

CS 295B/CS 395B
Systems for Knowledge
Discovery

Toward Sound Data
Collection



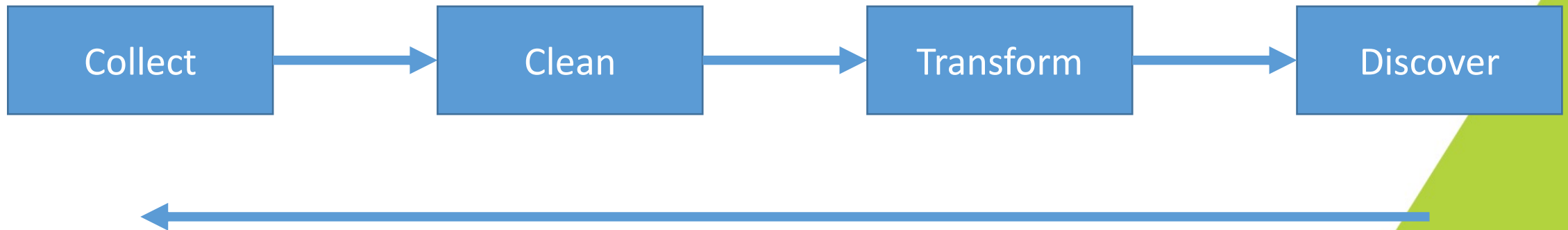
The University of Vermont

Topics for Today

- ~~Parisa's presentation~~
- Contextualizing Friday's papers
- Background for Friday's papers
- New method: ablation studies
- Things to reflect on for this week's reviews and presentations

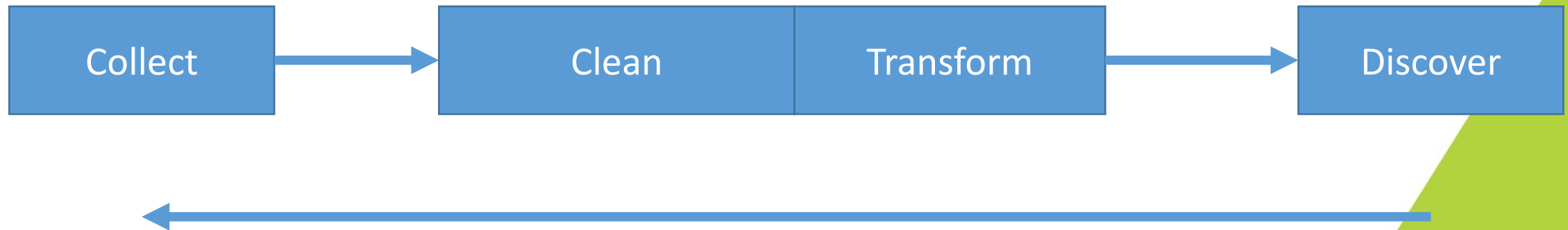
Context

Remember this diagram...

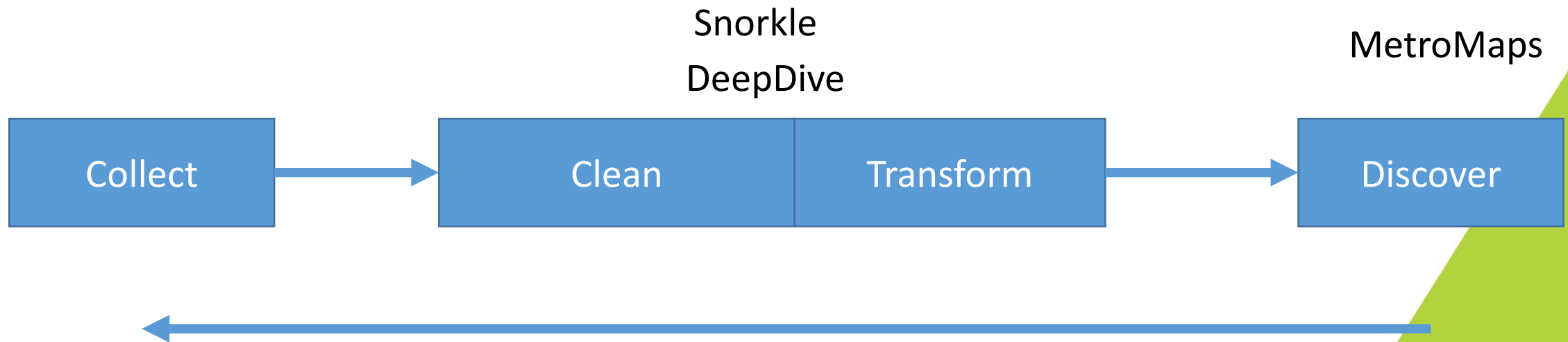


Remember this diagram...

Snorkle
DeepDive



Remember this diagram...





[Home](#) [Vision](#) [Slides](#) [People](#) [Members](#) [Projects](#) [Papers](#) [Seminar](#) [Blog](#)

DAWN is a five-year research project to democratize AI by making it dramatically easier to build AI-powered applications.

Our past research—from [Apache Spark](#) to [Mesos](#), [DeepDive](#), and [HogWild!](#)—already powers major functionality all over Silicon Valley and the world. Between [fighting against human trafficking](#), [assisting in cancer diagnosis](#) and [performing high-throughput genome sequencing](#), we’ve invested heavily in tools for AI and data product development.

The [next step is to make these tools more efficient and more accessible](#), from training set creation and model design to monitoring, efficient execution, and hardware-efficient implementation. This technology holds the power to change science and society—and we’re creating this change with partners throughout campus and beyond.

We’re proud to be supported by the following founding members:



Background: DAWN project

Goal: “To empower domain experts who are not ML experts”

Achieved through many initiatives

Our focus: **Developing new interfaces that**

- Make model specification easier
- Explain results to humans
- Make debugging easier
- Make improving data quality easier

Hardware Systems Algorithms Interfaces

Data Acquisition

Feature Extraction

Model Training

Productionizing

Snorkel

DeepDive

ModelSnap

ModelQA

MacroBase (Streaming Data)

Data Fusion

NoScope (Video)

AutoRec, SimDex (Recommendation)

Mulligan (SQL+graph+ML)

End-to-End Compilers: Weld, Delite

New Hardware: FuzzyBit, Plasticine CGRA



CPU



GPU



FPGA



Cluster



Mobile

...

Interfaces
Algorithms
Systems
Hardware

Data Acquisition

Feature Extraction

Model Training

Productionizing

Snorkel

DeepDive

ModelSnap

ModelQA

MacroBase (Streaming Data)

Data Fusion

NoScope (Video)

AutoRec, SimDex (Recommendation)

Mulligan (SQL+graph+ML)

End-to-End Compilers: Weld, Delite

New Hardware: FuzzyBit, Plasticine CGRA



CPU



GPU



FPGA



Cluster



Mobile

...

Interfaces
Algorithms
Systems
Hardware

Data Acquisition

Feature Extraction

Model Training

Productionizing

Snorkel

DeepDive

ModelSnap

ModelQA

MacroBase (Streaming Data)

Data Fusion

NoScope (Video)

AutoRec, SimDex (Recommendation)

Mulligan (SQL+graph+ML)

End-to-End Compilers: Weld, Delite

New Hardware: FuzzyBit, Plasticine CGRA



CPU



GPU



FPGA



Cluster

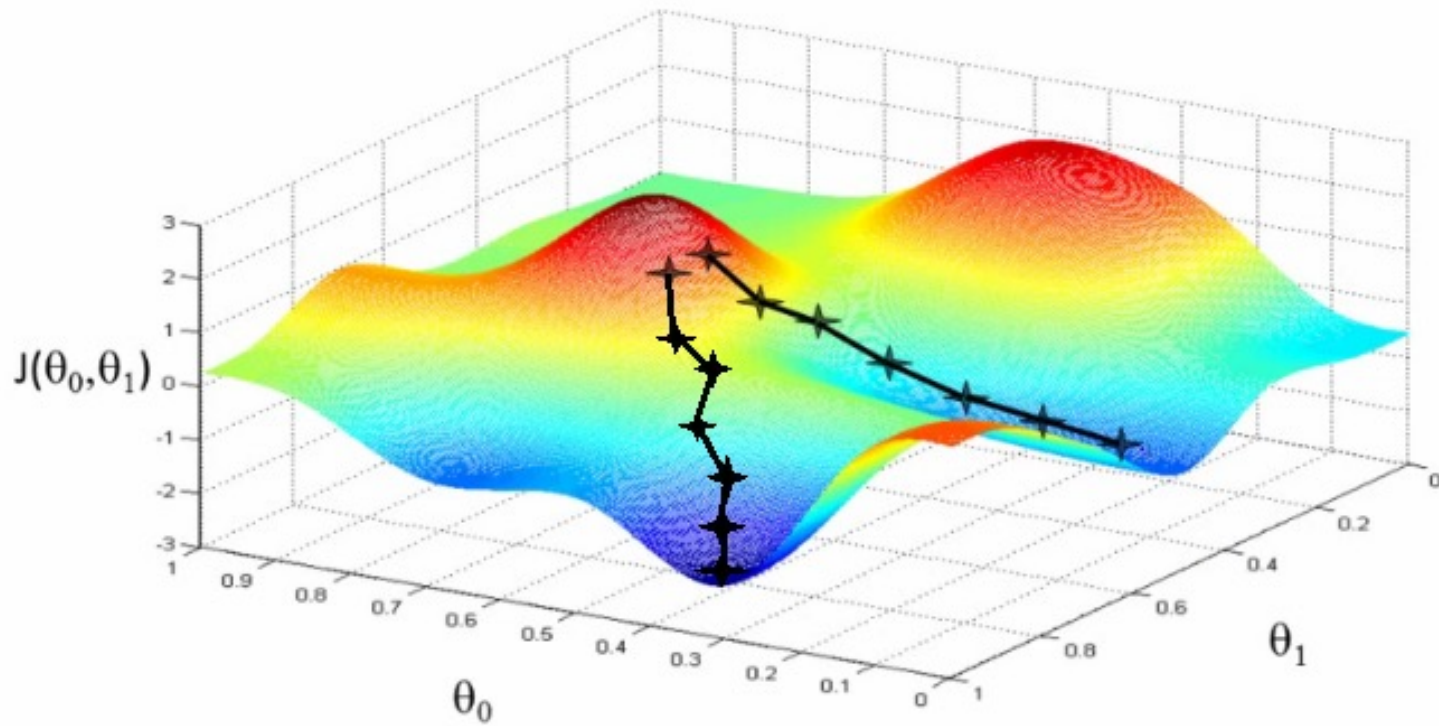


Mobile

...

Questions?

SGD and infinite data



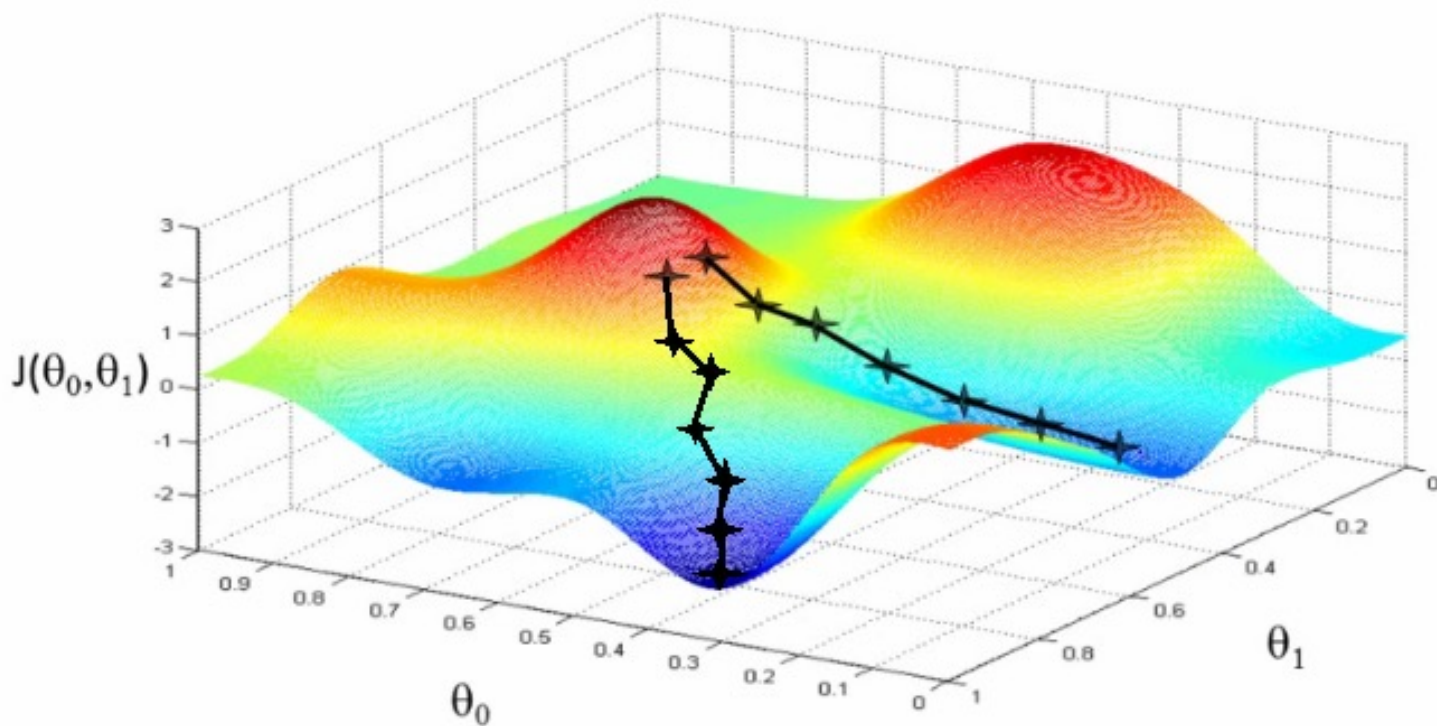
Stochastic randomly

Gradient follow the derivative

Descent toward the minimum

See Andrew Ng's videos for details

SGD and infinite data

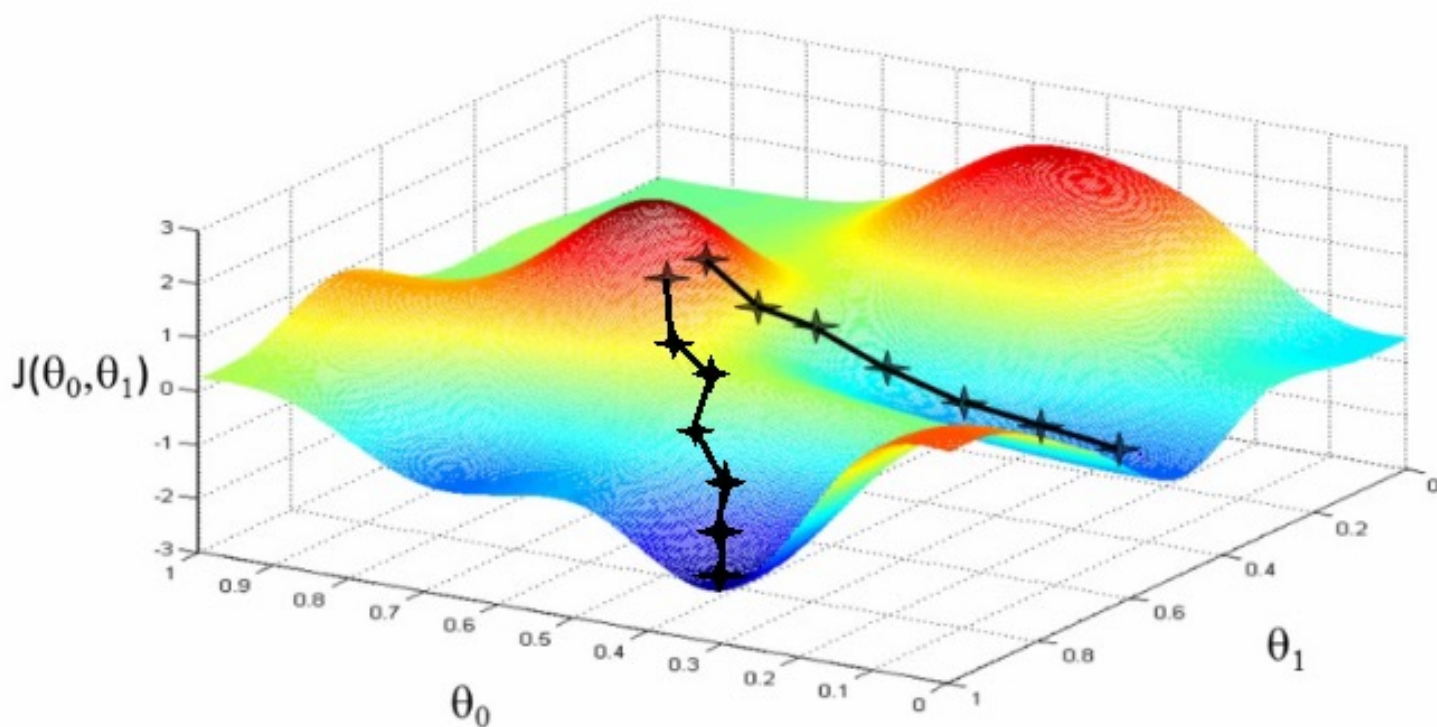


SGD powers all modern machine learning...

... especially deep learning ...

Idea: more data is better

SGD and infinite data



But how to get more data?

Better question:

How to get more quality data?

Gold standard:

Quality *labelled* data

- Expensive
- DeepDive + Snorkle

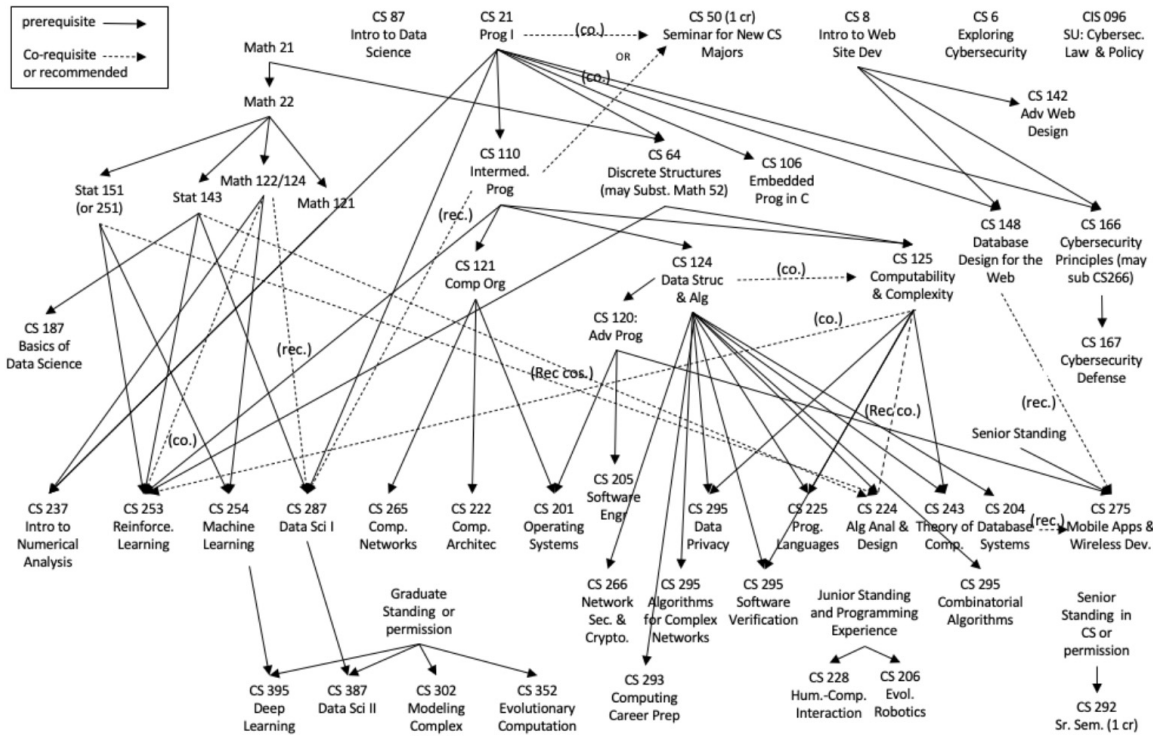
Questions?

Datalog

Snorkle/DeepDive interfaces include Datalog

- A domain-specific language (DSL) for *logic programming*
- Declaratively write out facts and relations
- Can then query this set of rules

Computer Science Course Prerequisite Graph (last updated 11/03/2020)



NOTE: This is not a complete listing, but includes most currently offered courses. Elective offerings, especially at the 2xx and 3xx level, may vary over time.

Datalog Example

This is the CEMS CS major pre-req chart

(It is out of date)

We can express with datalog:

Required("CS50")

Prereq("CS121", "CS222")

IsNot(Standing(student), "graduate") :-

HasPermission(student, "CS352")

?- CanTake(student, "CS204")

?- CanTake(student)

Questions?

Some prior work

Connecting the Dots Between News Articles

Dafna Shahaf
Carnegie Mellon University
dshahaf@cs.cmu.edu

Carlos Guestrin
Carnegie Mellon University
guestrin@cs.cmu.edu

ABSTRACT

The process of extracting useful knowledge from large datasets has become one of the most pressing problems in today's society. The problem spans entire sectors, from scientists to intelligence analysts and web users, all of whom are constantly struggling to keep up with the larger and larger amounts of content published every day. With this much data, it is often easy to miss the big picture.

In this paper, we investigate methods for automatically connecting the dots — providing a structured, easy way to navigate within a new topic and discover hidden connections within the news domain: given two news articles, our system automatically finds a coherent chain of events starting with the decline of health-care debate (2007), and ending with the decline of home prices in 2007. We formalize the characteristics of a good chain and propose an efficient algorithm (with theoretical guarantees) to find an efficient algorithm (with theoretical guarantees) to connect two fixed endpoints. We incorporate user feedback into our framework, allowing the algorithm to learn from personalized data. Our user studies demonstrate the effectiveness in helping users understanding the news.

Categories and Subject Descriptors: I.2.6 [Artificial Intelligence]: Learning; G.3 [Probability and Statistics]

General Terms: Algorithms, Experimentation

1. INTRODUCTION

"Can't a Group Credit Crisis? Join the Club", stated David Lesh's article in the New York Times. Credit crisis had been going on for seven months by that time, and had been extensively covered by every major media outlet throughout the world. Yet many people felt as if they did not understand what it was about.

Paradoxically, the extensive media coverage might have been a part of the problem. This is another instance of the information overload problem, long recognized in the

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.
KDD'10, July 25–28, 2010, Washington, DC, USA.
Copyright 2010 ACM 978-1-4503-0051-1/1007...\$10.00.

computing industry. Users are constantly struggling to keep up with the larger and larger amounts of content that is being published every day; with this much data, it is often easy to miss the big picture.

For this reason, there is an increasing need for techniques to present data in a meaningful and effective manner. In this paper, we investigate methods for automatically connecting the dots — providing a structured, easy way to navigate within the news domain: given two news articles, our system automatically finds a coherent chain of events starting with the decline of health-care debate (2007), and ending with the decline of home prices in 2007. We formalize the characteristics of a good chain and propose an efficient algorithm (with theoretical guarantees) to find an efficient algorithm (with theoretical guarantees) to connect two fixed endpoints. We incorporate user feedback into our framework, allowing the algorithm to learn from personalized data. Our user studies demonstrate the effectiveness in helping users understanding the news.

1.3.07 Home Prices Fall Just a Bit
3.4.07 Keeping Borrowers Afloat
3.5.07 A Mortgage Crisis Begins to Spiral...
8.10.07 Investors Grow Wary of Bank's Reliance on Debt
9.26.08 Markets Can't Wait for Congress to Act
10.4.08 Bailout Plan Wins Approval
1.26.09 Obama's Bailout Plan Moving Forward
9.1.09 ... and its effect on health benefits
9.1.09 Do Bank Bailouts Hurt Obama on Health?
9.22.09 (Bailout handling was undermine health-care reform)
9.22.09 Yes to Health-Care Reform, but Is This the Right Plan?

The chain mentions some of the key events connecting the mortgage crisis to healthcare, including the bailout plan. Most importantly, the chain should be coherent: after reading it, the user should gain a better understanding of the progression of the story. To the best of our knowledge, the problem of connecting the dots is novel. Previous research (e.g., [19, 13, 18, 17, 4, 6]) focused on organizing news articles into hierarchies or graphs, but did not address the notion of output coherence. Our main contributions are:

- Formalizing characteristics of a good story and the notion of coherence.
- Formalizing influence with no link structure.
- Providing an efficient algorithm for connecting two fixed endpoints while maximizing chain coherence (with theoretical guarantees).

Metro Maps of Science

Dafna Shahaf
Carnegie Mellon University
5000 Forbes Avenue
Pittsburgh, PA
dshahaf@cs.cmu.edu

Carlos Guestrin
Carnegie Mellon University
5000 Forbes Avenue
Pittsburgh, PA
guestrin@cs.cmu.edu

Eric Horvitz
Microsoft Research
One Microsoft Way
Redmond, WA
horvitz@microsoft.com

ABSTRACT

As the number of scientific publications soars, even the most enthusiastic reader can have trouble staying on top of the evolving literature. It is easy to focus on a narrow aspect of one's field and lose track of the big picture. Information overload is indeed a major challenge for scientists today, and is especially daunting for new investigators attempting to master a discipline and scientists who seek to cross disciplinary borders. In this paper, we propose metrics of influence, coverage, and connectivity for scientific literature. We use these metrics to create structured summaries of information, which we call *metro maps*. Most importantly, metro maps explicitly show the relations between papers in a way which captures developments in the field. Pilot user studies demonstrate that our method can help researchers acquire new knowledge efficiently: map users achieved better precision and recall scores and found more seminal papers while performing fewer searches.

Categories and Subject Descriptors

H.3.1 [Information Storage and Retrieval]: Content Analysis and Indexing; H.3.3 [Information Storage and Retrieval]: Information Search and Retrieval; H.5 [Information Interfaces and Presentation]

Keywords

Metro maps, Information, Summarization

1. INTRODUCTION

"Distringit liberum multitudine" (the abundance of books is a distraction), said Lucius Annaeus Seneca; he lived in the first century.

A lot has changed since the first century, but Lucius' problem has only become worse. The surge of the Web brought down the barriers of distribution, and the scientific community finds itself overwhelmed by the increasing numbers of publications; relevant data is often buried in an avalanche of publications, and locating it is difficult.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.
KDD'12, August 12–16, 2012, Beijing, China.
Copyright 2012 ACM 978-1-4503-1462-6/1208...\$10.00.

Trains of Thought: Generating Information Maps

April 16–20, 2012, Lyon, France

Dafna Shahaf
Carnegie Mellon University
5000 Forbes Avenue
Pittsburgh, PA
dshahaf@cs.cmu.edu

Carlos Guestrin
Carnegie Mellon University
5000 Forbes Avenue
Pittsburgh, PA
guestrin@cs.cmu.edu

Eric Horvitz
Microsoft Research
One Microsoft Way
Redmond, WA
horvitz@microsoft.com

ABSTRACT

When information is abundant, it becomes increasingly difficult to fit nuggets of knowledge into a single coherent picture. Complex stories spaghetti into branches, side stories, and intertwining narratives. In order to explore these stories, one needs a map to navigate unfamiliar territory. We propose a methodology for creating structured summaries of information, which we call *metro maps*. Our proposed algorithm generates a concise structured set of documents which maximizes coverage of salient pieces of information. Most importantly, metro maps explicitly show the relations among retrieved pieces in a way that captures story development. We first formalize characteristics of good maps and formulate their construction as an optimization problem. Then we provide efficient methods with theoretical guarantees for generating maps. Finally, we integrate user interaction into our framework, allowing users to alter the maps to better reflect their interests. Pilot user studies with a real-world dataset demonstrate that the method is able to produce maps which help users acquire knowledge efficiently.

Categories and Subject Descriptors

H.3.1 [Information Storage and Retrieval]: Content Analysis and Indexing; H.3.3 [Information Storage and Retrieval]: Information Search and Retrieval; H.5 [Information Interfaces and Presentation]

Keywords

Metro maps, Information, Summarization

1. INTRODUCTION

As data becomes increasingly ubiquitous, users are often overwhelmed by the flood of information available to them. Although search engines are effective in retrieving nuggets of knowledge, the task of fitting those nuggets into a single coherent picture remains difficult.

We are interested in methods for building more comprehensive views that explicitly show the relations among retrieved nuggets. We believe that such methods can enable people to navigate new, complex topics and discover previously unknown links. We shall focus on the news domain; for example, the system described in this paper can be used by a person who wishes to understand the debt crisis in Europe and its implications.

Copyright is held by the International World Wide Web Conference Committee (IW3C2). Distribution of these papers is limited to classroom use, and personal use by others.
WWW 2012, April 16–20, 2012, Lyon, France.
ACM 978-1-4503-1229-5/1204.



Figure 1: Greek debt crisis: a simplified metro map

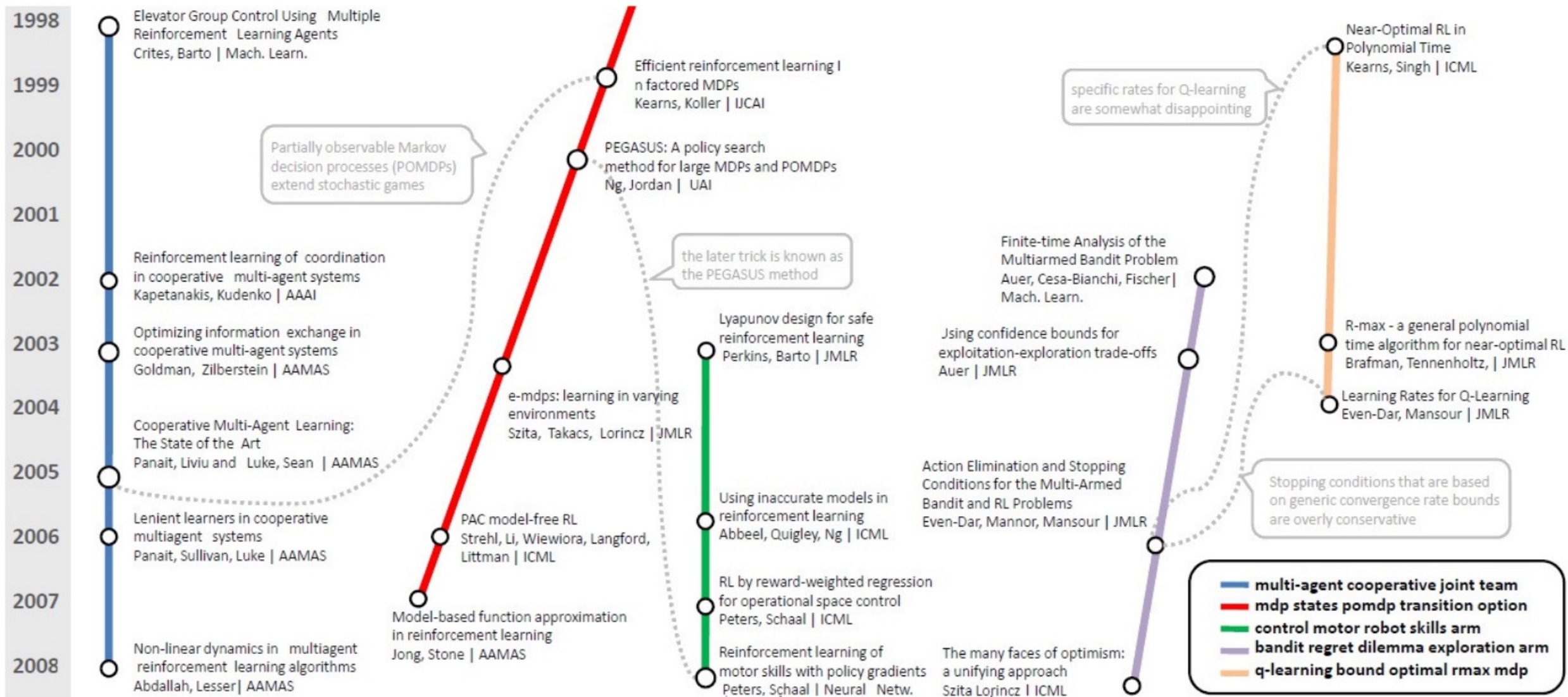


Figure 8: Part of the map computed for the query ‘Reinforcement Learning’. The map depicts multiple lines of research (see legend at the bottom). Interactions between the lines are depicted as dashed gray lines, and relevant citation text appears near them.

Why bring this up?

- DAWN project's version of this is for business
 - For them, explanation/discovery == business analytics
- Exercise: think about what you could do with these tools for general knowledge
- MetroMaps an example of general knowledge extraction in this context
- Similar in effort to other methods, e.g. topic models
 - Big DARPA project (FUSE)
 - DAWN also influenced by a big DARPA project
 - Discuss/reflect: funding <-> project choices

Questions?

Topics for Today

- ~~• Parisa's presentation~~
- ~~• Contextualizing Friday's papers~~
- Background for Friday's papers
- New method: ablation studies
- Things to reflect on for this week's reviews and presentations

Background

Reminder: ETL

Snorkle/DeepDive as an alternative to traditional ETL

Extract

Transform

Load

Often partially loaded
in DBs

The screenshot displays a patient's medical record in a web-based EHR system. The patient is identified as James Patient, 65 years old, male, with a primary care provider of Stephanie Provider. The interface includes a navigation sidebar with options like Home, Schedule, Tasks, Charts, Messages, and Reports. The main content area is divided into several sections: **Flowsheets** (Vitals), **Diagnoses** (Chronic and Acute), **Social history**, **Smoking status**, and **Past medical history**. The **Medications** section lists several drugs, with **Abilify 10 MG Oral Tablet** selected. A pop-up window titled 'Medication > Record medication' provides details for this medication, including its start date (05/04/2016), stop date, and associated diagnosis. Two red circles highlight the 'SIG' field (containing 'Take 2 tablets (20 mg) by mouth one time') and the 'MEDICATION COMMENT' field (containing 'Enter a comment about this medication').

Questions?

Generative vs. Discriminative Models

Suppose we want to predict class C from features X_1 , X_2 , and X_3 .

All classifiers learn $P(C \mid X_1, X_2, X_3)$

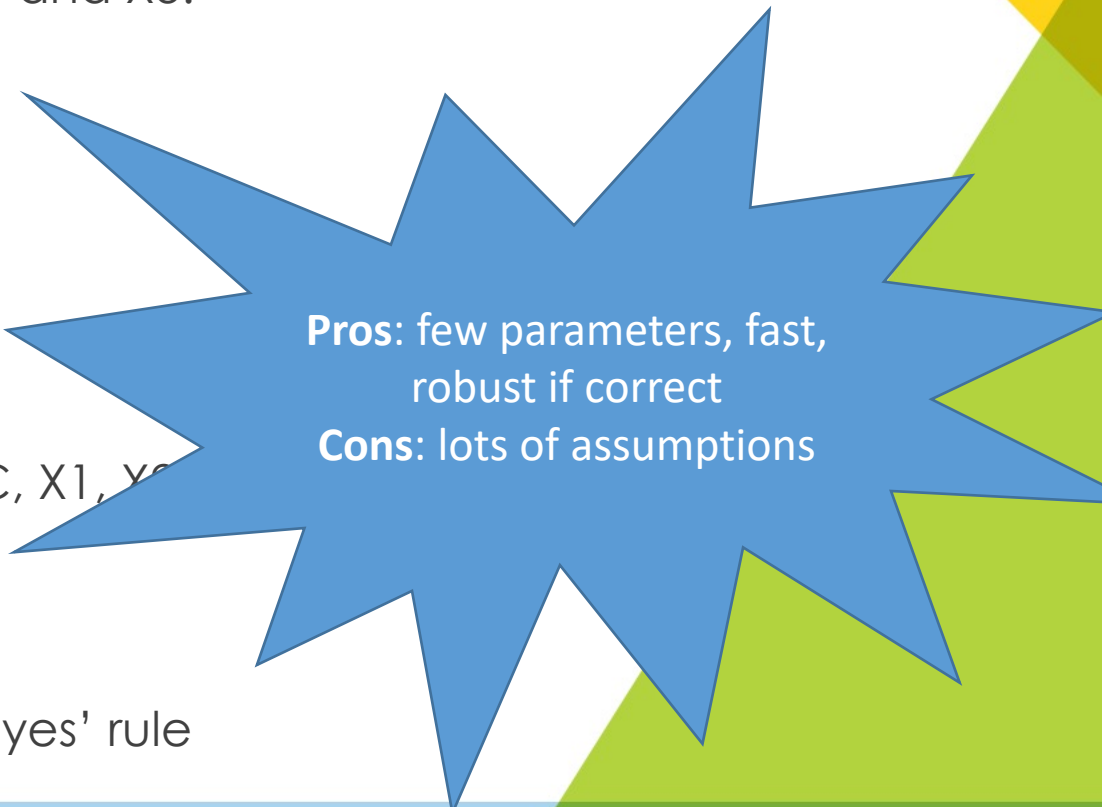
Generative:

Learn the “whole story”

A generative model learns the probability distribution $P(C, X_1, X_2, X_3)$

$P(C \mid X_1, X_2, X_3)$ can be derived

Requires labelled data to get $P(X_1, X_2, X_3 \mid C) \rightarrow$ Use Bayes' rule



Pros: few parameters, fast,
robust if correct
Cons: lots of assumptions

Generative vs. Discriminative Models

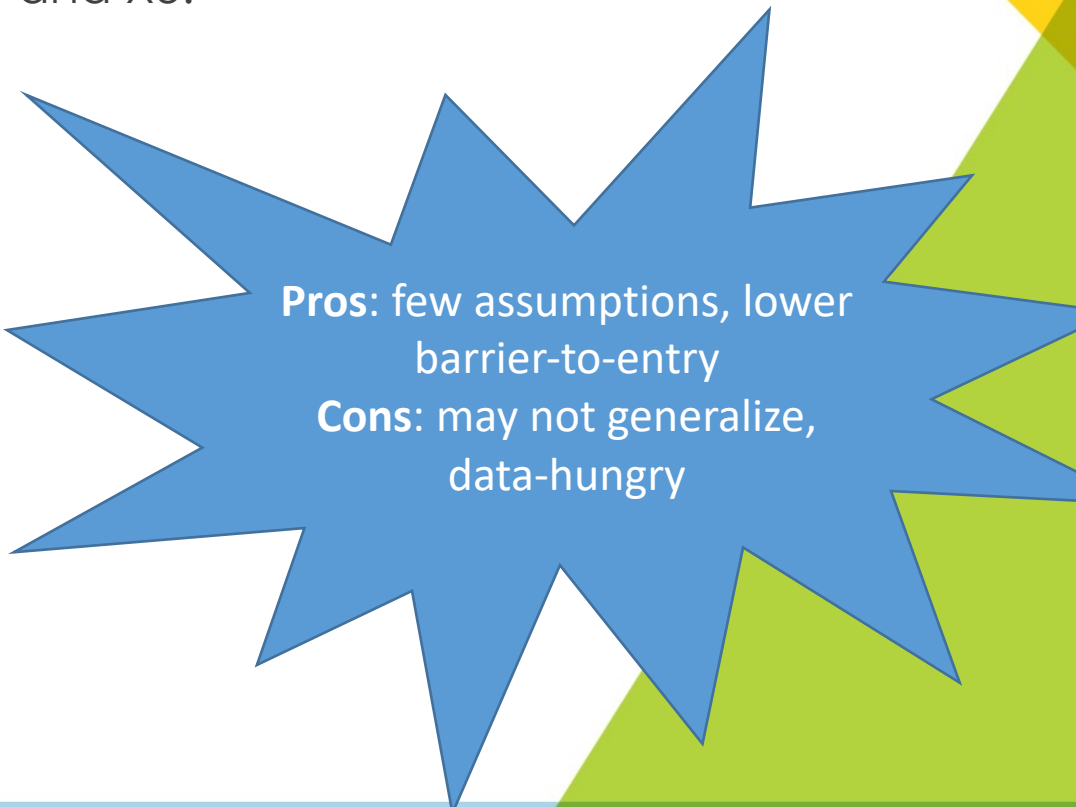
Suppose we want to predict class C from features X_1 , X_2 , and X_3 .

All classifiers learn $P(C \mid X_1, X_2, X_3)$

Discriminative:

Only learns what's necessary: $P(C \mid X_1, X_2, X_3)$

Requires labelled data to minimize error



Pros: few assumptions, lower barrier-to-entry
Cons: may not generalize, data-hungry

Examples of Generative vs. Discriminative Models

Generative	Discriminative
Naïve Bayes Classifier	Classic Deep Neural Networks
Hidden Markov Models	Regression
Latent Dirichlet Allocation (a kind of topic model)	Conditional Random Fields
Gaussian Mixture Models	Random Forests

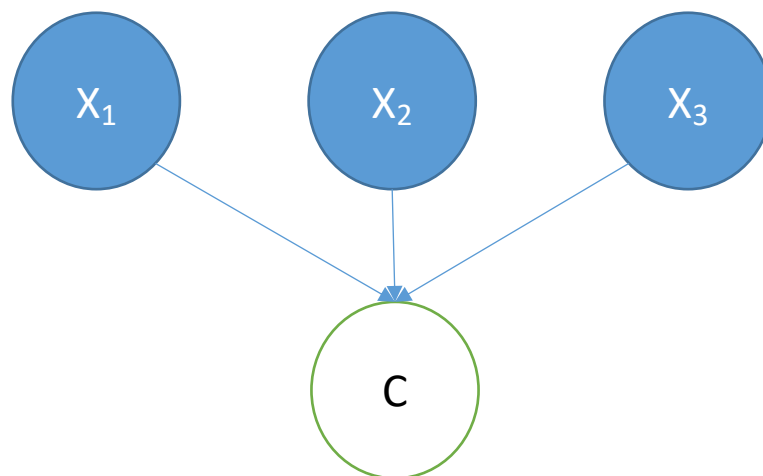


Pairing not semantic

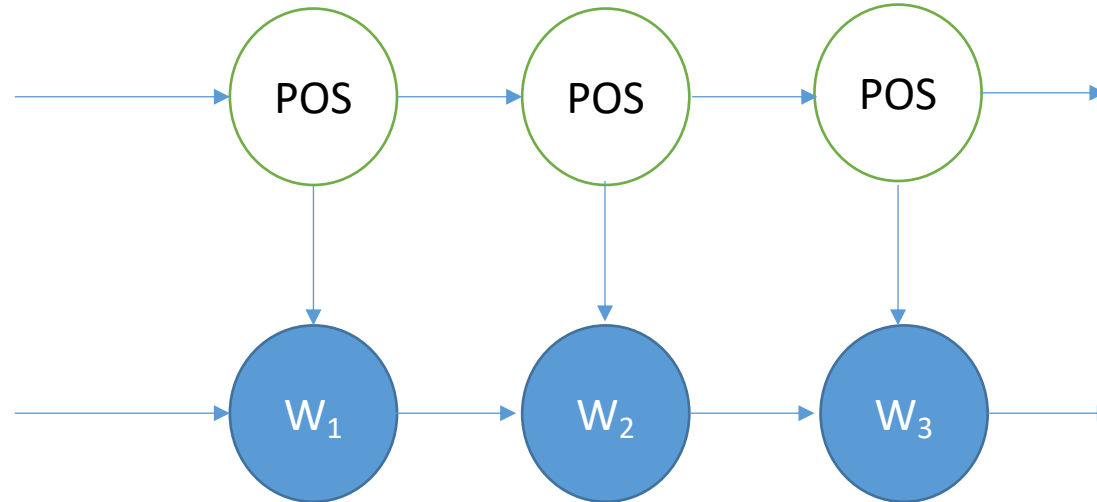
What is a latent variable?

Simple answer:

a variable whose values are not/cannot be observed



More interesting example



Questions?

Topics for Today

- ~~Parisa's presentation~~
- ~~Contextualizing Friday's papers~~
- ~~Background for Friday's papers~~
- New method: ablation studies
- Things to reflect on for this week's reviews and presentations

Ablation Studies

Ablation study: a method we haven't seen before!

We have not previously talked about ablation studies

Idea:

You've designed a solution that performs better. Yay!

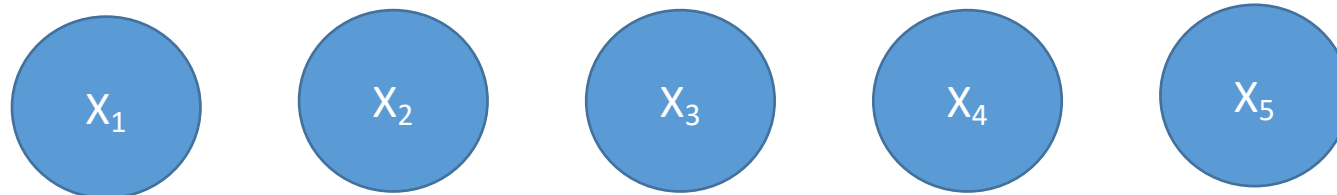
But that solution is *complex*. Which parts really contributed to the improved performance?

Run a kind of experiment (**control** is the complete solution; **treatment** is removing one part)

Very common in complex machine learning, especially deep learning.

Ablation study: what it can tell us

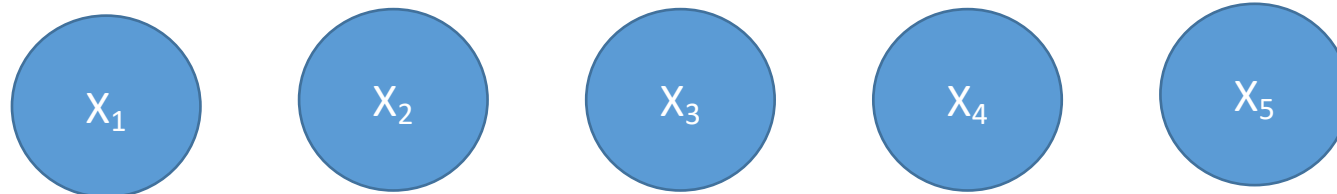
Only tests whether a feature **independently** does not contribute to the final performance



$$P_{\text{solution}} > P_{\text{SOTA}}$$

Ablation study: what it can tell us

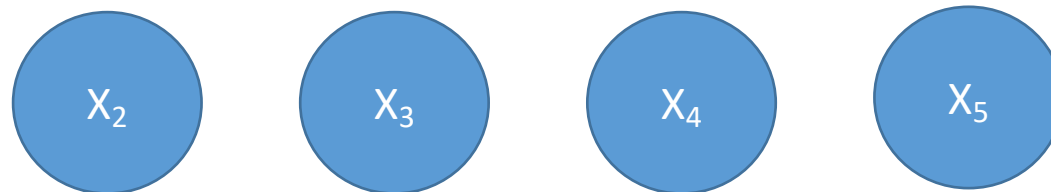
Only tests whether a feature **independently** does not contribute to the final performance



Remove or replace with
an appropriate
substitute

Ablation study: what it can tell us

Only tests whether a feature **independently** does not contribute to the final performance



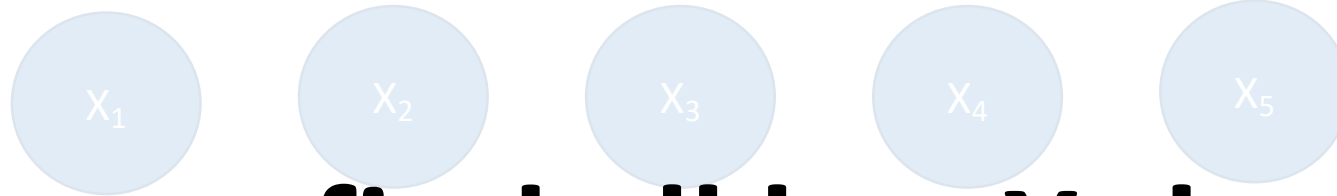
X_1 had no effect

$$P_{\text{solution}/x_1} \approx P_{\text{solution}}$$

Ablation study: what it can tell us

Only tests whether a feature **directly** contributes to the final performance

Continue for X_3, X_4, X_5



Assume you find all but X_1 had an effect.

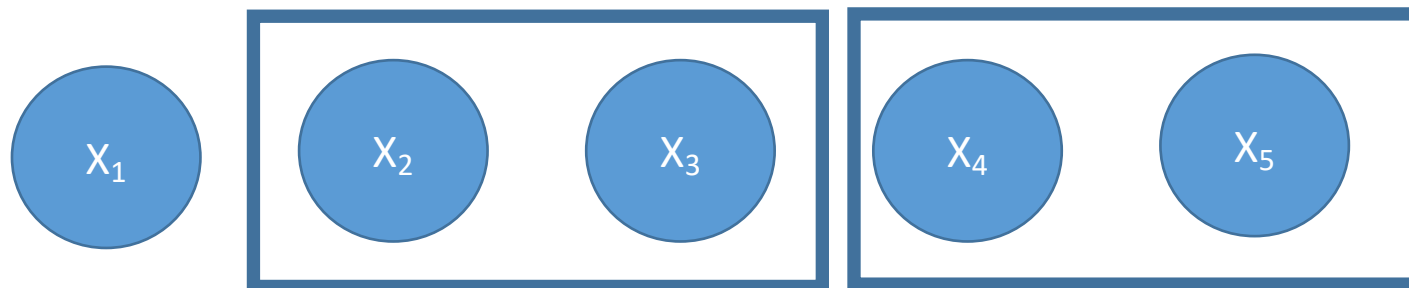


X_2 had an effect!

$$P_{\text{solution}/x_2} < P_{\text{solution}}$$

Ablation study: what it can tell us

Only tests whether a feature **independently** does not contribute to the final performance

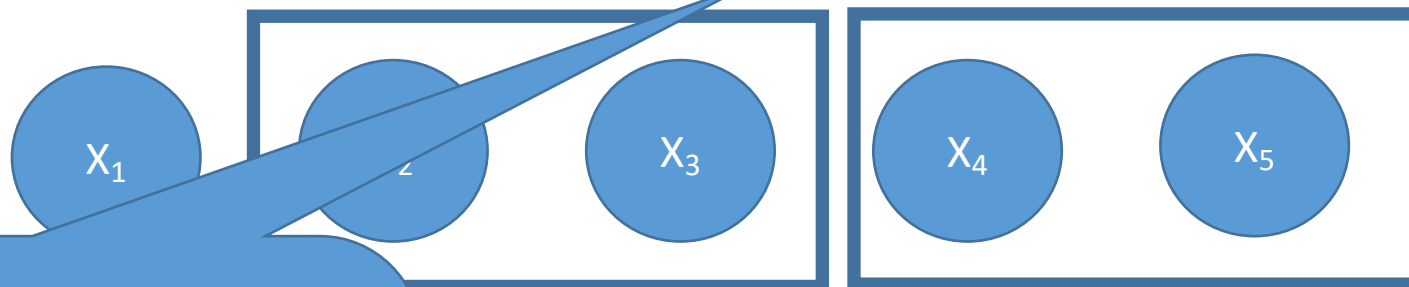


But what if certain parts/features only have an effect when they are together?

Can't detect this!

Ablation study: what it can tell us

Only tests whether a feature **independently** does not contribute to the final performance



This is what we mean by
“independently does not
contribute”

in parts/features only have an effect when they
are together?

Can't detect this!

Ablation study: example

Domain: classifying Wikipedia pages as literary or not

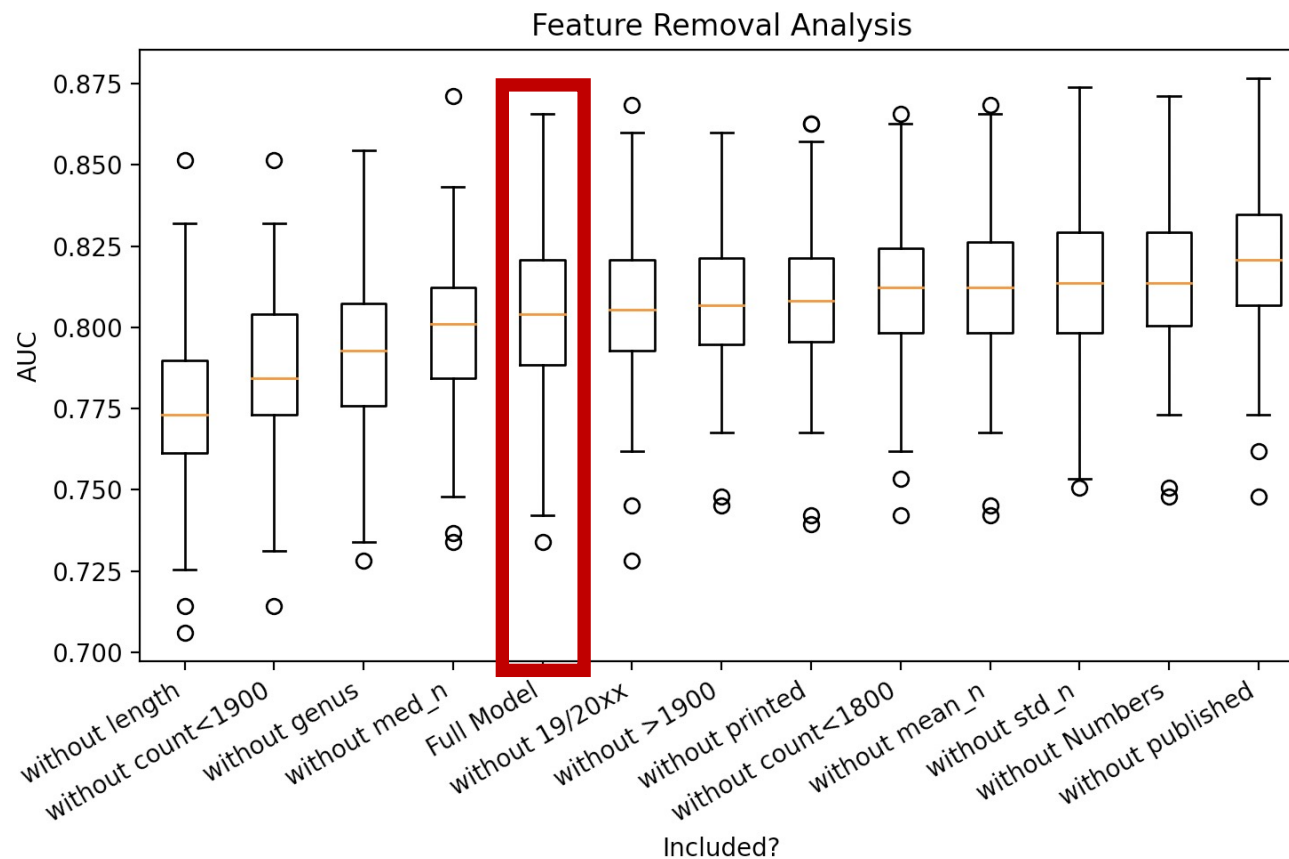
To left of Full Model:

important features that contribute to performance

To the right of Full Model:

Probably just randomness

Definitely possible to do better without!



Questions?

Topics for Today

- ~~Parisa's presentation~~
- ~~Contextualizing Friday's papers~~
- ~~Background for Friday's papers~~
- ~~New method: ablation studies~~
- Things to reflect on for this week's reviews and presentations

Things to think about

As you read...

DeepDive vs. Snorkle

- What's new?
- Which was published first?
- Following the citations: how many papers are by the same authors?

What is the ethical obligation of the annotator? Do you think human-in-the-loop solutions are a viable antidote to uncontrolled machine learning applications?