

# Natural Language Processing

AI51701/CSE71001

Lecture 2

08/31/2023

Instructor: Taehwan Kim

# Review of the Last Lecture

- ❑ Introduction to the course
- ❑ Examples of Natural Language Processing
- ❑ Word
  - Tokenization

# Lecture Plan

- ❑ Tips for Final Projects
- ❑ Word
  - Tokenization
  - Morphology
- ❑ Word representations

# Final Project

- Final Project: 40%
  - Proposal (5%): presentation (and written proposal) on Oct. 10 & Oct. 12, feedback will be given.
  - Final presentation (5%): on Dec.5 & Dec. 7
  - Final report (30%): by the end of last week of the semester (Dec. 15)
  - You will choose your project topics.
  - 3 people team
    - For team search, use BlackBoard discussion forum if needed.
  - You can use any language/framework for your project
    - Though we expect most of you to use Python (and recommend PyTorch for deep learning)

# Final Project

## ❑ What you would do:

- Defining a research goal
- Finding data and tools (models)
- Working out on the problem
- Evaluating your proposed solution/approach

## ❑ Must be human-language related.

- E.g., applying Transformers to a genetics dataset is **not** allowed unless related to human language pre-training dataset/tasks on downstream performance on the genetics task.

# Project Proposal

- Find a relevant (key) research paper for your topic
- Write a summary of that research paper and what you took away from it as key ideas that you hope to use
- Write what you plan to work on and how you can innovate in your final project work
- Describe as needed
  - A project plan, relevant existing literature, the kind(s) of models you will use/explore (usually deep learning models); the **data** you will use (and how it is obtained), and how you will **evaluate** success
  - Title, team members, and your planned contributions
- (Re-) implementing existing work is not enough.
- 3–4 pages (in Latex), presentation (and written proposal via BlackBoard) on Oct. 10 & Oct. 12.

# Project Proposal

- ❑ Need to think critically about a research paper
  - What were the main novel contributions or points?
  - Is what makes it work something general and reusable or a special case?
  - Are there flaws or neat details in what they did?
  - How does it fit with other papers on similar topics?
  - Does it provoke good questions on further or different things to try?

# Choosing Your Topic for Final Project

- ❑ Two basic starting points
  - Start with a (domain) problem of interest and try to find good/better ways to address it than are currently known/used
  - Start with a technical method/approach of interest, and work out good ways to extend it, improve it, understand it, or find new ways to apply it

# Choosing Your Topic for Final Project

- ❑ This is not an exhaustive list, but most projects are one of
  - Find an application/task of interest and explore how to approach/solve it effectively, often with an existing model
  - Implement a complex neural architecture and demonstrate its performance on some data
  - Come up with a new or variant neural network model or approach and explore its empirical success
  - Analysis project. Analyze the behavior of a model: how it represents linguistic knowledge or what kinds of phenomena it can handle or errors that it makes
  - Rare theoretical project: Show some interesting, non-trivial properties of a model type, data, or a data representation

# Choosing Your Topic for Final Project

- ❑ How to find an interesting place to start?
  - look at online proceedings of major NLP/ML conferences:
    - ACL anthology: <https://aclanthology.org/>
    - NeurIPS <https://papers.nips.cc>, ICML <https://proceedings.mlr.press/>,
    - ICLR <https://openreview.net/group?id=ICLR.cc>
  - online preprint servers
    - <https://arxiv.org>
  - look for an interesting problem in the world!
    - Hal Varian: How to Build an Economic Model in Your Spare Time  
<https://people.ischool.berkeley.edu/~hal/Papers/how.pdf>

# Finding data

- Some people collect their own data for a project
  - You may have a project that uses “unsupervised” data
  - You can annotate a small amount of data
  - You can find a website that effectively provides annotations, such as likes, stars, ratings, responses, etc.
  - But be careful on scoping things so that this doesn’t take most of your time
- **Most people make use of an existing, curated dataset built by previous researchers**
  - You get a fast start and there is obvious prior work and baselines

# Useful Resources

- ❑ Machine translation

- <http://statmt.org>
  - Look in particular at the various WMT shared tasks

- ❑ Dependency parsing: Universal Dependencies

- <https://universaldependencies.org>

- ❑ Huggingface Datasets

- <https://huggingface.co/datasets>

- ❑ Paperswithcode Datasets

- <https://www.paperswithcode.com/datasets?mod=texts&page=1>

# Useful Resources

- ❑ There are now many other datasets available online for all sorts of purposes
  - Look at Kaggle
  - Look at research papers to see what data they use
  - Look at lists of datasets
    - <https://machinelearningmastery.com/datasets-natural-language-processing/>
    - <https://github.com/niderhoff/nlp-datasets>
  - Lots of particular things:
    - <https://gluebenchmark.com/tasks>
    - <https://nlp.stanford.edu/sentiment/>
    - <https://research.fb.com/downloads/babi/> (Facebook bAbI-related)

# Useful Resources

## ❑ Sample projects

- Stanford 224n course final project reports

<https://web.stanford.edu/class/archive/cs/cs224n/cs224n.1214/project.html>

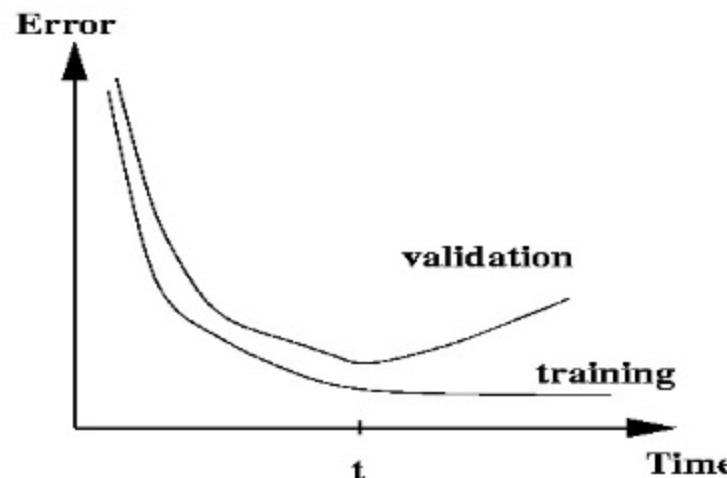
<https://web.stanford.edu/class/archive/cs/cs224n/cs224n.1204/project.html>

# Care with datasets and in model development

- ❑ Many publicly available datasets are released with a **train/dev/test** structure
- ❑ If there is no dev set or you want a separate tune set, then you create one by splitting the training data
- ❑ Having a fixed test set ensures that all systems are assessed against the same gold data.

# Training models and pots of data

- ❑ When training, models overfit to what you are training on
  - The model correctly describes what happened to occur in particular data you trained on, but the patterns are not general enough patterns to be likely to apply to new data
- ❑ The way to monitor and avoid problematic overfitting is using **independent validation** and test sets



From Stanford CS224n course

# Training models and pots of data

- ❑ You build (estimate/train) a model on a training set.
- ❑ Often, you then set further hyperparameters on another, independent set of data, the **tuning set**
- ❑ You measure progress as you go on a **dev set** (development test set or validation set)
  - If you do that a lot you overfit to the dev set so it can be good to have a second dev set, the **dev2** set
- ❑ **Only at the end**, you evaluate and present final numbers on a **test** set

# Training models and pots of data

- ❑ The train, tune, dev, and test sets need to be completely **distinct**
- ❑ It is invalid to give results testing on material you have trained on
  - You will get a falsely good performance.
  - We almost always overfit on train

# Experimental strategy

- ❑ Work incrementally!
- ❑ Start with a very simple model and get it to work!
  - It's hard to fix a complex but broken model
- ❑ Add bells and whistles one-by-one and get the model working with each of them (or abandon them)
- ❑ Initially run on a tiny amount of data
  - You will see bugs much more easily on a tiny dataset ... and they train really quickly
  - Something like 4–8 examples is good
  - Often synthetic data is useful for this
  - Make sure you can get 100% on this data

# Experimental strategy

- ❑ There are lots of things that can cause neural nets to not learn at all or to not learn very well
  - Finding and fixing them (“debugging and tuning”) can often take more time than implementing your model
  - But experience, experimental care, and rules of thumb help!
- ❑ Train and run your model on a large dataset
  - It should still score close to 100% on the training data after optimization.
    - Otherwise, you probably want to consider a more powerful model!
    - Overfitting to training data is **not** something to fear when doing deep learning
- ❑ You now want good generalization performance:
  - Regularize your model until it doesn’t overfit on dev data
  - Strategies like L2 regularization can be useful
  - But normally generous **dropout** is the secret to success

# Details matter!

- ❑ Look at your **data**, collect summary statistics
- ❑ Look at your model's outputs, do error **analysis**
- ❑ Tuning hyperparameters, learning rates, getting initialization right, etc. is **often** important to the successes of Neural Networks

Good luck with your projects!

# Word

# Tokenization (review)

- ❑ Tokenization: convert a character stream into words by adding spaces
  - For certain languages, highly nontrivial
  - E.g., Chinese word segmentation is a widely- studied NLP task
  - for other languages (English), tokenization is easier but is still not always obvious

# Another option for text tokenization

- ❑ Instead of
  - white-space segmentation
  - single-character segmentation
- ❑ Use the data to tell us how to tokenize
- ❑ **Subword tokenization** (because tokens can be parts of words as well as whole words)

# Subword tokenization

- ❑ Three common algorithms:
  - Byte-Pair Encoding (BPE) (Sennrich et al., 2016)
  - Unigram language modeling tokenization (Kudo, 2018)
  - WordPiece (Schuster and Nakajima, 2012)
- ❑ All have 2 parts:
  - A token learner that takes a raw training corpus and induces a vocabulary (a set of tokens).
  - A token segmenter that takes a raw test sentence and tokenizes it according to that vocabulary

# Byte Pair Encoding (BPE) token learner

- Let vocabulary be the set of all individual characters
  - = {A, B, C, D, ..., a, b, c, d, ...}
- Repeat:
  - Choose the two symbols that are most frequently adjacent in the training corpus (say 'A', 'B')
  - Add a new merged symbol 'AB' to the vocabulary
  - Replace every adjacent 'A' 'B' in the corpus with 'AB'.
- Until k merges have been done.

# BPE token learner

Original (very fascinating 😊) corpus:

low low low low lowest lowest newer newer newer  
newer newer newer wider wider wider new new

Add end-of-word tokens, resulting in this vocabulary:

**vocabulary**

\_, d, e, i, l, n, o, r, s, t, w

# BPE token learner

**corpus**

5 low \_  
2 lowest \_  
6 newer \_  
3 wider \_  
2 new \_

**vocabulary**

\_, d, e, i, l, n, o, r, s, t, w

Merge **e r** to **er**

**corpus**

5 low \_  
2 lowest \_  
6 newer \_  
3 wider \_  
2 new \_

**vocabulary**

\_, d, e, i, l, n, o, r, s, t, w, er

# BPE

**corpus**

5 low \_  
2 lowest \_  
6 newer \_  
3 wider \_  
2 new \_

**vocabulary**

\_, d, e, i, l, n, o, r, s, t, w, er

Merge er \_ to er\_

**corpus**

5 low \_  
2 lowest \_  
6 newer \_  
3 wider \_  
2 new \_

**vocabulary**

\_, d, e, i, l, n, o, r, s, t, w, er, er\_

# BPE

**corpus**

5 l o w \_  
2 l o w e s t \_  
6 n e w er\_  
3 w i d er\_  
2 n e w \_

**vocabulary**

\_, d, e, i, l, n, o, r, s, t, w, er, er\_

**Merge n e to ne****corpus**

5 l o w \_  
2 l o w e s t \_  
6 ne w er\_  
3 w i d er\_  
2 ne w \_

**vocabulary**

\_, d, e, i, l, n, o, r, s, t, w, er, er\_, ne

# BPE

The next merges are:

**Merge**

(ne, w)

(l, o)

(lo, w)

(new, er\_)

(low, \_\_)

**Current Vocabulary**

\_, d, e, i, l, n, o, r, s, t, w, er, er\_\_, ne, new

\_, d, e, i, l, n, o, r, s, t, w, er, er\_\_, ne, new, lo

\_, d, e, i, l, n, o, r, s, t, w, er, er\_\_, ne, new, lo, low

\_, d, e, i, l, n, o, r, s, t, w, er, er\_\_, ne, new, lo, low, newer\_

\_, d, e, i, l, n, o, r, s, t, w, er, er\_\_, ne, new, lo, low, newer\_\_, low\_\_

# BPE token segmenter algorithm

- ❑ On the test data, run each merge learned from the training data:
  - Greedily
  - In the order we learned them
  - (test frequencies don't play a role)
- ❑ So: merge every e r to er, then merge er \_ to er\_, etc.
- ❑ Result:
  - Test set "n e w e r \_" would be tokenized as a full word
  - Test set "l o w e r \_" would be two tokens: "low er\_

# Properties of BPE tokens

- ❑ Usually include frequent words
- ❑ And frequent subwords
  - Which are often morphemes like -est or –er
- ❑ A **morpheme** is the smallest meaning-bearing unit of a language
  - unlikeliest has 3 morphemes un-, likely, and -est

# Word structure and subword models

- ❑ Let's take a look at the assumptions we've made about a language's vocabulary.
- ❑ We assume a fixed vocab of tens of thousands of words, built from the training set.
- ❑ All *novel* words seen at test time are mapped to a single UNK.

	word	vocab mapping
Common words	hat	→ hat
	learn	→ learn
Variations	taaaaasty	→ UNK
	laern	→ UNK
misspellings		
novel items	Transformerify	→ UNK

# Word structure and subword models

- Finite vocabulary assumptions make even **less** sense in many languages.
- Many languages exhibit complex morphology, or word structure.
  - The effect is more word types, each occurring fewer times.

Example: Swahili verbs can have hundreds of conjugations, each encoding a wide variety of information. (Tense, mood, definiteness, negation, information about the object, ++)

Here's a small fraction of the conjugations for *ambia* – to tell.

Conjugation of <i>-ambia</i>												[less ▲]												
Polarity	Form Infinitive				Non-finite forms								[less ▲]											
	Positive form		Imperative Habitual		Simple finite forms				Complex finite forms															
	Sg.	Pl.	Sg.	Pl.	Sg. / 1	Pl. / 2	3	M-mi	4	5	Ma	6	7	Ki-vi	8	9	N	10	11 / 14	15 / 17	Pa	16	18	
Positive	niliambia naliambia	tuliambia twalliambia	uliambia walliambia	miliambia mwaliambia	allambia	waliambia	uliambia	iliambia	iliambia	yaliambia	kiliambia	viliambia	iliambia	ziliambia	uliambia	kuliambia	paliambia	muliambia	[less ▲]					
Negative	sikuambia sliambia	hatukumbia hatuambilia	hukumbia huambilia	hamkumbia hamambilia	hakuambia	hawakumbia hawauambilia	haukumbia hauambilia	haikumbia haiambilia	halikumbia halambilia	hayakuambi hayambilia	hakiakumbia hakiambilia	haikumbia haiambilia	hazikumbia hazambilia	haukuambia hauambilia	hakuakumbia hakuambilia	hapakuakumbia hapaambilia	hamukuambia hamambilia	hamukuambia hamambilia	[less ▲]					
Positive	ninaambia naambia	tunaambia	unaambia	mnaambia	anaambia	wanaambia	unaambia	inaambia	linaambia	yanaambia	kinaambia	vinaambia	inaambia	zinaambia	unaambia	kunaambia	panaambia	munaambia	[less ▲]					
Negative	siambia sliambia	hatuambilia hatutambilia	huambilia hutambilia	hamambilia hamtambilia	hataambia	wataambia wataambilia	utaambia utaambilia	itaambia itaambilia	yataambia	kitaambia	vitaambia	itaambia	zitaambia	utaambia	kutaambia	pataambia	mutaambia muambia	[less ▲]						
Positive	nitaambia sitaambia	tutaambia	utaambia	mtaambia	ataambia	wataambia wataambilia	utaambia utaambilia	itaambia itaambilia	hautaambia	hataambia	hataambia	hataambia	hataambia	hataambia	hataambia	hataambia	hataambia	hataambia	hataambia	hataambia	hataambia	hataambia	[less ▲]	
Negative	siambia sliambia	hatutambilia hatutambilia	hutambilia hutambilia	hamtambilia hamtambilia	hataambia	wataambilia wataambilia	utaambilia utaambilia	itaambilia itaambilia	hautambilia hautambilia	hataambilia hataambilia	hataambilia hataambilia	hataambilia hataambilia	hataambilia hataambilia	hataambilia hataambilia	hataambilia hataambilia	hataambilia hataambilia	hataambilia hataambilia	hataambilia hataambilia	hataambilia hataambilia	hataambilia hataambilia	hataambilia hataambilia	[less ▲]		
Positive	niambie nisiambie	tuambie tusiambie	uambie usiambie	mambie msiambie	aambie asiambie	waambie wasiambie	uambie usiambie	lambie isiambie	lamble isiamble	yaambie yasiambie	kiambie kisiambie	viambie visiambie	lambie lisiambie	ziambie zisiambie	uambie usiambie	kuambie kusiambie	paambie pasiambie	muambie musiambie	[less ▲]					
Negative	ningeambia nisingeambia singeambia	tungeambia tusingeambia hutungeambia	ungeambia usingeambia hungengeambia	mngeambia msingeambia hamngeambia	angeambia asingeambia hangenambia	wangeambia wasingeambia hawangeambia	ungeambia usingeambia hawangeambia	ingeambia isingeambia haingeambia	lingeambia lisingeambia haingeambia	yangeambia yasingeambia hayangeambia	kingeambia kisingeambia hayingeambia	vingeambia visingeambia havingeambia	ingeambia isingeambia hazingeambia	zingeambia zisingeambia hazingeambia	ungeambia usingeambia hauingeambia	kungeambia kusingeambia hakungeambia	pangeambia pasingeambia hapangeambia	mungeambia musingeambia hamungeambia	[less ▲]					
Positive	ningaliambia nisingliambia singliambia	tungaliambia tusingliambia hutungliambia	ungaliambia usingliambia hungaliambia	mgngaliambia msngliambia hamngaliambia	angaliambia asingliambia hangaliambia	wangaliambia wasingliambia hawangaliambia	ungaliambia usingliambia hawangaliambia	ingaliambia isngaliambia halingaliambia	ingaliambia lisngaliambia halingaliambia	yangaliambia yasingliambia hayangliambia	kingaliambia kisingliambia hayingliambia	vingaliambia visngliambia havngliambia	ingaliambia isngaliambia hakingliambia	zingaliambia zisngaliambia hazingliambia	ungaliambia usingliambia haungliambia	kungaliambia kusingliambia hakungliambia	pangaliambia pasingliambia hapangliambia	mgngaliambia musingliambia hamngaliambia	[less ▲]					
Negative	ningelambia nisingelambia singelambia	tungelambia tusingelambia hutungelambia	ungelambia usingelambia hungelambia	mgngelambia msngelambia hamngelambia	angelambia asingelambia hangelambia	wangelambia wasingelambia hawangelambia	ungelambia usingelambia hawangelambia	ingelambia isngelambia halingelambia	ingelambia lisngelambia halingelambia	yangelambia yasingelambia hayangelambia	kingelambia kisingelambia hayingelambia	vingelambia visngelambia havngelambia	ingelambia isngelambia hakingelambia	zingelambia zisngelambia hazingelambia	ungelambia usingelambia haungelambia	kungelambia kusingelambia hakungelambia	pangelambia pasingelambia hapangelambia	mgngelambia musingelambia hamngelambia	[less ▲]					
Positive	naambia	twaambia	waambia	mwaambia	aambia	waambia	waambia	yaambia	laambia	yaambia	chaambia	vyambia	yaambia	zaambia	waambia	kwaambia	paambia	mwambia	[less ▲]					

# The byte-pair encoding algorithm

- ❑ Subword modeling in NLP encompasses a wide range of methods for reasoning about structure below the word level. (Parts of words, characters, bytes.)
  - The dominant modern paradigm is to learn a vocabulary of parts of words (subword tokens).
  - At training and testing time, each word is split into a sequence of known subwords.
- ❑ **Byte-pair encoding** is a simple, effective strategy for defining a subword vocabulary.
- ❑ Originally used in NLP for machine translation; now a similar method (WordPiece) is used in pretrained models.

[[Sennrich et al., 2016](#), [Wu et al., 2016](#)]

# Word structure and subword models

- Common words end up being a part of the subword vocabulary, while rarer words are split into (sometimes intuitive, sometimes not) components.
- In the worst case, words are split into as many subwords as they have characters.

	word	vocab mapping
Common words	hat	→ hat
	learn	→ learn
Variations	taaaaasty	→ taa## aaa## sty
	laern	→ la## ern
misspellings		
novel items	Transformerify	→ Transformer## ify

# Morphology

- Words are not atoms
  - They have internal structure
  - They are composed (to a first approximation) of **morphemes**
  
- **Morphemes** are minimal meaningful components of words.
  - **Stems** (roots) : The central morphemes in words, which carry the main meaning
  - **Affixes**: bits and pieces that adhere to stems
    - Prefixes: **pre**-nuptial, **ir**-regular
    - Suffixes: determin-**ize**, iterat-**or**
    - Infixes: cup-**s**-ful
    - Circumfixes: **ge**-sammel-**t**

# Morphology

- ❑ Type/Token Ratio across Languages
  - high type/token ratio -> rich morphology
  - low type/token ratio -> poor morphology

# Representing words as discrete symbols

- ❑ In traditional NLP, we regard words as discrete symbols:
  - Such symbols for words can be represented by one-hot vectors:
    - E.g., motel = [0 0 0 0 0 0 0 0 1 0 0 0 0], hotel = [0 0 0 0 0 1 0 0 0 0 0 0 0]
  - Vector dimension = number of words in vocabulary (e.g., 500,000+)

# Problem with words as discrete symbols

- ❑ Example: in web search, if a user searches for “Seattle motel”, we would like to match documents containing “Seattle hotel”
  - But these two vectors are orthogonal
    - $\text{motel} = [0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 1 \ 0 \ 0 \ 0 \ 0]$ ,  $\text{hotel} = [0 \ 0 \ 0 \ 0 \ 0 \ 1 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0]$
- ❑ There is no natural notion of similarity for one-hot vectors!
- ❑ Instead, learn to encode similarity in the vectors themselves

# Intuitions of Distributional Models

- Suppose I gave you the following corpus:
  - A bottle of ***tesgüino*** is on the table
  - Everybody likes ***tesgüino***
  - ***Tesgüino*** makes you drunk
  - We make ***tesgüino*** out of corn.
- What is ***tesgüino***?
- From context, we can guess ***tesgüino*** is an alcoholic beverage like beer
- Intuition: two words are similar if they have similar word contexts

# Many ways to get word vectors

- Some based on counting, some based on prediction/learning
- Some sparse, some dense
- Some have interpretable dimensions, some don't
  
- Shared ideas:
  - model meaning of a word by “embedding” it in a vector space
  - these word vectors are also called “**embeddings**”
  
- Contrast: in traditional NLP, word meaning is represented by a vocabulary index (“word #545”)

# Distributional Word Vectors

- ❑ We'll start with the simplest way to create word vectors:
- ❑ Count occurrences of context words
  - so, vector for *pineapple* has counts of words in the context of *pineapple* in a dataset
  - one entry in vector for each unique context word
  - stack these vectors for all words in a vocabulary  $V$  to produce a count matrix  $C$
  - $C$  is called the **word-context matrix** (or **word-word co-occurrence matrix**)

# Counting Context Words

sugar, a sliced lemon, a tablespoonful of  
r enjoyment. Cautiously she sampled her first  
well suited to programming on the digital  
for the purpose of gathering data and

**apricot** preserve or jam, a pinch each of,  
**pineapple** and another fruit whose taste she likened  
**computer.** In finding the optimal R-stage policy from  
**information** necessary for the study authorized in the

	aardvark	computer	data	pinch	result	sugar	...
apricot	0	0	0	1	0	1	...
pineapple	0	0	0	1	0	1	...
digital	0	2	1	0	1	0	...
information	0	1	6	0	4	0	...
...							

# Word-Context Matrix

- ❑ We showed  $4 \times 6$ , but actual matrix is  $|V| \times |V|$ 
  - Very large, but very **sparse** (mostly zeroes)
  - Lots of efficient algorithms for sparse vectors and matrices
- ❑ Context Window Size
  - Size of context window affects word vectors

# Measuring similarity

- Given 2 word vectors, how should we measure their similarity?
- Most measure of vector similarity are based on **dot product** (or **inner product**):

$$\langle \mathbf{u}, \mathbf{v} \rangle = \mathbf{u} \cdot \mathbf{v} = \mathbf{u}^\top \mathbf{v} = \sum_i u_i v_i$$

- High when vectors have large values in same dimensions
- Notation

$\mathbf{u}$  = a vector

$u_i$  = entry i in the vector

$\mathbf{u}^\top \mathbf{v}$  = dot (inner) product

# Problem with dot product?

$$\mathbf{u}^\top \mathbf{v} = \sum_i u_i v_i$$

- Dot product is larger if vector is longer
- Vector length:

$$\|\mathbf{u}\| = \sqrt{\sum_i u_i^2}$$

- Frequent words → larger counts → larger dot products
- This is bad: we don't want a similarity metric to be overly sensitive to word frequency

# Solution: cosine similarity

- ❑ Divide dot product by lengths of the vectors

$$\frac{\mathbf{u}^\top \mathbf{v}}{\|\mathbf{u}\| \|\mathbf{v}\|}$$

- ❑ Turns out to be the cosine of the angle between them!

# Problems with raw counts

- ❑ Raw word counts are not a great measure of association between words
  - very skewed: *the* and *of* are frequent, but not the most discriminative

# Top co-occurrence counts with ``cooked''

123	,	13	as
92	and	12	for
79	the	12	food
71	.	11	which
68	<s>	11	that
66	</s>	11	meat
53	in	11	can
39	a	11	by
38	is	10	when
35	of	9	rice
30	with	9	raw
28	are	9	beef
25	to	7	they
23	or	7	their
23	it	7	on
20	(	7	not
19	be	7	from
15	)	6	leaves
14	"	6	has

# Top co-occurrence counts with ``cooked''

123	,	13	as
92	and	12	for
79	the	12	<b>food</b>
71	.	11	which
68	<s>	11	that
66	</s>	11	<b>meat</b>
53	in	11	can
39	a	11	by
38	is	10	when
35	of	9	<b>rice</b>
30	with	9	<b>raw</b>
28	are	9	<b>beef</b>
25	to	7	they
23	or	7	their
23	it	7	on
20	(	7	not
19	be	7	from
15	)	6	<b>leaves</b>
14	"	6	has

# Problems with raw counts

- Raw word counts are not a great measure of association between words
  - very skewed: *the* and *of* are frequent, but not the most discriminative
  
- Rather have a measure that asks whether a context word is informative about the center word
  - **pointwise mutual information (PMI)**

# Pointwise Mutual Information (PMI)

- PMI has been used for finding **collocations** and **associations** between words

$$\text{pmi}(x; y) = \log_2 \frac{p(x, y)}{p(x)p(y)}$$

- Here,  $x$  is the center word and  $y$  is the word in the context window
- Each of these probabilities can be estimated from the counts collected from the corpus
- Replace raw counts with PMI scores

# Context words of “cooked” with highest PMIs

9.30533	beef	7.66406	chili
8.88418	shrimp	7.56264	rice
8.63397	potatoes	7.56167	soup
8.61946	ate	7.45315	flour
8.56584	dishes	7.43874	steamed
8.50945	eaten	7.43715	crushed
8.4931	beans	7.41193	meals
8.33137	texture	7.39793	digest
8.29489	vegetables	7.39175	rockies
8.25088	soda	7.34773	ramsay
8.20831	meat	7.33211	honey
8.15708	sauce	7.32253	toxicity
8.08345	consuming	7.29057	cared
7.9532	cuisine	7.28626	tomatoes
7.94043	raw	7.27912	boiling
7.78435	curry	7.27769	dal
7.7563	juice	7.27485	citrus
7.74444	vegetable	7.25649	doncaster

# Positive Pointwise Mutual Information (PPMI)

- ❑ PMI ranges from –infinity to +infinity
- ❑ But negative values are problematic:
  - things are co-occurring **less than** we expect by chance
  - unreliable without enormous corpora
- ❑ So we sometimes replace negative PMI values by 0, calling it positive PMI (PPMI)

# Alternative to PPMI

- ❑ TF-IDF
- ❑ Product of two factors:
  - term frequency (TF; Luhn, 1957): count of word (or possibly log of count)
  - **inverse document frequency** (IDF; Sparck Jones, 1972)
    - $N$ : total number of documents
    - $\text{df}(x)$ : # of documents with word  $x$

$$\text{idf}(x) = \log \frac{N}{\text{df}(x)}$$

# Sparse vs. dense vectors

- ❑ So far, our vectors are
  - **long** (length = 25,000)
  - **sparse** (mostly zero)
  
- ❑ Why might we want to reduce vector dimensionality?

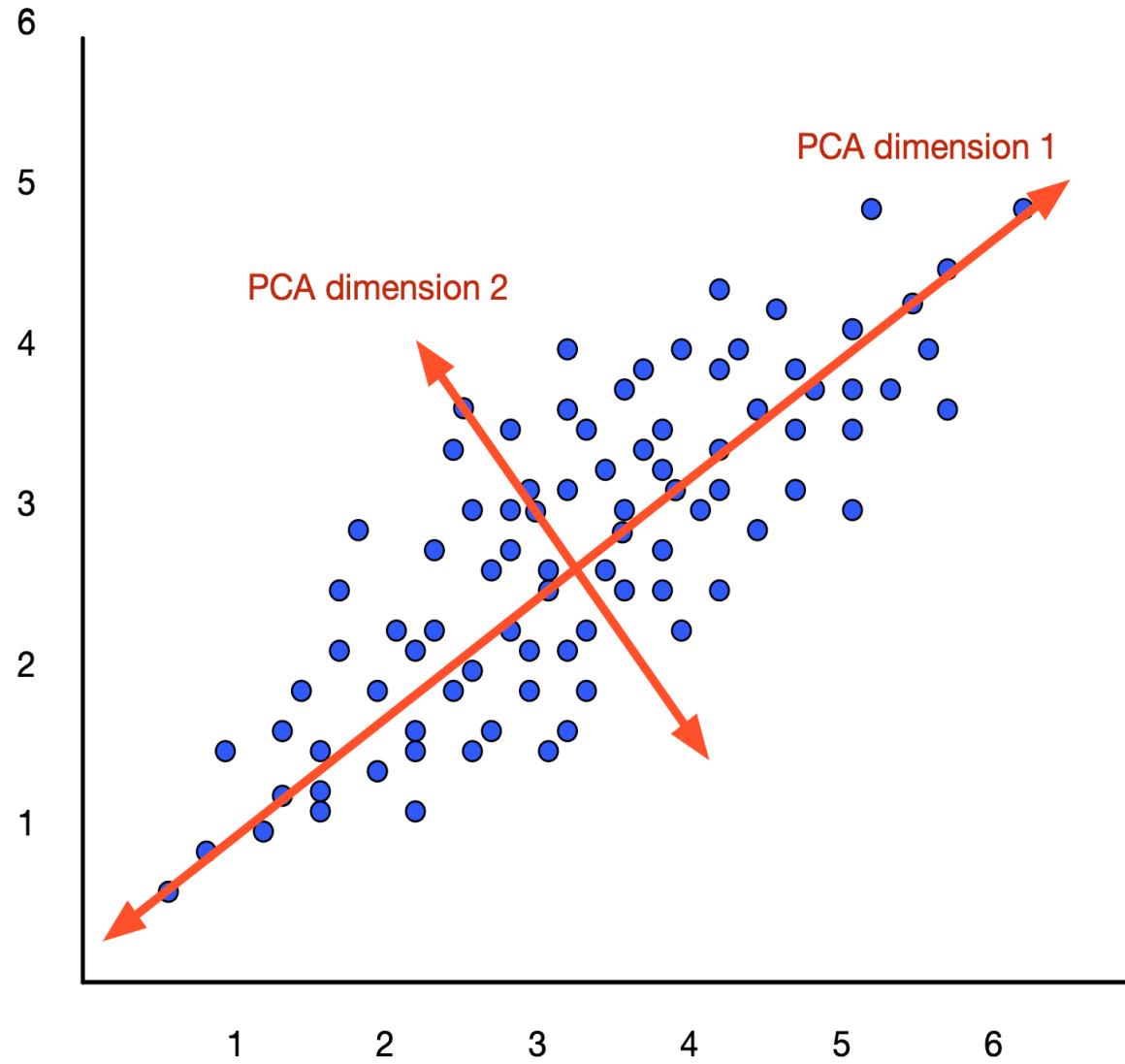
# Why reduce dimensionality?

- Short vectors may be easier to use as features (fewer weights to tune)
- Reducing dimensionality may better handle variability in natural language due to synonymy:
  - car and *automobile* are synonyms, but are distinct dimensions
  - fails to capture similarity between a word with *car* as a neighbor and one with *automobile* as a neighbor

# Dimensionality Reduction: Intuition

- Approximate an  $N$ -dimensional dataset using fewer dimensions:
  - rotate axes into a new space
  - in which first dimension captures most variance in original dataset
  
- Many such (related) methods:
  - principal component analysis (PCA)
  - factor analysis
  - singular value decomposition (SVD)

# Dimensionality Reduction



# SVD embeddings vs. sparse vectors

- ❑ Dense SVD embeddings sometimes work better than sparse PMI vectors at tasks (like word similarity)
  - denoising: low-order dimensions may represent unimportant information
  - truncation may help the models generalize better to unseen data
  - smaller number of dimensions may make it easier for classifiers to effectively assign weights to dimensions for the task
  - dense models may do better at capturing higher order co-occurrence

# How should we evaluate word vectors?

- ❑ WordSim353 (Finkelstein et al., 2002)

word pair	similarity
journey	voyage
king	queen
computer	software
law	lawyer
forest	graveyard
rooster	voyage

# How should we evaluate word vectors?

- ❑ WordSim353 (Finkelstein et al., 2002)

**Instructions:**

Assign a numerical similarity score between 0 and 10  
(0 = words are totally unrelated,  
10 = words are VERY closely related).

**When estimating similarity of antonyms, consider them "similar" (i.e., belonging to the same domain or representing features of the same concept), rather than "dissimilar".**

forest	graveyard	
rooster	voyage	

# How should we evaluate word vectors?

- ❑ WordSim353 (Finkelstein et al., 2002)

word pair		similarity
journey	voyage	9.3
king	queen	8.6
computer	software	8.5
law	lawyer	8.4
forest	graveyard	1.9
rooster	voyage	0.6

# How should we evaluate word vectors?

- SimLex-999 (Hill et al., 2014)

word pair		similarity
insane	crazy	9.6
attorney	lawyer	9.4
author	creator	8.0
diet	apple	1.2
new	ancient	0.2

measures **paraphrastic** similarity:

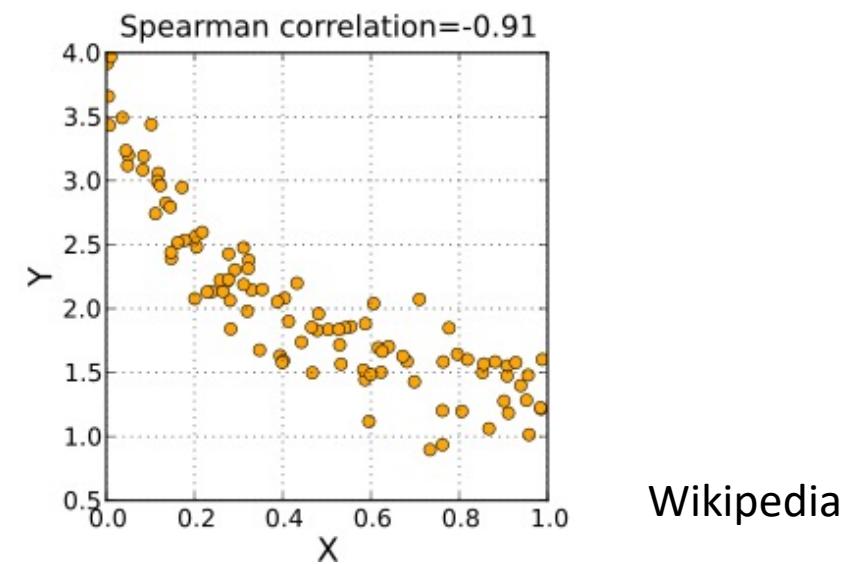
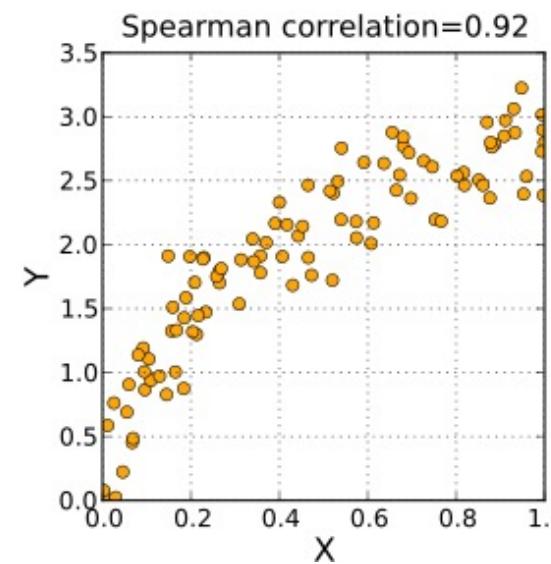
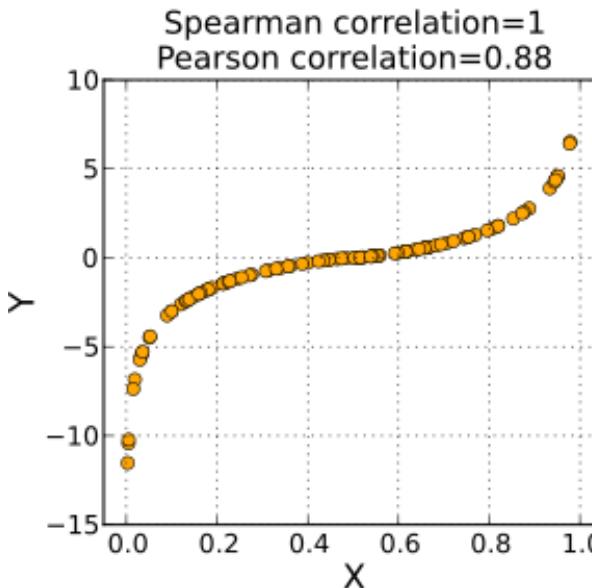
two words are “similar” if they have similar meanings

# How should we evaluate word vectors?

- ❑ There are many word similarity datasets
- ❑ Some focus on topical **relatedness**, others focus on similarity in **meaning**
- ❑ In assignment 1, you will evaluate your word vectors using MEN (relatedness) and SimLex-999 (meaning)

# Evaluation Metrics for Word Similarity

- ❑ Spearman rank correlation coefficient
- ❑ Measures correlation between two variables:
  - variable 1: human-annotated similarities for word pairs
  - variable 2: cosine similarities computed with your word vectors for the same word pairs



Wikipedia

# Text Classification

# Text Classification

- ❑ Simplest user-facing NLP application
- ❑ Email (spam, priority, categories):



- ❑ Sentiment:
- ❑ Topic classification
- ❑ Others?

# Text Classification

- ❑ Datasets
- ❑ Classification
  - Modeling
  - Inference
  - Learning

# NLP Datasets

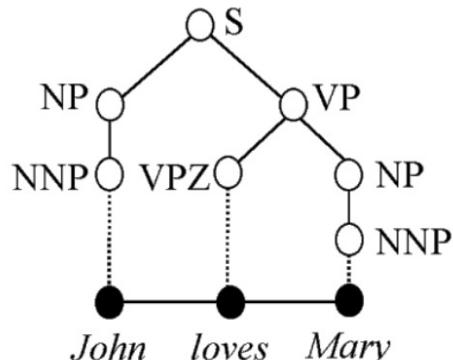
- ❑ NLP datasets include inputs (usually text) and outputs (usually some sort of **annotation**)

# Annotation

- ❑ Supervised machine learning needs labeled datasets, where labels are called **ground truth**
- ❑ In NLP, labels are annotations provided by humans
- ❑ There is always some disagreement among annotators, even for simple tasks
- ❑ These annotations are called a **gold standard**, not ground truth

# How are NLP datasets developed?

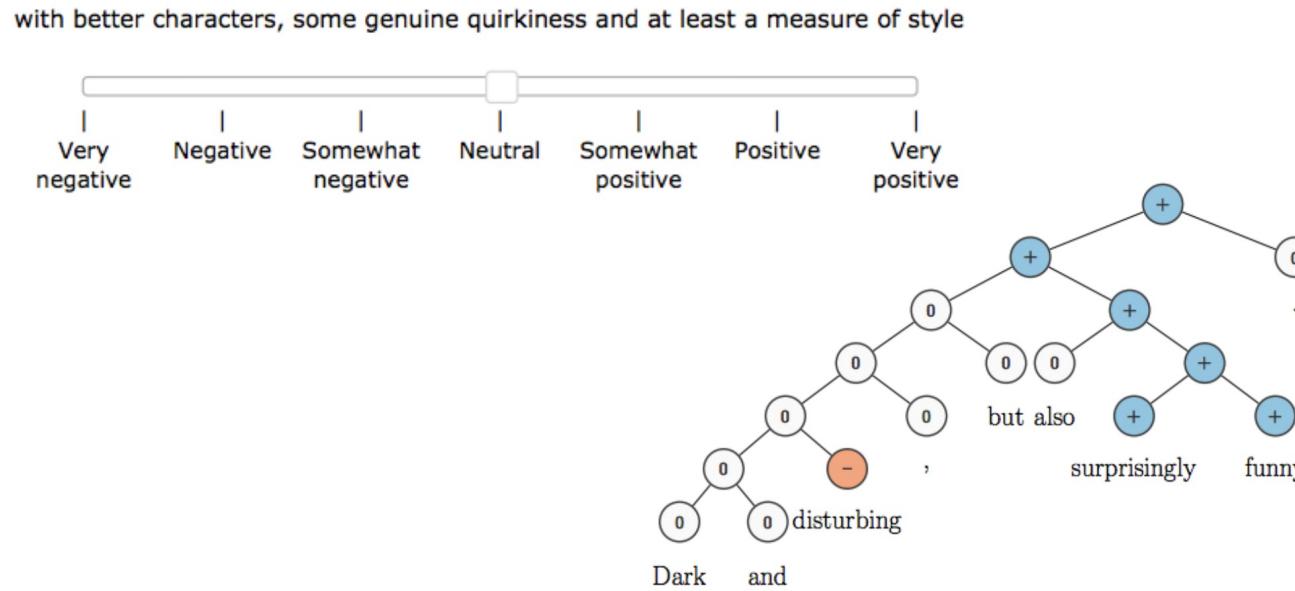
- ❑ Paid, trained human annotators
  - traditional approach
  - researchers write annotation guidelines, recruit & pay annotators (often linguists)
  - more consistent annotations, but costly to scale
  - e.g., Penn Treebank (1993)
    - 1 million words, mostly Wall Street Journal, annotated with part-of-speech tags and syntactic parse trees



# Crowdsourcing

## ❑ Crowdsourcing (e.g., Amazon Mechanical Turk)

- more recent trend
- can't really train annotators, but easier to get multiple annotations for each input (which can then be averaged)
- e.g., Stanford Sentiment Treebank



# Naturally-occurring annotation

- ❑ Long history: used by IBM for speech recognition

There's No Data Like More Data	
• Dick Garwin's correspondence	~2.5M words
• Associated Press	20M words
• Oil company	25M words
• Federal Register	??M words
• American Printing House for the Blind	60M words
• IBM Deposition	100M words
• Canadian Hansard English	100M words

credit: Brown & Mercer, 20 Years of  
Bitext Workshop, 2013

- ❑ How might you find naturally-occurring data for:
  - conversational agents
  - summarization
  - coreference resolution

# Annotator Agreement

- ❑ given annotations from two annotators, how should we measure inter-annotator agreement?
  - percent agreement?
  - Cohen's Kappa (Cohen, 1960) accounts for agreement by chance
  - generalizations exist for more than two annotators (Fleiss, 1971)

# Text Classification Data

- ❑ There are many annotated datasets
  - Stanford Sentiment Treebank: fine-grained sentiment analysis of movie reviews
  - subjectivity/objectivity sentence classification
  - binary sentiment analysis of customer reviews

# Subjectivity/objectivity classification:

the hulk is an anger fueled monster with incredible strength and resistance to damage .	objective
in trying to be daring and original , it comes off as only occasionally satirical and never fresh .	subjective
solondz may well be the only one laughing at his own joke	subjective
obstacles pop up left and right , as the adventure gets wilder and wilder .	objective

- ❑ How was this dataset generated?
  - IMDB plot summaries: objective
  - Rotten Tomatoes snippets: subjective

# customer review sentiment classification:

it works with a minimum of fuss .	positive
i 've had this thing just over a month and the headphone jack has already come loose .	negative
size - bigger than the ipod	negative
you can manage your profile , change the contrast of backlight , make different type of display , either list or tabbed .	positive
i replaced it with a router raizer and it works much better .	negative

# Question classification:

Who invented baseball ?	human
CNN is an acronym for what ?	abbreviation
Which Latin American country is the largest ?	location
How many small businesses are there in the U.S .	number
What would you add to the clay mixture to produce bone china ?	entity
What is the root of all evil ?	description

# Text Classification

- ❑ Datasets
- ❑ **Classification**
  - Modeling
  - Inference
  - Learning

# What is a classifier?

- A function from inputs  $x$  to classification labels  $y$
- One simple type of classifier:
  - for any input  $x$ , assign a score to each label  $y$ , parameterized by parameters  $w$ :
$$\text{score}(x, y, w)$$
- Notation
  - $\mathbf{u}$  = a vector
  - $u_i$  = entry i in the vector
  - $\mathbf{x}$  = a structured object
  - $x_i$  = entry i in the structured object

# What is a classifier?

- A function from inputs  $x$  to classification labels  $y$
- One simple type of classifier:
  - for any input  $x$ , assign a score to each label  $y$ , parameterized by parameters  $w$ :  
$$\text{score}(x, y, w)$$
  - classify by choosing highest-scoring label:

$$\text{classify}(x, w) = \operatorname{argmax}_y \text{score}(x, y, w)$$

# Modeling, Inference, Learning

$$\text{classify}(\mathbf{x}, \mathbf{w}) = \operatorname{argmax}_y \text{score}(\mathbf{x}, y, \mathbf{w})$$

# Modeling, Inference, Learning

modeling: define score function



$$\text{classify}(\mathbf{x}, \mathbf{w}) = \underset{y}{\operatorname{argmax}} \text{ score}(\mathbf{x}, y, \mathbf{w})$$

- ❑ **Modeling:** How do we assign a score to an (x,y) pair using parameters w?

# Modeling, Inference, Learning

**inference:** solve  $\operatorname{argmax}$

**modeling:** define score  $\mathbf{function}$

$$\text{classify}(\mathbf{x}, \mathbf{w}) = \operatorname{argmax}_y \text{score}(\mathbf{x}, y, \mathbf{w})$$

- ❑ **Inference:** How do we efficiently search over the space of all labels?

# Modeling, Inference, Learning

**inference:** solve  $\operatorname{argmax}$

**modeling:** define score function

$$\text{classify}(x, w) = \operatorname{argmax}_y \text{score}(x, y, w)$$

**learning:** choose  $w$

- Learning: How do we choose the weights  $w$ ?

# Modeling, Inference, Learning

**inference:** solve  $\operatorname{argmax}$

**modeling:** define score function

$$\text{classify}(x, w) = \operatorname{argmax}_y \text{score}(x, y, w)$$

**learning:** choose  $w$

- We will use this formulation throughout
  - even when output space is exponentially large or unbounded (e.g., machine translation)

# Binary Sentiment Classification

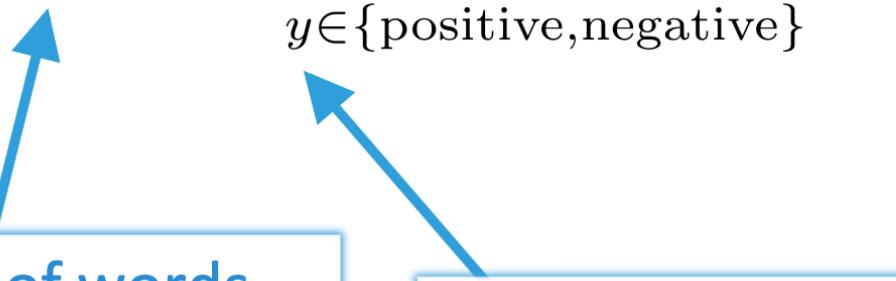
$$\text{classify}(\mathbf{x}, \mathbf{w}) = \operatorname{argmax}_y \text{score}(\mathbf{x}, y, \mathbf{w})$$

# Binary Sentiment Classification

$$\text{classify}(\mathbf{x}, \mathbf{w}) = \underset{y \in \{\text{positive, negative}\}}{\operatorname{argmax}} \text{score}(\mathbf{x}, y, \mathbf{w})$$

a sequence of words  
(a sentence or document)

sentiment label



# Binary Sentiment Classification

modeling: define score function

$$\text{classify}(\mathbf{x}, \mathbf{w}) = \underset{y \in \{\text{positive, negative}\}}{\operatorname{argmax}} \text{score}(\mathbf{x}, y, \mathbf{w})$$

# Linear Models

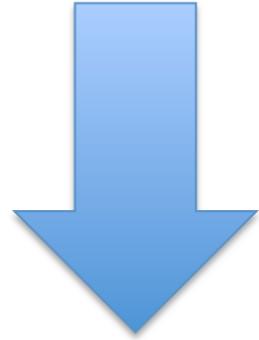
- ❑ Parameters are arranged in a vector  $\mathbf{w}$
- ❑ score function is linear in  $\mathbf{w}$ :

$$\text{score}(\mathbf{x}, y, \mathbf{w}) = \mathbf{w}^\top \mathbf{f}(\mathbf{x}, y) = \sum_i w_i f_i(\mathbf{x}, y)$$

- ❑  $\mathbf{f}$ : vector of feature functions

# Linear Models for Binary Sentiment Classification

$$\text{classify}(\mathbf{x}, \mathbf{w}) = \operatorname{argmax}_{y \in \{\text{positive, negative}\}} \text{score}(\mathbf{x}, y, \mathbf{w})$$



$$\text{classify}(\mathbf{x}, \mathbf{w}) = \operatorname{argmax}_{y \in \{\text{positive, negative}\}} \mathbf{w}^\top \mathbf{f}(\mathbf{x}, y)$$

❑ How do we define  $\mathbf{f}$  ?