

CS 124/LINGUIST 180

From Languages to Information

DAN JURAFSKY

PROFESSOR OF COMPUTER SCIENCE

PROFESSOR OF LINGUISTICS

STANFORD UNIVERSITY

WINTER 2025

INTRODUCTION AND COURSE OVERVIEW

What is this class?

Interacting with humans via language

- Answering questions
- Searching the web
- Recommending things

And extracting meaning and structure from:

- Human language (news, social media, websites, etc.)
- Social networks

What is this class?

The very broad undergrad intro to (at least) 12 grad classes!

cs224C: NLP for Computational Social Science (Yang)

cs224N: Natural Language Processing with Deep Learning (Hashimoto/Yang)

cs224U: Natural Language Understanding (Potts)

cs224V: Conversational Virtual Assistants with Deep Learning (Lam)

cs224S: Spoken Language Processing (Maas)

cs246: Mining Massive Data Sets (Leskovec)

cs224W: Graph Neural Networks (Leskovec)

cs276: Information Retrieval (Manning)

cs329R: Race and Natural Language Processing (Jurafsky/Eberhardt)

cs329X: Human-Centered LLMs (Yang)

cs336: Language modeling from scratch (Hashimoto/Liang)

cs384: Social and Ethical Issues in NLP (Jurafsky)

What is this class?

The rise of LLMs has completely changed everything in

- Natural Language Processing (NLP)
- AI
- Information Retrieval (IR)
- Recommendation Systems

This class starts from scratch and builds up how LLMs work and how they are applied!

What is this class?

Intro to the algorithmic components of LLMs

Logistic Regression

Programming
Assignment 3

Word embeddings

Programming
Assignment 5

Neural Networks

Programming
Assignment 6

Gradient Descent/Backprop

Cross-entropy loss

Quiz 7

Transformers

Language Modeling Loss

Sampling and Temperature

Large language models!

What can LLMs do?
What can't they do?

Chat with Gemini



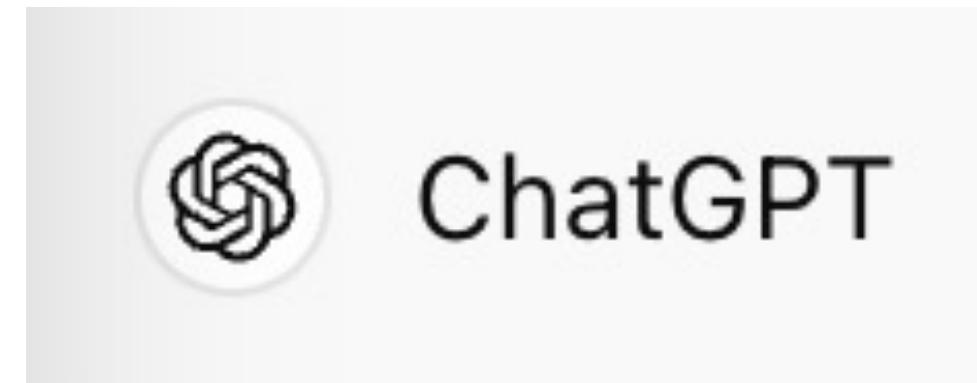
Claude



Llama 3

Lab 4
PA 7

LLMs and Personal Assistants



Message ChatGPT



Listening..



PA 7 Chatbot!



What is this class?

Intro to crucial algorithms that don't make up LLMs

Regular Expressions

Minimum Edit Distance

Supervised classifiers

- Naive Bayes
- Logistic Regression
- Neural Networks

Linguistic tools

- Sentiment/Emotion lexicons

Information Retrieval/
WebSearch

Network algorithms

- PageRank & Centrality
- Power Laws & Clustering

Recommendation engines

- Collaborative filtering

Example Topic: Information Retrieval

6,586,013,574 web searches every day (by one estimate)

Text-based information retrieval is one of the most frequently used algorithms in the world

How does it work? Can you build an IR engine?

Programming Assignment 3: Search!

Computational Biology: Comparing Sequences

AGGCTATCACCTGACCTCCAGGCCGATGCC

TAGCTATCACGACCGCGGGTCGATTGCCCGAC

-AGGCTATCAC_{CT}GACCTCC_AGGCGA--TGCCC---

TAG-CTATCAC--GACCGC--GGTCGA_{TT}TGCCCG

Sequence comparison is key to

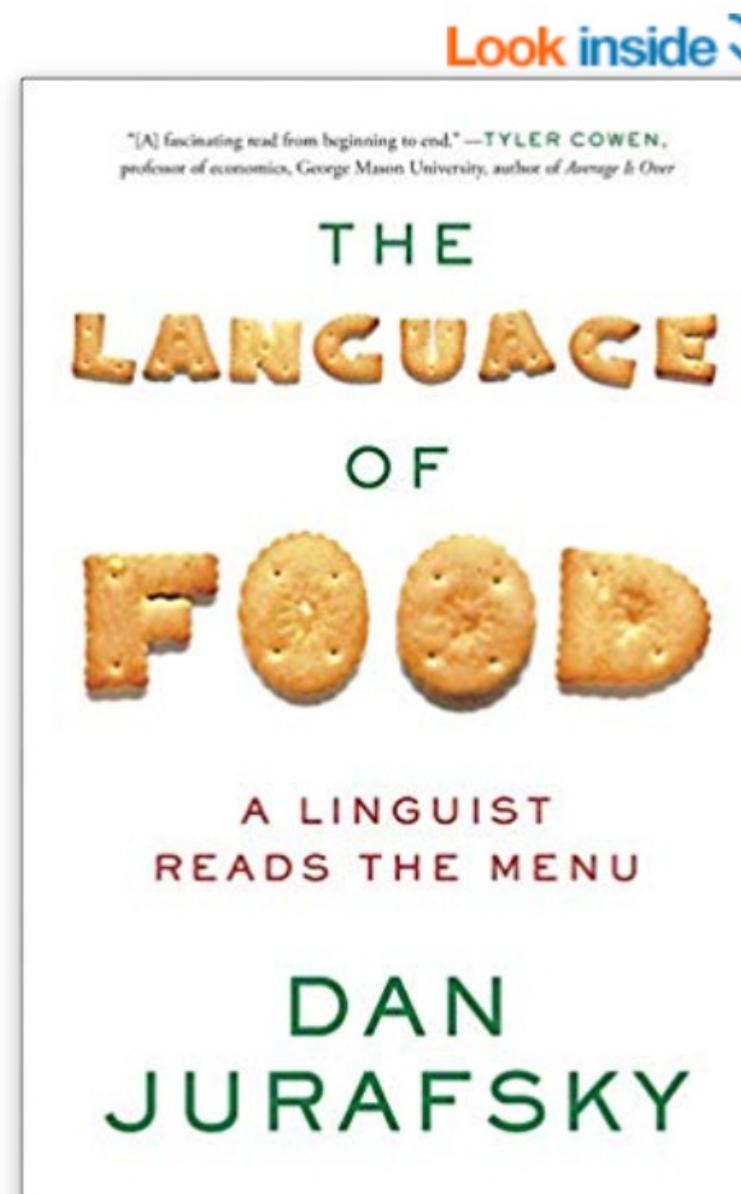
- Finding genes
 - Determining function
 - Uncovering evolutionary processes

This is also how simple spell checkers work!

Minimum edit distance (Quiz 1)

Recommendation Engines: The Good

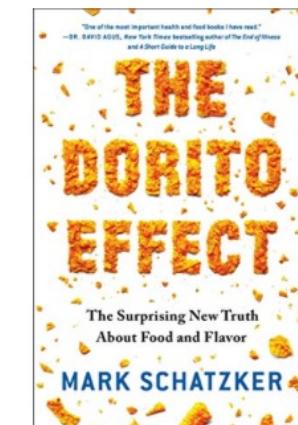
If you bought....



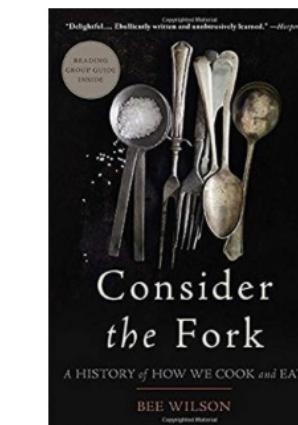
Customers who bought this item also bought



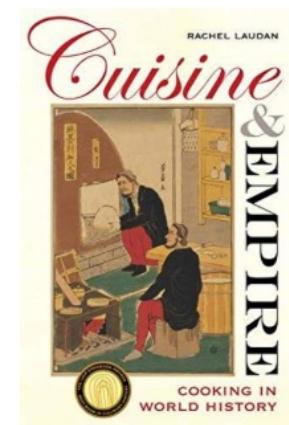
[First Bite: How We Learn to Eat](#)
by Bee Wilson
★★★★★ 46
Paperback
\$11.37 prime



[The Dorito Effect: The Surprising New Truth About Food and Flavor](#)
by Mark Schatzker
★★★★★ 193
Paperback
\$9.48 prime



[Consider the Fork: A History of How We Cook and Eat](#)
by Bee Wilson
★★★★★ 253
Paperback
\$15.65 prime



[Cuisine and Empire: Cooking in World History](#)
(California Studies in...
by Rachel Laudan
★★★★★ 35
Paperback
\$16.20 prime

PA 7 and Quiz 8

And the dark side: YouTube Radicalization



What is this class? Introduction to Social NLP

NLP and LLMs now interact with the **social world**.

We'll cover applications to social questions like

- Disaster Relief
- Helping teachers in the classroom
- Detecting latent meaning in political speech
- Analyzing language on policy body-worn camera to improve Police-Community relations

Example Topic: Text Classification

Disaster Response!

Haiti Earthquake 2010

Classifying SMS messages

Mwen thomassin 32 nan pyron
mwen ta renmen jwen yon ti dlo
gras a dieu bo lakay mwen anfom
se sel dlo nou bezwen

I am in Thomassin number 32, in the area named Pyron. I would like to have some water. Thank God we are fine, but we desperately need water.



*Programming
Assignment 2: Triage!*

Example Topic: Extracting Sentiment and Social Meaning

Lots of meaning is in **connotation**

"connotation: an idea or feeling that a word invokes in addition to its literal or primary meaning."

Extracting connotation is generally called
sentiment analysis

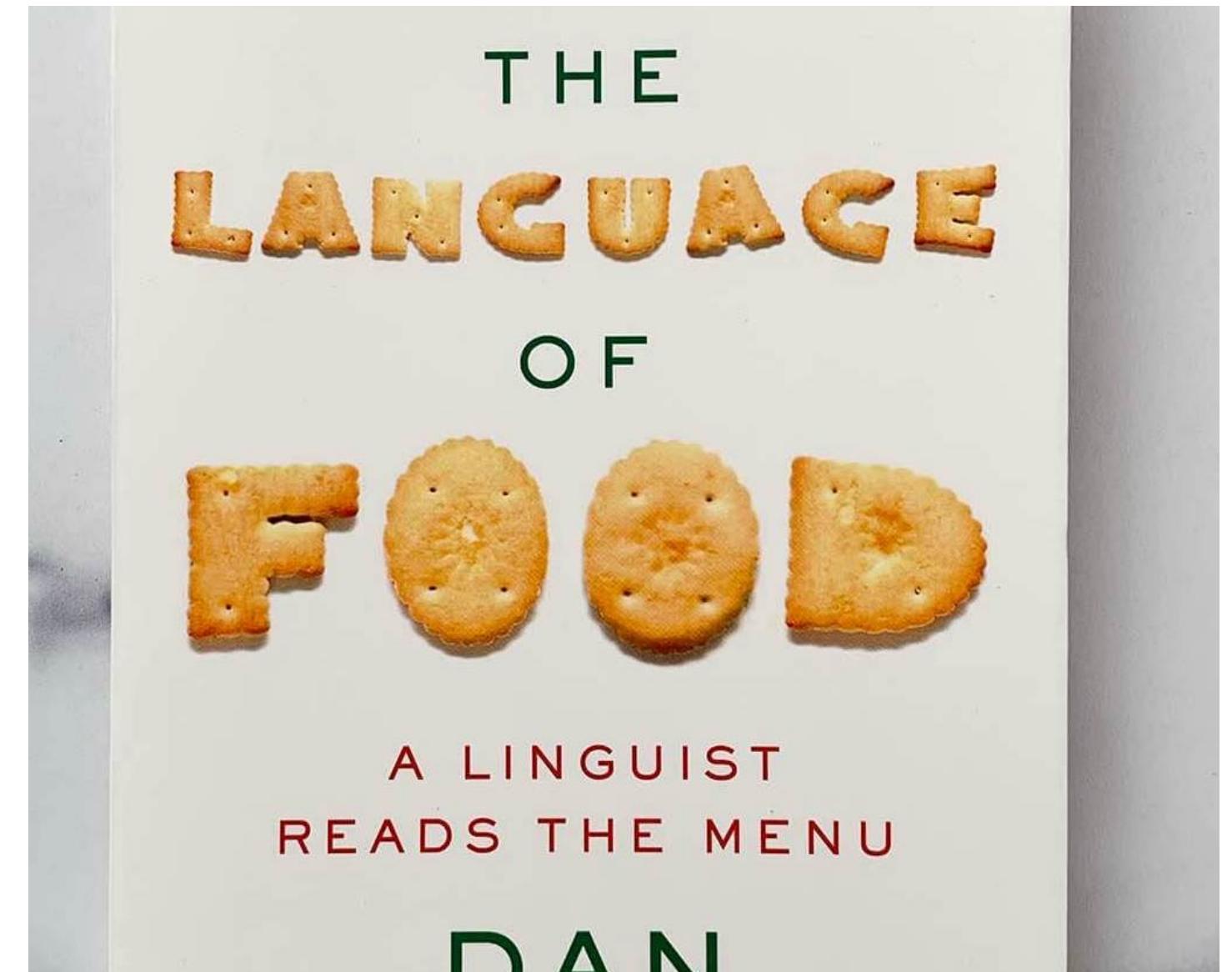
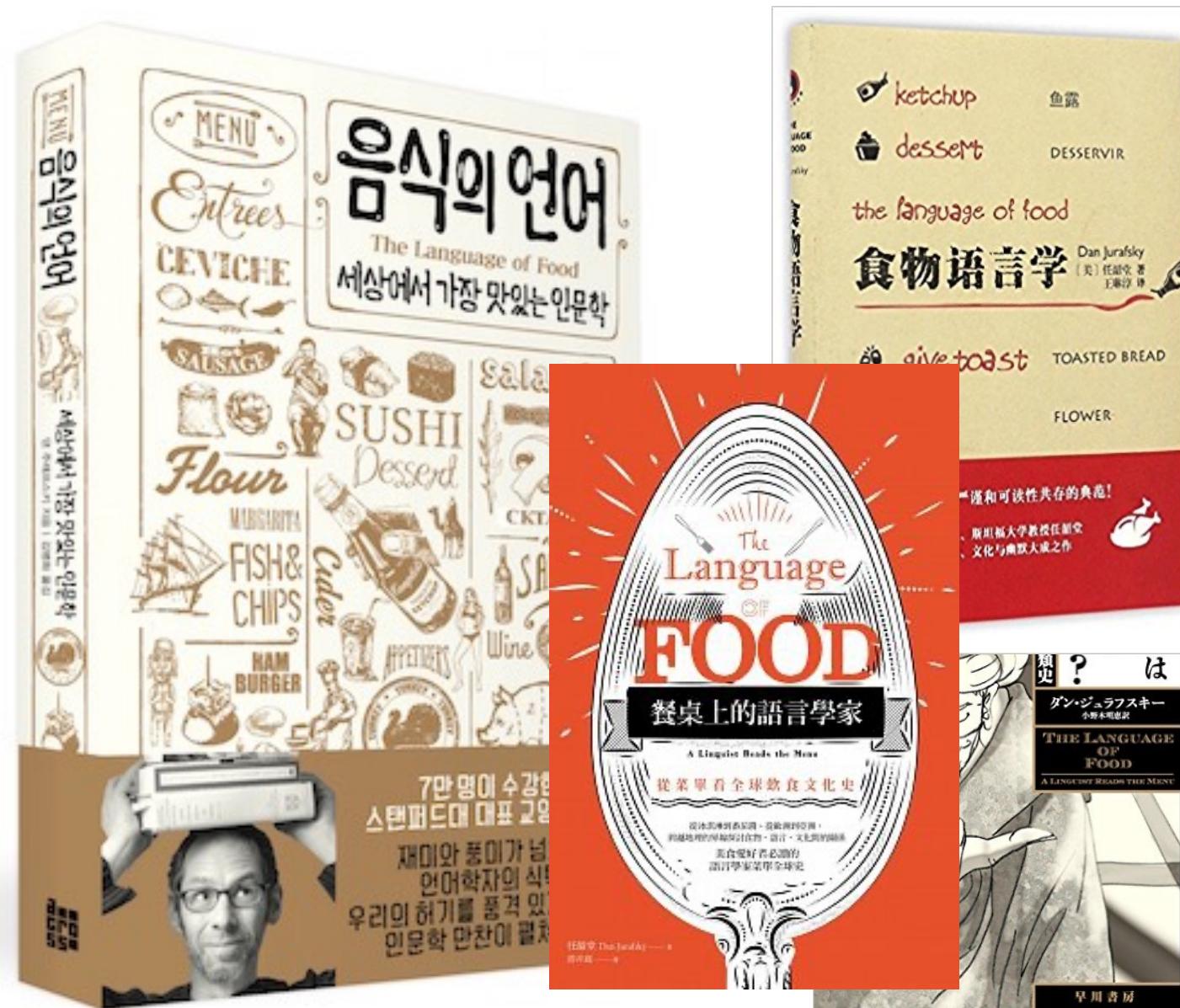
Programming Assignment 2: Sentiment

BTW, can we extract affective/social meaning with LLMs?

Mostly, but not completely!

- **Toxicity detection:** detecting hate speech, abuse, harassment, or other kinds of toxic language
 - Widely used in online **content moderation**
 - LLMs can find obvious toxicity (at least in English)
 - Problem: toxicity classifiers misfire on words for minority identities:
 - The word "blind": (Hutchinson et al., 2020)
 - The word "gay" (Dixon et al. 2018, Oliva et al., 2021)
 - Result: censorship of speech by disabled people and others

Example topic: Applying social NLP to humanities, social science, and cultural analytics!



Sentiment in Restaurant Reviews

Dan Jurafsky, Victor Chahuneau, Bryan R. Routledge, and Noah A. Smith. 2014. Narrative framing of consumer sentiment in online restaurant reviews. First Monday 19:4

900,000 Yelp reviews online

A very bad (one-star) review:

The bartender... absolutely horrible... we waited 10 min before we even got her attention... and then we had to wait 45 - FORTY FIVE! - minutes for our entrees... stalk the waitress to get the cheque... she didn't make eye contact or even break her stride to wait for a response ...

What is the language of bad reviews?

Negative sentiment language

horrible awful terrible bad disgusting

Past narratives about people

waited, didn't, was

he, she, his, her,

manager, customer, waitress, waiter

Frequent mentions of **we** and **us**

... **we** were ignored until **we** flagged down a waiter to get **our** waitress ...

Other narratives with this language

A genre using:

Past tense, we/us, negative, people narratives

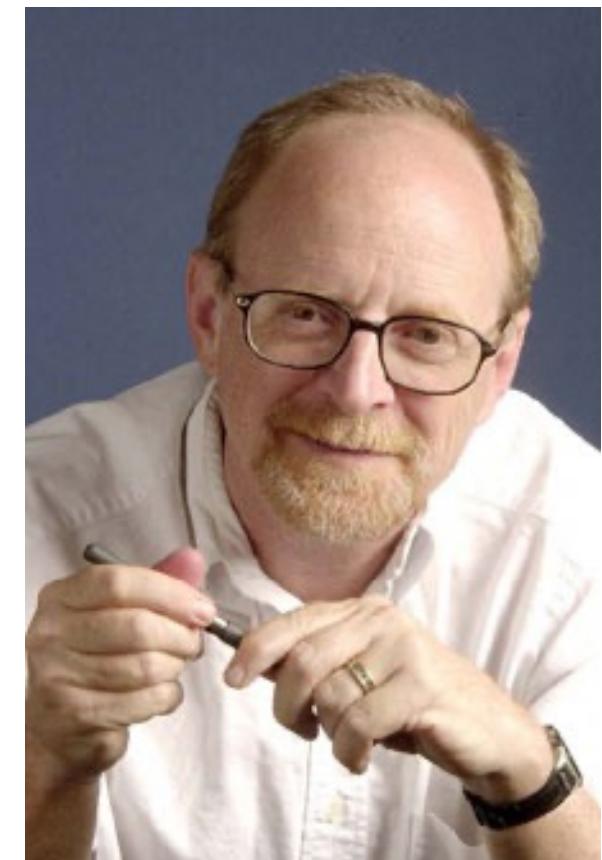
Texts written by **people suffering trauma**

- James Pennebaker lab at UT Austin
- Past tense is used for "distancing"
- Use of "we": seeking solace in community

1-star reviews are trauma narratives!

The lesson of reviews:

It's all about personal interaction



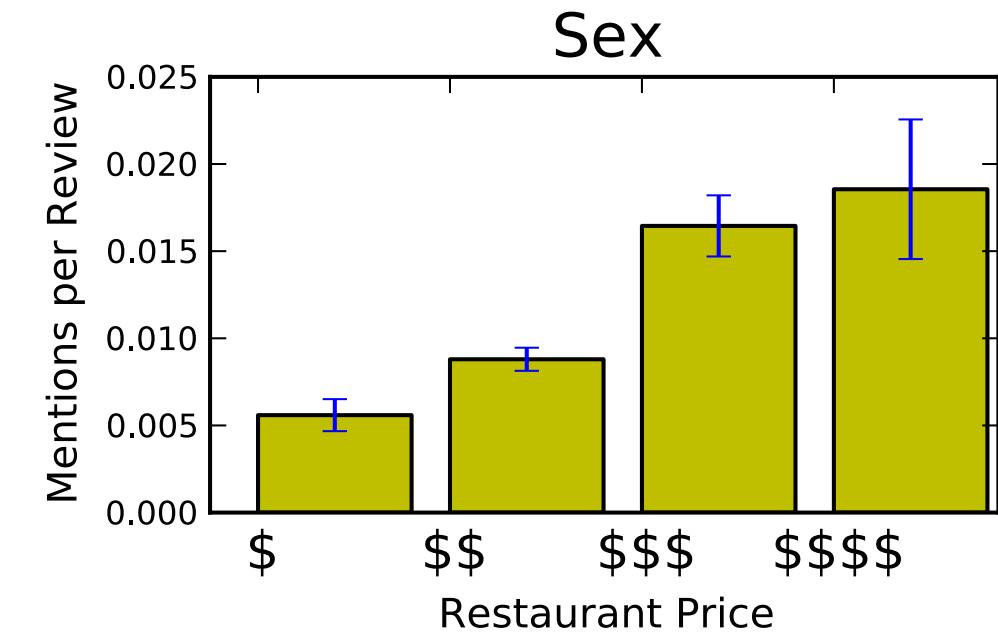
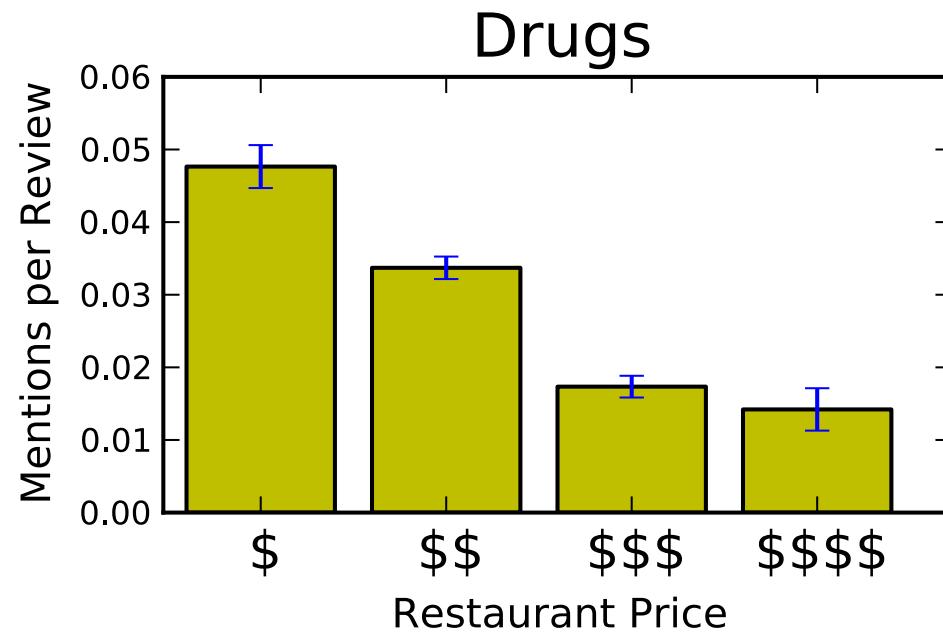
What about positive reviews? Sex, Drugs, and Dessert

addicted to pepper shooters

garlic noodles... my drug of choice

the fries are like crack

orgasmic pastry
sexy food
seductively seared fois gras



Help improve Police-Community Interaction (week 6)

Problems:

- A flood of videos show inappropriate use of force by police
- Black Americans suffer from much more negative interactions with police



Could natural language processing help?

- Quantify police-community interactions using body-worn cameras?
- Help develop officer training?
- Reduce the chances of violence?
- I'll talk about my work with Prof. Jennifer Eberhardt



Yet another topic: Social Networks

The network formed by your friends or other relations offline or online

- Can we compute properties of these networks?
- Extract information from them?
 - *Network algorithms (Quiz 9)*

What is this class? The Commercial World



OpenAI



YouTube

Google

∞ Meta

Microsoft[®]



amazon

Let's think about the language modeling task.

Why is it so remarkable?

What makes language interpretation hard?

Ambiguity

Language is ambiguous

Often as language users we don't even notice this

Resolving ambiguity is hard

Some very simple kinds of ambiguity

There are at least half a dozen meanings of this sentence:

The chef made her duck

Go here and type (and vote for) some definitions

<https://pollev.com/danjurafsky451>



Ambiguity

create
the chef
someone else
cook identify
The chef made her duck waterfowl
lower

The cook cooked waterfowl for a different woman X (person using "she/her" pronouns) to eat

The cook cooked waterfowl belonging to X

The cook cooked waterfowl belonging to the cook

The cook created the (plaster?) waterfowl that X owns

The cook caused X to quickly lower X's head or body

The cook uncovered the true identity of the cook's spy waterfowl

The cook waved their magic wand and turned X into undifferentiated waterfowl

How we deal with these rich meanings in LLMs: Neural "word embeddings"

A word's meaning in each sentence/context:
A point or region in 1000-dimensional space! In 2D:



LLM representation are rich but also reflect human biases!

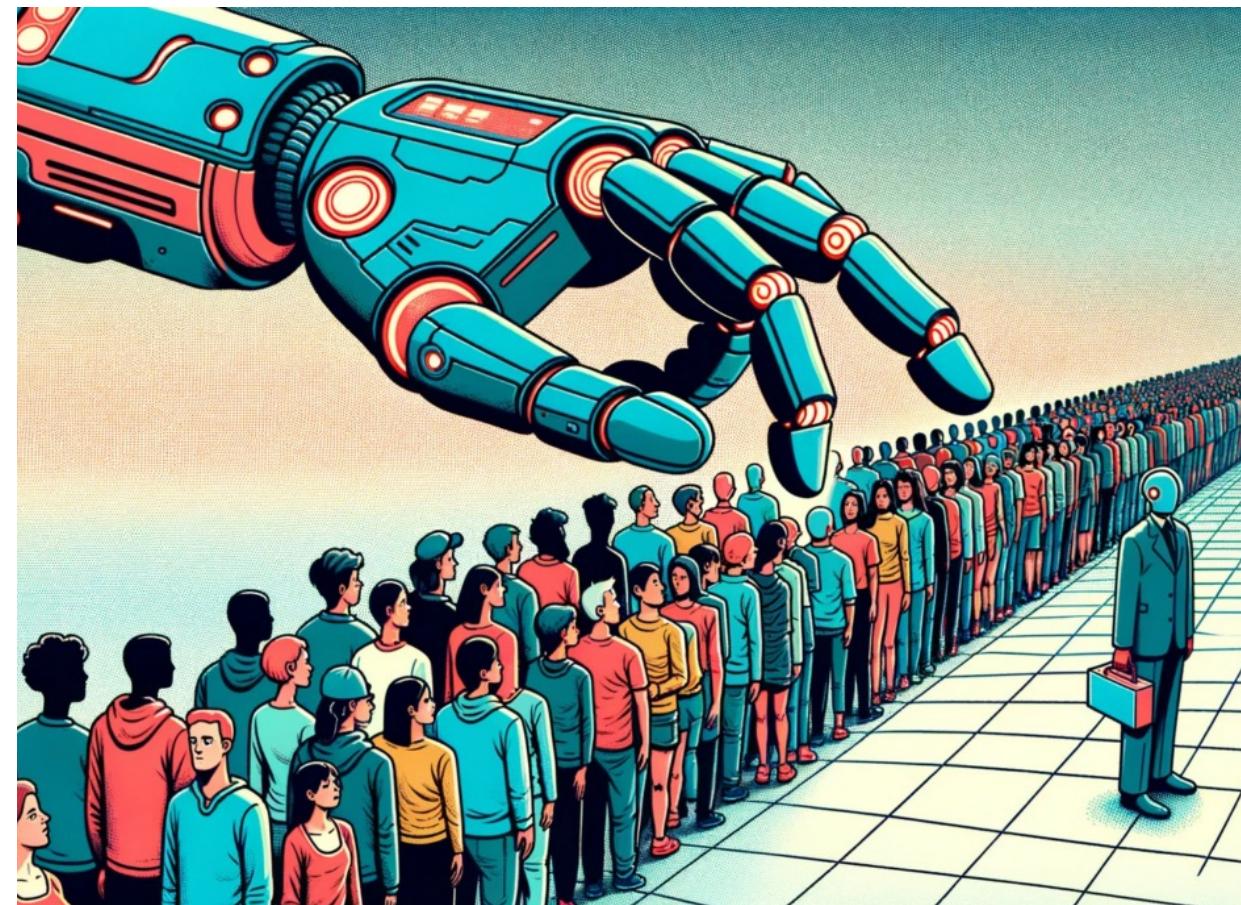
LLMs display stereotypes about pretty much every group (Asians, Blacks, Muslims, women, etc)

LLM internal representations also show these biases

These representations lead LLMs to also take biased actions

These biases have been clear from the very earliest papers studying LLM representations:

- Bolukbasi, Tolga, Kai-Wei Chang, James Y. Zou, Venkatesh Saligrama, and Adam T. Kalai. 2016. "Man is to computer programmer as woman is to homemaker? Debiasing word embeddings." In *NeurIPS 2016*, pp. 4349-4357.



The Decoder, Matthias Bastian, created by Dall-E

What can we do about this problem? PA/Quiz 5!

What is this class?
Evidence Based Pedagogy!

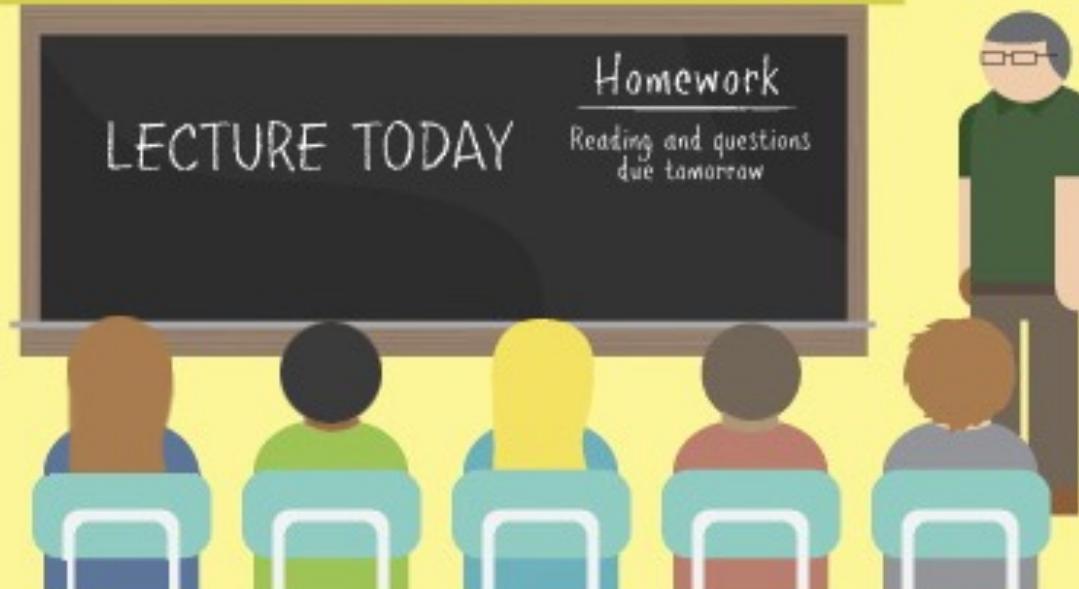
WHAT IS THE FLIPPED CLASSROOM?

The flipped classroom inverts traditional teaching methods, delivering instruction online outside of class and moving “homework” into the classroom.

THE INVERSION

The Traditional Classroom

Teacher's Role: Sage on the Stage



The Flipped Classroom

Teacher's Role: Guide on the Side



Why the flipped classroom (1)

Mastery learning: Learn until you master

Benjamin Bloom, 1968



Bloom's mastery learning

Personalized, **goal-driven practice**, driven by **feedback**

1. Watch (and re-watch) lectures at your own pace and learn when it's best for you
2. Videos have embedded miniquizzes. If you get it wrong, it gives you feedback about why you misunderstood.
3. You have **infinite** chances at each weekly Tuesday Quiz, so you can go back to the lecture and retake them.
4. With programming assignments you can see your performance on the training and dev set to see what you might be doing wrong on the test set!

Why the videos have embedded quizzes: “summative” vs “formative” assessment

Summative assessment

- Final exams/midterms: goal is grading

Formative assessment

- Along the way: goal is for **you** to find out what you don’t know so you can learn

Why I don't have a midterm or final

Multiple-choice timed tests don't reflect real life tasks

Scores don't correlate well with ability to do the task

They are stressful and annoying

They invite cheating

They waste an entire week that we instead use for content

Why the flipped classroom (2)

Attention span: everyone spaces out during long lectures

- Middendorf and Kalish, 1995, Johnstone and Percival 1976, Burns 1985

“the class started 1:00. The student sitting in front of me took copious notes until 1:20. Then he just nodded off... motionless, with eyes shut for about a minute and a half, pen still poised. Then he awoke and continued his rapid note-taking as if he hadn’t missed a beat.”

Student remembered only the first 15-20 minutes

Why the flipped classroom (3)

Active learning: Be in charge of your learning

- Most important: programming assignments
- Active learning (“constructivism”), learning by doing

Collaborative learning: Learn from each other

- Use class **lab** time for group problem-solving
- “Small group active learning”
- You must do **PA7** in groups of 3-4
- We encourage pair programming on PA1-6 and quizzes

cs124: Flipped classroom

1. Prerecorded video lectures on Canvas:

- About 80 ~10-minute lectures by me
- About ~90 minutes/week of video lectures
- Another 10 lectures by the TAs

2. Live sessions: (none are recorded)

- 5 required in-person lectures
- 5 required in-class labs (“active learning”)
 - Lab #1 (Unix text tools) next Tuesday is required in person
 - Labs #2, #3, #4 are required but attendance is extra credit (you can do at home).
 - Lab #5 March 3 (Git and PA7) is required in person

Logistics More Specifically

Online Video **Lectures** w/embedded non-graded questions (watch **before** relevant class/lab/quiz)

20 pages of **reading** a week (up to you when to read)

Weekly online **quizzes** (due Tue of following week)

7 Python programming assignments (PAs) (due Fri of following week)

- Except PA 7 you get extra time, 2+ weeks

Learning Goals

At the end of this course, you will be able to:

Learning goals

Write efficient regular expressions to solve any kind of text-based extraction task

Learning goals

Apply the edit distance algorithm to all sorts of text sequence problems

Learning goals

Build a supervised classifier to solve problems like sentiment classification

Learning goals

Build a neural network and train it using stochastic gradient descent

Learning goals

Build a search engine

Learning goals

Build a recommendation engine

Learning goals

Build a computational model of word meaning
(using lexicons and neural word embeddings)

Learning goals

Build a chatbot, both by prompting of a large language model, and by building dedicated components

Learning goals

Understand and implement PageRank and other social network functions

Learning goals

Understand the internal components of large language models like the transformer

Learning goals

Be able to prompt large language models, reason about what they can do and about their social implications

Learning goals

Work in our field is rarely done alone!

Goal: Learn to work together on computational projects and use group tools like github

PA1-6: Pair programming is encouraged

PA7: Must be done in groups of 3-4

What is this class?

The very broad undergrad intro to (at least) 12 grad classes!

cs224C: NLP for Computational Social Science (Yang)

cs224N: Natural Language Processing with Deep Learning (Hashimoto/Yang)

cs224U: Natural Language Understanding (Potts)

cs224V: Conversational Virtual Assistants with Deep Learning (Lam)

cs224S: Spoken Language Processing (Maas)

cs246: Mining Massive Data Sets (Leskovec)

cs224W: Graph Neural Networks (Leskovec)

cs276: Information Retrieval (Manning)

cs329R: Race and Natural Language Processing (Jurafsky/Eberhardt)

cs329X: Human-Centered LLMs (Yang)

cs336: Language modeling from scratch (Hashimoto/Liang)

cs384: Social and Ethical Issues in NLP (Jurafsky)

Should I take 124 or 224N or something else?

CS124 is designed for sophomores or juniors

- It's gentle (I explain everything) and broad (covering many topics, not just NLP/LLMs but also recommendation engines, IR, social networks, social computing)
- Mastery learning, quizzes, programming assignments with starter code and scaffolding.
- No research project, but a fun chatbot final homework

CS224N is a deeper, laser focused, grad course

- They assume you are very familiar with ML; 1st homework jumps right into optimization
- More focus on systems/implementation/scaling, you code more advanced things

CS224N/U/V/S/W, 246, 336, 329R

- Learning via research: novel research projects as a large component

CS324X (Human Centered NLP), CS346 (Social and Ethical Issues in NLP) require 224N or 224U

CS224C: more applied focus, applying NLP to social science: (NLP for Computational Social Science)

(You should of course take all of them!!)

Logistics: Instructor

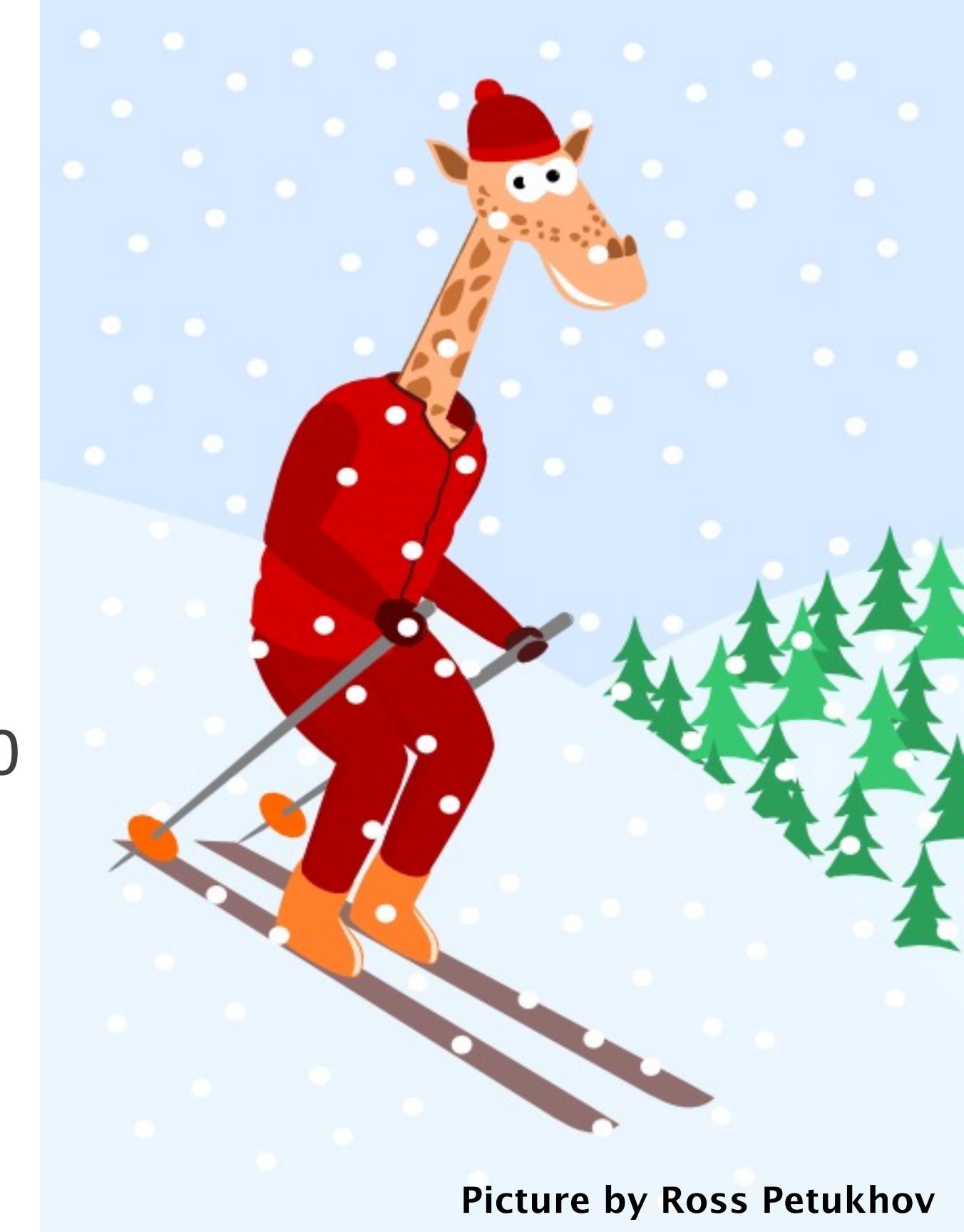
Instructor: Dan Jurafsky (he/him)

Professor in CS and Linguistics

My office hours:

- This week and next
 - Tuesday after class 4:30-5:40
- Then every Thursday classtime 3-4:20
- Margaret Jacks Hall 117
- Book times at calendly.com/jurafsky

How to pronounce my name:



Picture by Ross Petukhov

Course Staff



Dan Jurafsky
Professor



Priti Rangnekar
Head TA



Veronica Rivera
Ethics Postdoc



Xuheng Cai
TA



Adam Chun
TA



Gabriela Cortes Arias
TA



Kate Eselius
TA



Daniel Guo
TA



Sri Jaladi
TA



Jonathan Lee
TA



Kasey Luo
TA



Gabe Magaña
TA



Elena Recaldini
TA



Jeong Shin
Ethics TA



Savitha Srinivasan
TA



Rachel Yixing Wang
TA



Pannisy Zhao
TA

Syllabus

cs124.stanford.edu

Where do I find all the programming assignments and quizzes and readings?

Everything is on the webpage `cs124.stanford.edu`

Except the videos which are on Canvas Modules!

In other words:

- Lectures slides: webpage
- Lab instructions: webpage (points to git where they live)
- Tutorial information: webpage
- Readings: webpage
- Programming assignments: webpage (points to git where they live)
- Weekly quizzes: webpage (points to gradescope where they live)
- Videos: canvas

Coming up this week: Thursday

Optional tutorial on jupyter notebooks and PA0, getting ready for PA1

Come to class **with your laptops** and we'll go through PA0 together!

This tutorial will be led by amazing TA Sri Jaladi!!! But I and many other CAs will be there!

Action Items Before Thursdays class!

- 1) Read the syllabus webpage at cs124.stanford.edu
- 2) Look at PAO (you can find it from the webpage)
- 3) Watch Canvas Videos on "PAO Mac Setup" (or "PAO Windows Setup"), also pointed to by webpage

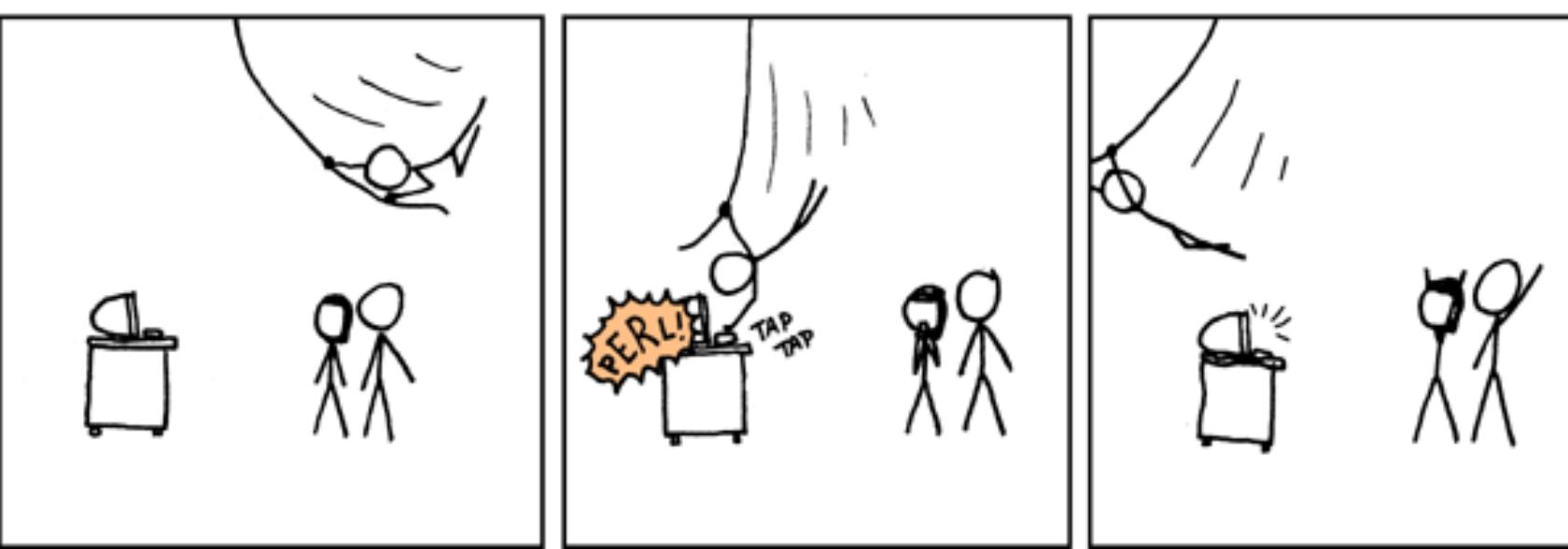
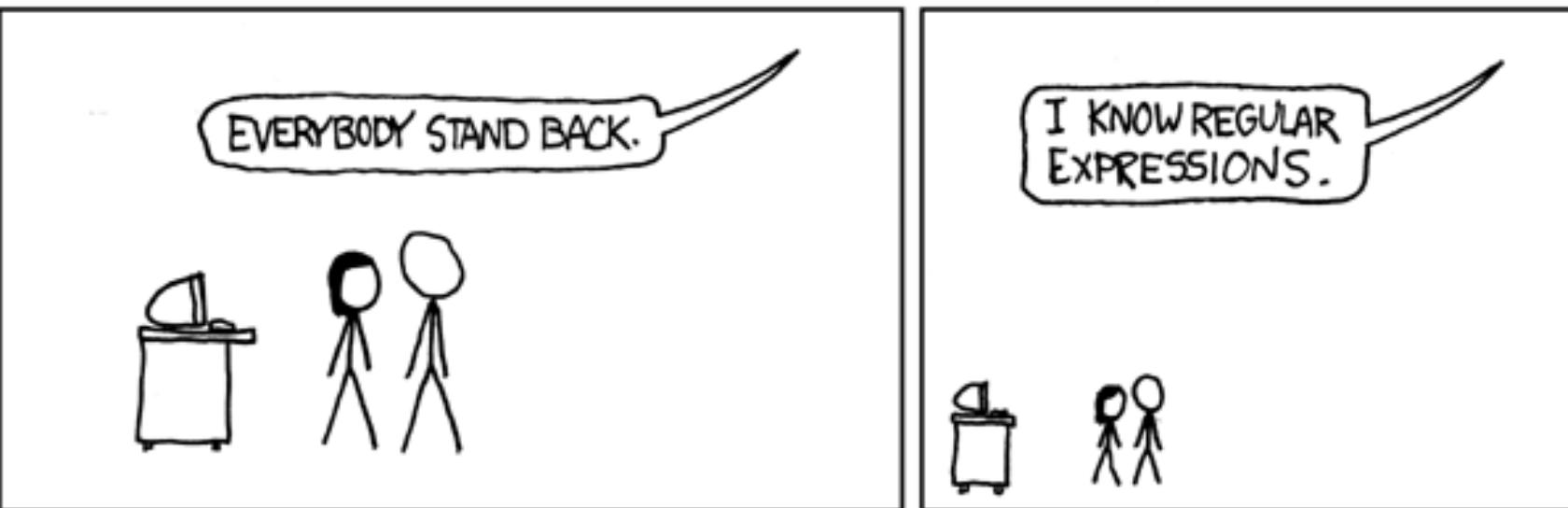
Coming up next week (Tuesday)

"Unix for poets":

grep

sort

Key UNIX tools for
dealing with text files and
regular expressions.



Action Items before next Tuesday's class!

1) Watch the "week 1" videos on Canvas by Sunday (since the quiz is also due Tuesday)

3) Download this file to your laptop

http://cs124.stanford.edu/nyt_200811.txt

4) If you don't know UNIX yet (haven't had cs107):

- For people using a Windows 10 machine, if you don't have Ubuntu on your machine:
 - Watch the pa0 Windows video about how to download and install Ubuntu (it's pointed to from the website)
 - Watch Chris Gregg's excellent UNIX videos here: Logging in, first 7 File System, and first 8 useful commands

<https://web.stanford.edu/class/archive/cs/cs107/cs107.1186/unixref/>

PA1: Spam Lord!

Write regular expressions to spread evil* throughout the galaxy!

By extracting email addresses and phone numbers from the web!

jur a fs ky at st anford dot e d u

Goes live Friday 5pm!

*Just kidding; don't be evil

YOU KNOW HOW SOMETIMES PEOPLE
PUT A SPACE IN THEIR EMAIL ADDRESS
TO MAKE IT HARDER TO HARVEST?

YEAH?

THEY HAVE A TOOL THAT
CAN DELETE THE SPACE!

OH MY GOD.



LESS-DRAMATIC REVELATIONS
FROM THE CIA HACKING DUMP