts in a contingency table

• Test whether an observation is consistent with the data

$$\sigma^2 = \sum_{i=1}^n (x_i - \bar{x})^2 / n$$

Fisher's exact test works when counts are under different conditions

• ANOVA (Analysis Of Variance)

• All the above tests are parametric

• Assume the data are

• normally distributed

• Samples are independent of each other and all have the same distribution (IID).

• Non-parametric tests (assume nothing about distribution)

• Wilcoxon signed rank: useful test for checking whether two (continuous or discrete) distributions are the same.

• Such as testing if a distribution is normal

Types of Error

• Type I error: P(x | H0) < p

• Type II error: P(x | H1) > 1 - p

• Power: 1 - P(x | H1) > p

• Type III error: P(x | H0) > 1 - p

• Type IV error: P(x | H1) < p

• Type V error: P(x | H0) > 1 - p

• Type VI error: P(x | H1) < p

• Type VII error: P(x | H0) < 1 - p

• Type VIII error: P(x | H1) > 1 - p

• Type IX error: P(x | H0) > 1 - p

• Type X error: P(x | H1) < 1 - p

• Type XI error: P(x | H0) < 1 - p

• Type XII error: P(x | H1) > 1 - p

• Type XIII error: P(x | H0) > 1 - p

• Type XIV error: P(x | H1) < 1 - p

• Type XV error: P(x | H0) < 1 - p

• Type XVI error: P(x | H1) > 1 - p

• Type XVII error: P(x | H0) > 1 - p

• Type XVIII error: P(x | H1) < 1 - p

• Type XVIX error: P(x | H0) < 1 - p

• Type XX error: P(x | H1) > 1 - p

• Type XXI error: P(x | H0) > 1 - p

• Type XXII error: P(x | H1) < 1 - p

• Type XXIII error: P(x | H0) < 1 - p

• Type XXIV error: P(x | H1) > 1 - p

• Type XXV error: P(x | H0) > 1 - p

• Type XXVI error: P(x | H1) < 1 - p

• Type XXVII error: P(x | H0) < 1 - p

• Type XXVIII error: P(x | H1) > 1 - p

• Type XXIX error: P(x | H0) > 1 - p

• Type XXX error: P(x | H1) < 1 - p

• Type XXXI error: P(x | H0) < 1 - p

• Type XXXII error: P(x | H1) > 1 - p

• Type XXXIII error: P(x | H0) > 1 - p

• Type XXXIV error: P(x | H1) < 1 - p

• Type XXXV error: P(x | H0) > 1 - p

• Type XXXVI error: P(x | H1) < 1 - p

• Type XXXVII error: P(x | H0) > 1 - p

• Type XXXVIII error: P(x | H1) < 1 - p

• Type XXXIX error: P(x | H0) > 1 - p

• Type XXXX error: P(x | H1) < 1 - p

• Type XXXXI error: P(x | H0) > 1 - p

• Type XXXXII error: P(x | H1) < 1 - p

• Type XXXXIII error: P(x | H0) > 1 - p

• Type XXXXIV error: P(x | H1) < 1 - p

• Type XXXXV error: P(x | H0) > 1 - p

• Type XXXXVI error: P(x | H1) < 1 - p

• Type XXXXVII error: P(x | H0) > 1 - p

• Type XXXXVIII error: P(x | H1) < 1 - p

• Type XXXXIX error: P(x | H0) > 1 - p

• Type XXXXX error: P(x | H1) < 1 - p

• Type XXXXXI error: P(x | H0) > 1 - p

• Type XXXXXII error: P(x | H1) < 1 - p

• Type XXXXXIII error: P(x | H0) > 1 - p

• Type XXXXXIV error: P(x | H1) < 1 - p

• Type XXXXXV error: P(x | H0) > 1 - p

• Type XXXXXVI error: P(x | H1) < 1 - p

• Type XXXXXVII error: P(x | H0) > 1 - p

• Type XXXXXVIII error: P(x | H1) < 1 - p

• Type XXXXXIX error: P(x | H0) > 1 - p

• Type XXXXXI error: P(x | H1) < 1 - p

• Type XXXXXII error: P(x | H0) > 1 - p

• Type XXXXXIII error: P(x | H1) < 1 - p

• Type XXXXXIV error: P(x | H0) > 1 - p

• Type XXXXXV error: P(x | H1) < 1 - p

• Type XXXXXVI error: P(x | H0) > 1 - p

• Type XXXXXVII error: P(x | H1) < 1 - p

• Type XXXXXVIII error: P(x | H0) > 1 - p

• Type XXXXXIX error: P(x | H1) < 1 - p

- Samples**
- We are most interested in models of the population, but we have access only to a sample
 - We train on a training sample D and denote the model as $\hat{f}(D)$
 - Bias is the error model $\hat{f}(D)$ on many training sets. $\hat{f}(D)$ is the expected difference between their predictions and the true y .
 - Variance is the error models $\hat{f}(D)$ on many training sets. $\hat{f}(D)$ is the variance of the estimate.
 - There is usually a bias-variance tradeoff caused by model complexity.
 - Complex models (many parameters) have lower bias, but higher variance.
 - Simple models (few parameters) have higher bias, but lower variance.
 - The total expected error is $B^2 + V$ known as the total error.
 - If variance strongly dominates, it means there is too much variation between models. This is called over-fitting.
 - If variance is dominant, then the models are not fitting the data well enough. This is called under-fitting.

K-nearest Neighbors (kNN)

- Given a query item. Find k closest matches in a labeled dataset
- The Data is the Model: No training needed.
- Classification

 - Model is a $\hat{y} = \text{vote}$, y_i from a discrete set (labels).
 - Given X , compute y – majority vote of the k nearest neighbors.
 - Can also use a weighted ‘vote’ of the neighbors.

- Regression

 - Model is a $\hat{y} = \bar{y}$, y_i from the data.
 - Given X , compute \bar{y} – average value of the k nearest neighbors.
 - Can also use weighted average of the neighbors.
 - Weight function usually the inverse distance.
 - Small k → low bias, high variance.
 - Large k → high bias, low variance.
 - This is a very simple approach to phenomena that occur in high dimensions (100s to millions) that do not occur in low-dimensional (e.g. 3-dimensional) space.
 - Data in high dimensions tend to be very sparse.
 - Means there are few points that are close to the point we want to predict.

Naive Bayes

D = Data, e = some event

$$P(e|D) = \frac{P(D|e)P(e)}{P(D)}$$

- P(e)** is called the prior probability of e. It's what we know (or think we know) about e with no other evidence.
- P(D|e)** is the conditional probability of D given that e happened, or just the likelihood of D. This can often be measured or computed precisely – it follows from your model assumptions.
- P(D)** is the marginal probability of D given D. It's the answer we were looking for when we chose a best hypothesis.

You can see that the posterior is heavily colored by the prior, so Bayes' is not used to test hypotheses

- Key Assumption:** The features are generated independently given e .
 - Simple and fast, depends only on term frequency data for the classes. One pass over the training dataset, no iteration.
 - Compact model. Size is proportional to number of n features.
 - Very well behaved metric. Test weight depends only on frequency of that term.
 - Decomposes into terms.
 - Can work very well with sparse data, where combinations of dependent terms are rare.
 - Subject to error and bias when term probabilities are not independent (e.g. URL prediction).
 - Can't model patterns in the data.
 - Typically not as accurate as other methods for the above reasons.
 - If you care, here is the formula we are trying to maximize with Naive Bayes (The classifier is trying to fold a category y that maximizes)

We want the category c_j that maximizes:

$$\Pr(X_1, \dots, X_k | c_j) \Pr(c_j) = \Pr(c_j) \prod_{i=1}^k \Pr(X_i | c_j)$$

Lecture 7: Machine Learning 2

Linear Regression

$$X^1 = X^0 - H_f^{-1}(X^0) G_f(X^0)$$

Where $G_f = \frac{df}{dx_i}$ is the gradient and

$$H_f = \frac{d^2 f}{dx_i dx_j} \text{ (a matrix)} \text{ is the Hessian.}$$

- Converges very fast when feasible
- Challenges:
 - Both the gradient and Hessian are normally computed on the entire dataset. This may require many passes over the data to complete a series of Newton steps.
 - The Hessian has size $O(M^2)$ if there are M features. This is impractical for large feature spaces, e.g. text or event data.

Stochastic Gradient

- A very important set of iterative algorithms use stochastic gradient updates.
- They use a small subset or mini-batch of the data, and use it to compute a gradient update.
- Stochastic gradient has some serious limitations however, especially if the gradients vary widely in magnitude. Some coefficients change very fast, others very slowly.
- For harder (non-convex) learning problems, its common to use learning rate schedules that decay over time.

Decision Tree

- Walk from root to a class-labeled leaf.
- Attack node, branch based on the value of some feature.
- At each node, we choose the feature f which maximizes the information gain, which is the information loss.
- Information Gain:
 - Entropy is defined as each node based on the class breakdown as a function over the probability distribution
 - Equation for entropy by feature f (keep in mind we want to maximize information gain while avoiding overfitting)

Before the split by f , entropy is

$$E = -\sum_{i=1}^n p_i \log p_i$$

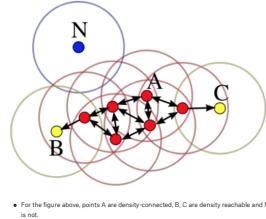
After split by f , the entropy is

$$E_f = -p^f \sum_{i=1}^m p_i^f \log p_i^f - p^{1-f} \sum_{i=1}^m p_i^{1-f} \log p_i^{1-f}$$

The information gain = $E - E_f$ (information = -entropy)

Ensemble Methods

- Crowdsourced machine learning: take a bunch of weak learning algorithms and combine them to make a good one.
 - Types:
 - Bagging: trains learners in parallel on different samples of the data, then combine by voting (noisy effect to avoid averaging to consensus result).
 - Stacking: combine model outputs using a second stage learner like linear regression.
 - Boosting: train learners on the filtered output of other learners.
 - Random Forests
- The basic premise is that we don't want the models to be similar so we split them across different learners, train them on that data, and decided on the processing of our actual data in the ensemble based on the results.
- Given N trees on datasets sampled from the original dataset with replacement (Bootstrap sampled) [$p = \text{number of features}$]
 - Draw a bootstrap sample of size m . Each tree is trained on a different dataset.
 - Draw a bootstrap sample of size m from a random set of m of all features at each node, and choosing the best feature to split on (corresponds features in a nearby tree can't split on the same feature).
 - Typically m will be \sqrt{p} but can be smaller.
 - Random forest is very popular on data with fewer features (dense). Easily parallelizable (fast), not the most accurate, needs at least the max depth of the tree passed in and is easy to overfit.



- For the figure above, points A are density-connected, B, C are density reachable and N is not.
- The algorithm overall can be run in $O(n \log n)$.

Matrix Factorization

- Matrix factorization is the curse of dimensionality when the model has more degrees of freedom (variables) in the output than we have in the data.
- When having to deal with thousands or millions of responses, as well as thousands or millions of features.
- High dimensional data will often have simpler structure, e.g.
 - Topics in documents
 - General product characteristics (availability, cost etc.)
- We can often find a low-rank matrix affecting a “lower-dimensional” model.
- The subspace allows us to predict the value of other features from known features.
- Usually, data vary more strongly in some dimensions than others. Fitting a linear model to the (K) strongest directions gives the best approximation to the data in the least-squares sense.
- Algebraically, we have a high-dimensional data matrix (typically sparse) S , which we approximate as a product of two denumerable A and B with low “inner” dimension:
 - “Matrix Completion”
- Some complex math and SGD, or Alternating Least Squares can be used to find the

factorization (think we should focus more on high level and applications than the math)

Performance

- There are two different dimensions of performance:
 - Training time vs. number of features.
 - Best possible model accuracy.
- Online Model: Predicting items, limited resources (memory and machines). Goal is usually to minimize latency and reduce memory usage.
- There are at least three approaches to improving performance:
 - Algorithmic improvements: e.g. adaptive SGD methods.
 - Feasible Order of magnitude improvements possible.
 - Compounds algorithmic improvements. Compounds algorithmic improvements.
- Clustered: Distribute model or data across computers in a network.
- Supervised cost: A full field requires speedup, $f(x)$ + k additional cost. Usually does not compound algorithmic and computational improvements (network bottleneck).

Learning 9: Deep Learning

Deep Learning is a branch of machine learning specifically for modeling complex and high-level data by decomposing layers of learned data built upon the previous layer.

Community is it used for Computer Vision and Voice recognition.

Loss (cost) is a function that specifies the goal of a learning task and marks the values against some inadequacy metric called badness.

Risk is measuring loss over the input. Since we can't try all the values in an unbounded domain, we estimate it using a summation function and call the resulting value **Empirical Risk**.

Risk is measuring loss over the input. Since we can't try all the values in an unbounded domain, we estimate it using a summation function and call the resulting value **Empirical Risk**.

Recall that we are trying to learn a function f that minimizes the loss function and call the resulting function \hat{f} .

Loss function is a function that specifies the goal of a learning task and marks the values against some inadequacy metric called badness.

Risk is measuring loss over the input. Since we can't try all the values in an unbounded domain, we estimate it using a summation function and call the resulting value **Empirical Risk**.

Risk is measuring loss over the input. Since we can't try all the values in an unbounded domain, we estimate it using a summation function and call the resulting value **Empirical Risk**.

Recall that we are trying to learn a function f that minimizes the loss function and call the resulting function \hat{f} .

Loss function is a function that specifies the goal of a learning task and marks the values against some inadequacy metric called badness.

Risk is measuring loss over the input. Since we can't try all the values in an unbounded domain, we estimate it using a summation function and call the resulting value **Empirical Risk**.

Recall that we are trying to learn a function f that minimizes the loss function and call the resulting function \hat{f} .

Loss function is a function that specifies the goal of a learning task and marks the values against some inadequacy metric called badness.

Risk is measuring loss over the input. Since we can't try all the values in an unbounded domain, we estimate it using a summation function and call the resulting value **Empirical Risk**.

Recall that we are trying to learn a function f that minimizes the loss function and call the resulting function \hat{f} .

Loss function is a function that specifies the goal of a learning task and marks the values against some inadequacy metric called badness.

Risk is measuring loss over the input. Since we can't try all the values in an unbounded domain, we estimate it using a summation function and call the resulting value **Empirical Risk**.

Recall that we are trying to learn a function f that minimizes the loss function and call the resulting function \hat{f} .

Loss function is a function that specifies the goal of a learning task and marks the values against some inadequacy metric called badness.

Risk is measuring loss over the input. Since we can't try all the values in an unbounded domain, we estimate it using a summation function and call the resulting value **Empirical Risk**.

Recall that we are trying to learn a function f that minimizes the loss function and call the resulting function \hat{f} .

Loss function is a function that specifies the goal of a learning task and marks the values against some inadequacy metric called badness.

Risk is measuring loss over the input. Since we can't try all the values in an unbounded domain, we estimate it using a summation function and call the resulting value **Empirical Risk**.

Recall that we are trying to learn a function f that minimizes the loss function and call the resulting function \hat{f} .

Loss function is a function that specifies the goal of a learning task and marks the values against some inadequacy metric called badness.

Risk is measuring loss over the input. Since we can't try all the values in an unbounded domain, we estimate it using a summation function and call the resulting value **Empirical Risk**.

Recall that we are trying to learn a function f that minimizes the loss function and call the resulting function \hat{f} .

Loss function is a function that specifies the goal of a learning task and marks the values against some inadequacy metric called badness.

Risk is measuring loss over the input. Since we can't try all the values in an unbounded domain, we estimate it using a summation function and call the resulting value **Empirical Risk**.

Recall that we are trying to learn a function f that minimizes the loss function and call the resulting function \hat{f} .

Loss function is a function that specifies the goal of a learning task and marks the values against some inadequacy metric called badness.

Risk is measuring loss over the input. Since we can't try all the values in an unbounded domain, we estimate it using a summation function and call the resulting value **Empirical Risk**.

Recall that we are trying to learn a function f that minimizes the loss function and call the resulting function \hat{f} .

Loss function is a function that specifies the goal of a learning task and marks the values against some inadequacy metric called badness.

Risk is measuring loss over the input. Since we can't try all the values in an unbounded domain, we estimate it using a summation function and call the resulting value **Empirical Risk**.

Recall that we are trying to learn a function f that minimizes the loss function and call the resulting function \hat{f} .

Loss function is a function that specifies the goal of a learning task and marks the values against some inadequacy metric called badness.

Risk is measuring loss over the input. Since we can't try all the values in an unbounded domain, we estimate it using a summation function and call the resulting value **Empirical Risk**.

Recall that we are trying to learn a function f that minimizes the loss function and call the resulting function \hat{f} .

Loss function is a function that specifies the goal of a learning task and marks the values against some inadequacy metric called badness.

Risk is measuring loss over the input. Since we can't try all the values in an unbounded domain, we estimate it using a summation function and call the resulting value **Empirical Risk**.

Recall that we are trying to learn a function f that minimizes the loss function and call the resulting function \hat{f} .

Loss function is a function that specifies the goal of a learning task and marks the values against some inadequacy metric called badness.

Risk is measuring loss over the input. Since we can't try all the values in an unbounded domain, we estimate it using a summation function and call the resulting value **Empirical Risk**.

Recall that we are trying to learn a function f that minimizes the loss function and call the resulting function \hat{f} .

Loss function is a function that specifies the goal of a learning task and marks the values against some inadequacy metric called badness.

Risk is measuring loss over the input. Since we can't try all the values in an unbounded domain, we estimate it using a summation function and call the resulting value **Empirical Risk**.

Recall that we are trying to learn a function f that minimizes the loss function and call the resulting function \hat{f} .

Loss function is a function that specifies the goal of a learning task and marks the values against some inadequacy metric called badness.

Risk is measuring loss over the input. Since we can't try all the values in an unbounded domain, we estimate it using a summation function and call the resulting value **Empirical Risk**.

Recall that we are trying to learn a function f that minimizes the loss function and call the resulting function \hat{f} .

Loss function is a function that specifies the goal of a learning task and marks the values against some inadequacy metric called badness.

Risk is measuring loss over the input. Since we can't try all the values in an unbounded domain, we estimate it using a summation function and call the resulting value **Empirical Risk**.

Recall that we are trying to learn a function f that minimizes the loss function and call the resulting function \hat{f} .

Loss function is a function that specifies the goal of a learning task and marks the values against some inadequacy metric called badness.

Risk is measuring loss over the input. Since we can't try all the values in an unbounded domain, we estimate it using a summation function and call the resulting value **Empirical Risk**.

Recall that we are trying to learn a function f that minimizes the loss function and call the resulting function \hat{f} .

Loss function is a function that specifies the goal of a learning task and marks the values against some inadequacy metric called badness.

Risk is measuring loss over the input. Since we can't try all the values in an unbounded domain, we estimate it using a summation function and call the resulting value **Empirical Risk**.

Recall that we are trying to learn a function f that minimizes the loss function and call the resulting function \hat{f} .

Loss function is a function that specifies the goal of a learning task and marks the values against some inadequacy metric called badness.

Risk is measuring loss over the input. Since we can't try all the values in an unbounded domain, we estimate it using a summation function and call the resulting value **Empirical Risk**.

Recall that we are trying to learn a function f that minimizes the loss function and call the resulting function \hat{f} .

Loss function is a function that specifies the goal of a learning task and marks the values against some inadequacy metric called badness.

Risk is measuring loss over the input. Since we can't try all the values in an unbounded domain, we estimate it using a summation function and call the resulting value **Empirical Risk**.

Recall that we are trying to learn a function f that minimizes the loss function and call the resulting function \hat{f} .

Loss function is a function that specifies the goal of a learning task and marks the values against some inadequacy metric called badness.

Risk is measuring loss over the input. Since we can't try all the values in an unbounded domain, we estimate it using a summation function and call the resulting value **Empirical Risk**.

Recall that we are trying to learn a function f that minimizes the loss function and call the resulting function \hat{f} .

Loss function is a function that specifies the goal of a learning task and marks the values against some inadequacy metric called badness.

Risk is measuring loss over the input. Since we can't try all the values in an unbounded domain, we estimate it using a summation function and call the resulting value **Empirical Risk**.

Recall that we are trying to learn a function f that minimizes the loss function and call the resulting function \hat{f} .

Loss function is a function that specifies the goal of a learning task and marks the values against some inadequacy metric called badness.

Risk is measuring loss over the input. Since we can't try all the values in an unbounded domain, we estimate it using a summation function and call the resulting value **Empirical Risk**.

Recall that we are trying to learn a function f that minimizes the loss function and call the resulting function \hat{f} .

Loss function is a function that specifies the goal of a learning task and marks the values against some inadequacy metric called badness.

Risk is measuring loss over the input. Since we can't try all the values in an unbounded domain, we estimate it using a summation function and call the resulting value **Empirical Risk**.

Recall that we are trying to learn a function f that minimizes the loss function and call the resulting function \hat{f} .

Loss function is a function that specifies the goal of a learning task and marks the values against some inadequacy metric called badness.

Risk is measuring loss over the input. Since we can't try all the values in an unbounded domain, we estimate it using a summation function and call the resulting value **Empirical Risk**.

Recall that we are trying to learn a function f that minimizes the loss function and call the resulting function \hat{f} .

Loss function is a function that specifies the goal of a learning task and marks the values against some inadequacy metric called badness.

Risk is measuring loss over the input. Since we can't try all the values in an unbounded domain, we estimate it using a summation function and call the resulting value **Empirical Risk**.

Recall that we are trying to learn a function f that minimizes the loss function and call the resulting function \hat{f} .

Loss function is a function that specifies the goal of a learning task and marks the values against some inadequacy metric called badness.

Risk is measuring loss over the input. Since we can't try all the values in an unbounded domain, we estimate it using a summation function and call the resulting value **Empirical Risk**.

Recall that we are trying to learn a function f that minimizes the loss function and call the resulting function \hat{f} .

Loss function is a function that specifies the goal of a learning task and marks the values against some inadequacy metric called badness.

Risk is measuring loss over the input. Since we can't try all the values in an unbounded domain, we estimate it using a summation function and call the resulting value **Empirical Risk**.

Recall that we are trying to learn a function f that minimizes the loss function and call the resulting function \hat{f} .

Loss function is a function that specifies the goal of a learning task and marks the values against some inadequacy metric called badness.

Risk is measuring loss over the input. Since we can't try all the values in an unbounded domain, we estimate it using a summation function and call the resulting value **Empirical Risk**.

Recall that we are trying to learn a function f that minimizes the loss function and call the resulting function \hat{f} .

Loss function is a function that specifies the goal of a learning task and marks the values against some inadequacy metric called badness.

Risk is measuring loss over the input. Since we can't try all the values in an unbounded domain, we estimate it using a summation function and call the resulting value **Empirical Risk**.

Recall that we are trying to learn a function f that minimizes the loss function and call the resulting function \hat{f} .

Loss function is a function that specifies the goal of a learning task and marks the values against some inadequacy metric called badness.

Risk is measuring loss over the input. Since we can't try all the values in an unbounded domain, we estimate it using a summation function and call the resulting value **Empirical Risk**.

Recall that we are trying to learn a function f that minimizes the loss function and call the resulting function \hat{f} .

Loss function is a function that specifies the goal of a learning task and marks the values against some inadequacy metric called badness.

Risk is measuring loss over the input. Since we can't try all the values in an unbounded domain, we estimate it using a summation function and call the resulting value **Empirical Risk**.

Recall that we are trying to learn a function f that minimizes the loss function and call the resulting function \hat{f} .

Loss function is a function that specifies the goal of a learning task and marks the values against some inadequacy metric called badness.

Risk is measuring loss over the input. Since we can't try all the values in an unbounded domain, we estimate it using a summation function and call the resulting value **Empirical Risk**.

Recall that we are trying to learn a function f that minimizes the loss function and call the resulting function \hat{f} .

Loss function is a function that specifies the goal of a learning task and marks the values against some inadequacy metric called badness.

- Transaction Networks
- Biological Networks
- Co-Authorship Networks
- etc.

Degree distribution follows a power law. (there are very few high degree nodes, but many low degree nodes).

Giant Connected Component:

- Typically only one GIANT CC
- CC are components that can reach each other.

Desirability - Average distance between nodes decreases over time.

Things to look at when given a graph:

- Density (typically ~2%)
- Degree distribution
 - Skewed?
 - Number of high degree vertices (who's the person at the middle of the social graph)
- Connected Component Sizes
- Community Profile Plots (does it fit some community pattern)
- Structure around large, important vertices (something like a bridge between two connected components)

Page Rank

- Recursive Relationship:

$$R[i] = \alpha + (1 - \alpha) \sum_{(j,i) \in E} \frac{1}{L[j]} R[j]$$

- α is the random reset probability (typically 0.15)
- $L[i]$ is the number of links on page

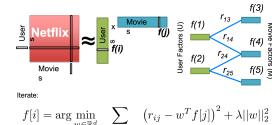
Predicting Behavior

We can use a similar recursive structure to predict behavior by analyzing behavior of connections.

- Low-Rank Matrix Factorization:
- Used for recommending products.
- Here's the function for it.

Recommending Products

Low-Rank Matrix Factorization:



Iterate:

$$f[i] = \arg \min_{w \in \mathbb{R}^d} \sum_{j \in \text{Nbrs}(i)} (r_{ij} - w^T f[j])^2 + \lambda \|w\|_2^2$$

Clustering Coefficient

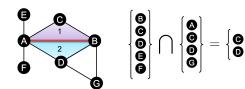
- How tight are the communities?
- Specifically how tight knit is a community around a single node.

$$\text{Measure "cohesiveness" of local community}$$

$$\text{ClusterCoeff}[i] = \frac{2 * \# \text{Triangles}[i]}{\text{Deg}[i] * (\text{Deg}[i] - 1)}$$

- You need to be able to count triangles.

- Counting Triangles tells you how many shared connections a particular node's connections have. Take the intersection of two nodes out degrees to determine their shared connections.



Connected Components:

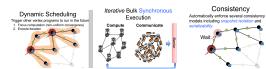
- Describes nodes connected to each other by some sequence of connections.

Dynamic Working Sets:

- Many entries (SG) only run once!
- But there's a long tail, a few entries run many times (change many times)

Summary:

- Graph algorithms address transitive structure
 - Iterative
 - Dynamic working sets
 - Many graph algorithms can be specified locally
 - Define a set at one level or a vertex and its neighbors
- Graph Algorithm implementation
 - user defined vertex program has access to the properties of neighboring edges and vertices. They interact by sending messages.
 - When messages come in they might update their vertex's property.
 - These messages can be sent to neighboring vertices with messages.
 - Dynamic Scheduling → Assign jobs to be calculated in bulk. → Bulk Synchronous Execution ensures consistency!



Commutative Associative Vertex Programs:

- Gather - Get messages
- Apply - Compute a new vertex property
- Scatter - Scatter that new information (in the form of messages)

Power law graphs are difficult to partition:

- They don't have low cost balanced cuts.
 - Cut Edges
 - Cut Vertices
 - Instead of partitioning by cutting edges, cut vertices!

You can treat Map Reduce (Spark) jobs as graph problems (joins can be solved more efficiently as a graph problem)