

CS60092: Information Retrieval

Text Categorization and Distributed Representations

Prof. Sourangshu Bhattacharya

CSE, IIT Kharagpur

Standing queries

- The path from IR to text classification:
 - You have an information need to monitor, say:
 - [Unrest in the Niger delta region](#)
 - You want to rerun an appropriate query periodically to find new news items on this topic
 - You will be sent new documents that are found
 - I.e., it's not [ranking](#) but [classification](#) (relevant vs. not relevant)
- Such queries are called **standing queries**
 - Long used by “information professionals”
 - A modern mass instantiation is **Google Alerts**
- Standing queries are (hand-written) text classifiers

From: Google Alerts

Subject: **Google Alert - stanford -neuro-linguistic nlp OR "Natural Language Processing" OR parser OR tagger OR ner OR "named entity" OR segmenter OR classifier OR dependencies OR "core nlp" OR corenlp OR phrasal**

Date: May 7, 2012 8:54:53 PM PDT

To: Christopher Manning

Web

3 new results for stanford -neuro-linguistic nlp OR "Natural Language Processing" OR parser OR tagger OR ner OR "named entity" OR segmenter OR classifier OR dependencies OR "core nlp" OR corenlp OR phrasal

[Twitter / Stanford NLP Group: @Robertoross If you only n ...](#)

@Robertoross If you only need tokenization, java -mx2m edu.stanford.nlp.process.PTBTOKENIZER file.txt runs in 2MB on a whole file for me.... 9:41 PM Apr 28th ...

twitter.com/stanfordnlp/status/196459102770171905

[\[Java\] LexicalizedParser lp = LexicalizedParser.loadModel\("edu ...](#)

loadModel("edu/stanford/nlp/models/lexparser/englishPCFG.ser.gz");. String[] sent = { "This", "is", "an", "easy", "sentence", "." };. Tree parse = lp.apply(Arrays.

pastebin.com/az14R9nd

[More Problems with Statistical NLP || kuro5hin.org](#)

Tags: nlp, ai, coursera, stanford, nlp-class, cky, nltk, reinventing the wheel, ... Programming Assignment 6 for Stanford's nlp-class is to implement a CKY parser .

www.kuro5hin.org/story/2012/5/5/11011/68221

Tip: Use quotes ("like this") around a set of words in your query to match them exactly. [Learn more.](#)

[Delete](#) this alert.

[Create](#) another alert.

[Manage](#) your alerts.

Spam filtering

Another text classification task

From: "" <takworld@hotmail.com>

Subject: real estate is the only way... gem oalvgkay

Anyone can buy real estate with no money down

Stop paying rent TODAY !

There is no need to spend hundreds or even thousands for similar courses

I am 22 years old and I have already purchased 6 properties using the methods outlined in this truly INCREDIBLE ebook.

Change your life NOW !

=====

Click Below to order:

<http://www.wholesaledaily.com/sales/nmd.htm>

=====

Categorization/Classification

- Given:
 - A representation of a document d
 - Issue: how to represent text documents.
 - Usually some type of high-dimensional space – bag of words
 - A fixed set of classes:
$$C = \{c_1, c_2, \dots, c_J\}$$
- Determine:
 - The category of d : $\gamma(d) \in C$, where $\gamma(d)$ is a **classification function**
 - We want to build classification functions (“classifiers”).

Classification Methods (1)

- Manual classification
 - Used by the original Yahoo! Directory
 - Looksmart, about.com, ODP, **PubMed**
 - Accurate when job is done by experts
 - Consistent when the problem size and team is small
 - Difficult and expensive to scale
 - Means we need automatic classification methods for big problems

Classification Methods (2)

- Hand-coded rule-based classifiers
 - One technique used by news agencies, intelligence agencies, etc.
 - Widely deployed in government and enterprise
 - Vendors provide “IDE” for writing such rules

Classification Methods (2)

- Hand-coded rule-based classifiers
 - Commercial systems have complex query languages
 - Accuracy can be high if a rule has been carefully refined over time by a subject expert
 - Building and maintaining these rules is expensive

A Verity topic

A complex classification rule: art

```

comment line      # Beginning of art topic definition
top-level topic  art ACCRUE
                 /author = "fsmith"
topic definition modifiers {
                 /date  = "30-Dec-01"
                 /annotation = "Topic created
                             by fsmith"

subtopic topic    * 0.70 performing-arts ACCRUE
evidencetopic    ** 0.50 WORD
topic definition modifier /wordtext = ballet
evidencetopic    ** 0.50 STEM
topic definition modifier /wordtext = dance
evidencetopic    ** 0.50 WORD
topic definition modifier /wordtext = opera
evidencetopic    ** 0.30 WORD
topic definition modifier /wordtext = symphony
subtopic         * 0.70 visual-arts ACCRUE
                 ** 0.50 WORD
                 /wordtext = painting
                 ** 0.50 WORD
                 /wordtext = sculpture

subtopic         * 0.70 film ACCRUE
                 ** 0.50 STEM
                 /wordtext = film

subtopic         ** 0.50 motion-picture PHRASE
                 *** 1.00 WORD
                 /wordtext = motion
                 *** 1.00 WORD
                 /wordtext = picture
                 ** 0.50 STEM
                 /wordtext = movie

subtopic         * 0.50 video ACCRUE
                 ** 0.50 STEM
                 /wordtext = video
                 ** 0.50 STEM
                 /wordtext = vcr
# End of art topic

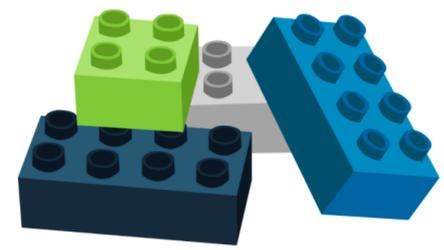
```

■ Note:

- maintenance issues (author, etc.)
- Hand-weighting of terms

[Verity was bought by Autonomy in 2005, which was bought by HP in 2011 – a mess; I think it no longer exists ...]

At the heart of eContext is its taxonomy of the commercial and social web. This taxonomy provides a consistent framework that groups conversations, behavior, and digital interactions into topical hierarchies up to 21 tiers, with clear connections between those topics. Clients push large amounts of text data from multiple sources through eContext, creating a single data set for measurement and analytics. With this holistic perspective, they can distill more comprehensive, meaningful insights, in less time.



Schedule a Demo

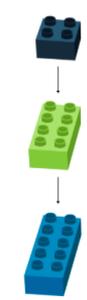


Enhancing Virtual Agents with Structured Knowledge



As clients feed information through our API, our custom pre-processing parses the text into linguistically appropriate units that eContext can work with efficiently.

Then eContext checks each unit against its 55 million positive and negative rules to determine the correct topic, or topics, that content exhibits. This happens hundreds of thousands of times per second, as we enrich the highest velocity data streams available in our industry — everything up to and including the full firehose of Tweets. Thanks to the rules that govern eContext’s decision-making, and the 450,000 different concepts it recognizes, eContext is able to deliver deep, broad, and accurate results.



Download White Paper

Classification Methods (3): Supervised learning

- Given:
 - A document d
 - A fixed set of classes:
 $C = \{c_1, c_2, \dots, c_J\}$
 - A training set D of documents each with a label in C
- Determine:
 - A learning method or algorithm which will enable us to learn a classifier γ
 - For a test document d , we assign it the class
 $\gamma(d) \in C$

Classification Methods (3)

- Supervised learning
 - Naive Bayes (simple, common) – see video, cs229
 - k-Nearest Neighbors (simple, powerful)
 - Support-vector machines (newer, generally more powerful)
 - Decision trees → random forests → gradient-boosted decision trees (e.g., xgboost)
 - ... plus many other methods
 - No free lunch: need hand-classified training data
 - But data can be built up by amateurs
- Many commercial systems use a mix of methods

Features

- Supervised learning classifiers can use any sort of feature
 - URL, email address, punctuation, capitalization, dictionaries, network features
- In the simplest bag of words view of documents
 - We use **only** word features
 - we use **all** of the words in the text (not a subset)

The bag of words representation

Y (

I love this movie! It's sweet, but with satirical humor. The dialogue is great and the adventure scenes are fun.. It manages to be whimsical and romantic while laughing at the conventions of the fairy tale genre. I would recommend it to just about anyone. I've seen it several times, and I'm always happy to see it again whenever I have a friend who hasn't seen it yet.

)

= C

The bag of words representation

Y (

great	2
love	2
recommend	1
laugh	1
happy	1
...	...

) = C

Feature Selection: Why?

- Text collections have a large number of features
 - 10,000 – 1,000,000 unique words ... and more
- Selection may make a particular classifier feasible
 - Some classifiers can't deal with 1,000,000 features
- Reduces training time
 - Training time for some methods is quadratic or worse in the number of features
- Makes runtime models smaller and faster
- Can improve generalization (performance)
 - Eliminates noise features
 - Avoids overfitting

Feature Selection: Frequency

- The simplest feature selection method:
 - Just use the commonest terms
 - No particular foundation
 - But it make sense why this works
 - They' re the words that can be well-estimated and are most often available as evidence
 - In practice, this is often 90% as good as better methods
- Smarter feature selection:
 - chi-squared, etc.

Naïve Bayes: See IIR 13 or cs124 lecture on Coursera or cs229

- Classify based on prior weight of class and conditional parameter for what each word says:

$$c_{NB} = \operatorname{argmax}_{c_j \in C} \left[\log P(c_j) + \sum_{i \in \text{positions}} \log P(x_i | c_j) \right]$$

- Training is done by counting and dividing:

$$P(c_j) \leftarrow \frac{N_{c_j}}{N} \quad P(x_k | c_j) \leftarrow \frac{T_{c_j x_k} + \alpha}{\sum_{x_i \in V} [T_{c_j x_i} + \alpha]}$$

- Don't forget to smooth

SpamAssassin

- Naïve Bayes has found a home in spam filtering
 - Paul Graham's A Plan for Spam
 - Widely used in spam filters
 - But many features beyond words:
 - black hole lists, etc.
 - particular hand-crafted text patterns

SpamAssassin Features:

- Basic (Naïve) Bayes spam probability
- Mentions: Generic Viagra
- Regex: millions of (dollar) ((dollar) NN,NNN,NNN.NN)
- Phrase: impress ... girl
- Phrase: 'Prestigious Non-Accredited Universities'
- From: starts with many numbers
- Subject is all capitals
- HTML has a low ratio of text to image area
- Relay in RBL, http://www.mail-abuse.com/enduserinfo_rbl.html
- RCVD line looks faked
- http://spamassassin.apache.org/tests_3_3_x.html

Naive Bayes is Not So Naive

- Very fast learning and testing (basically just count words)
- Low storage requirements
- Very good in domains with many equally important features
- More robust to irrelevant features than many learning methods

Irrelevant features cancel out without affecting results

Naive Bayes is Not So Naive

- More robust to concept drift (changing class definition over time)
- Naive Bayes won 1st and 2nd place in KDD-CUP 97 competition out of 16 systems
 - Goal: Financial services industry direct mail response prediction: Predict if the recipient of mail will actually respond to the advertisement – 750,000 records.
- A good dependable baseline for text classification (but not the best)!

Evaluating Categorization

- Evaluation must be done on test data that are independent of the training data
 - Sometimes use cross-validation (averaging results over multiple training and test splits of the overall data)
- Easy to get good performance on a test set that was available to the learner during training (e.g., just memorize the test set)

Evaluating Categorization

- Measures: precision, recall, F1, classification accuracy
- **Classification accuracy**: r/n where n is the total number of test docs and r is the number of test docs correctly classified

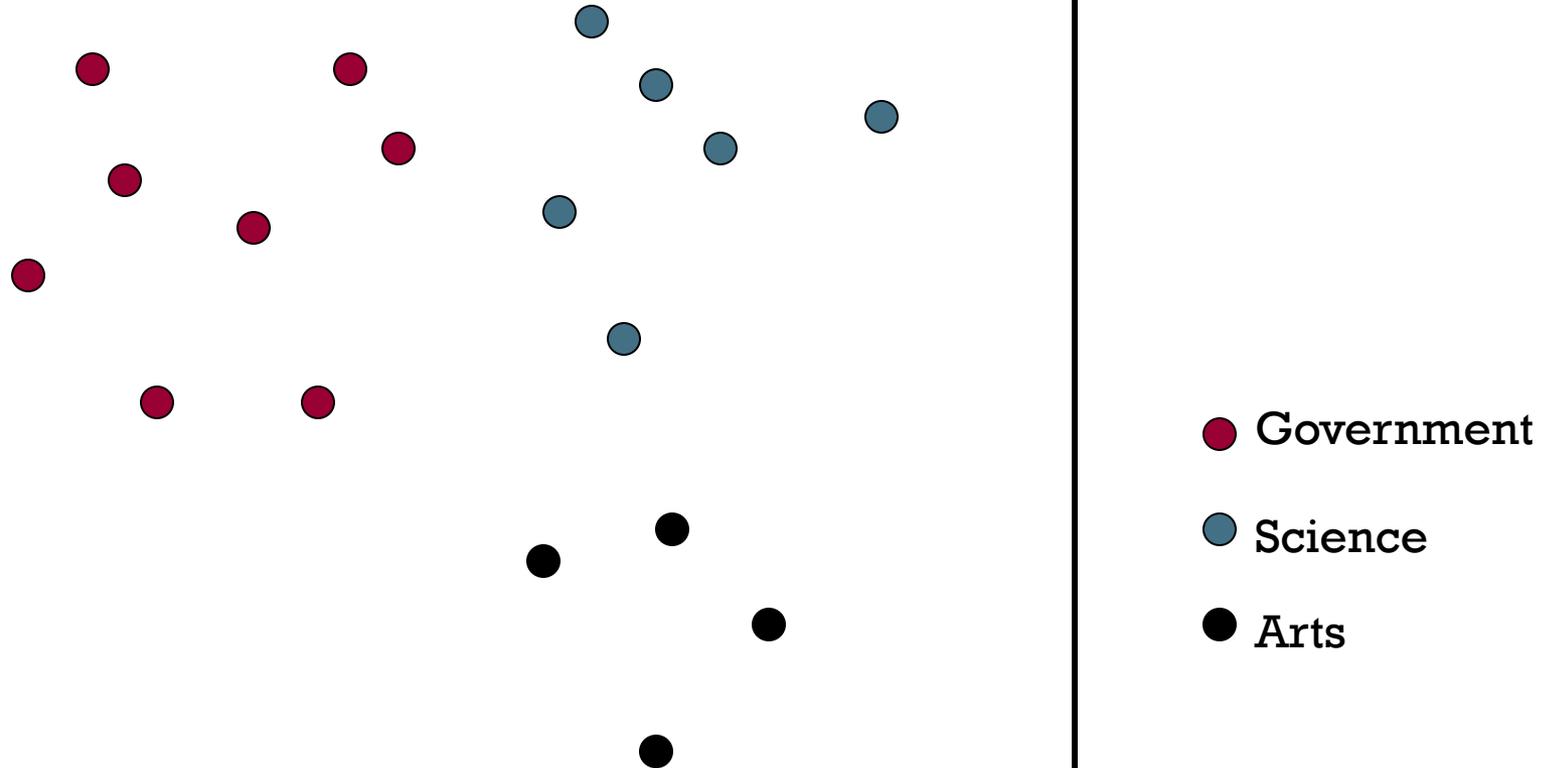
Remember: Vector Space Representation

- Each document is a vector, one component for each term (= word).
- Normally normalize vectors to unit length.
- High-dimensional vector space:
 - Terms are axes
 - 10,000+ dimensions, or even 100,000+
 - Docs are vectors in this space
- How can we do classification in this space?

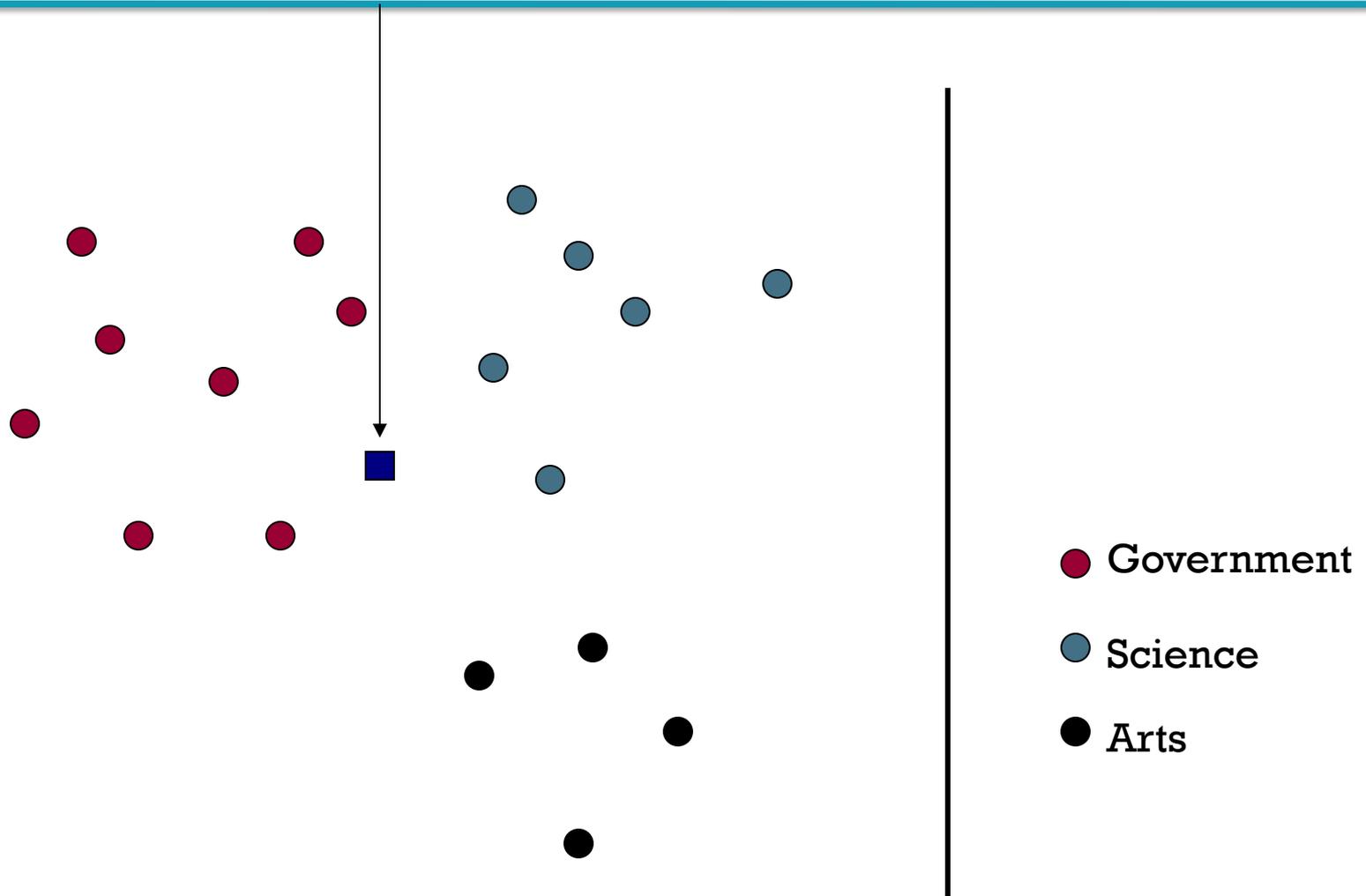
Classification Using Vector Spaces

- In vector space classification, training set corresponds to a labeled set of points (equivalently, vectors)
- **Premise 1:** Documents in the same class form a contiguous region of space
- **Premise 2:** Documents from different classes don't overlap (much)
- Learning a classifier: build surfaces to delineate classes in the space

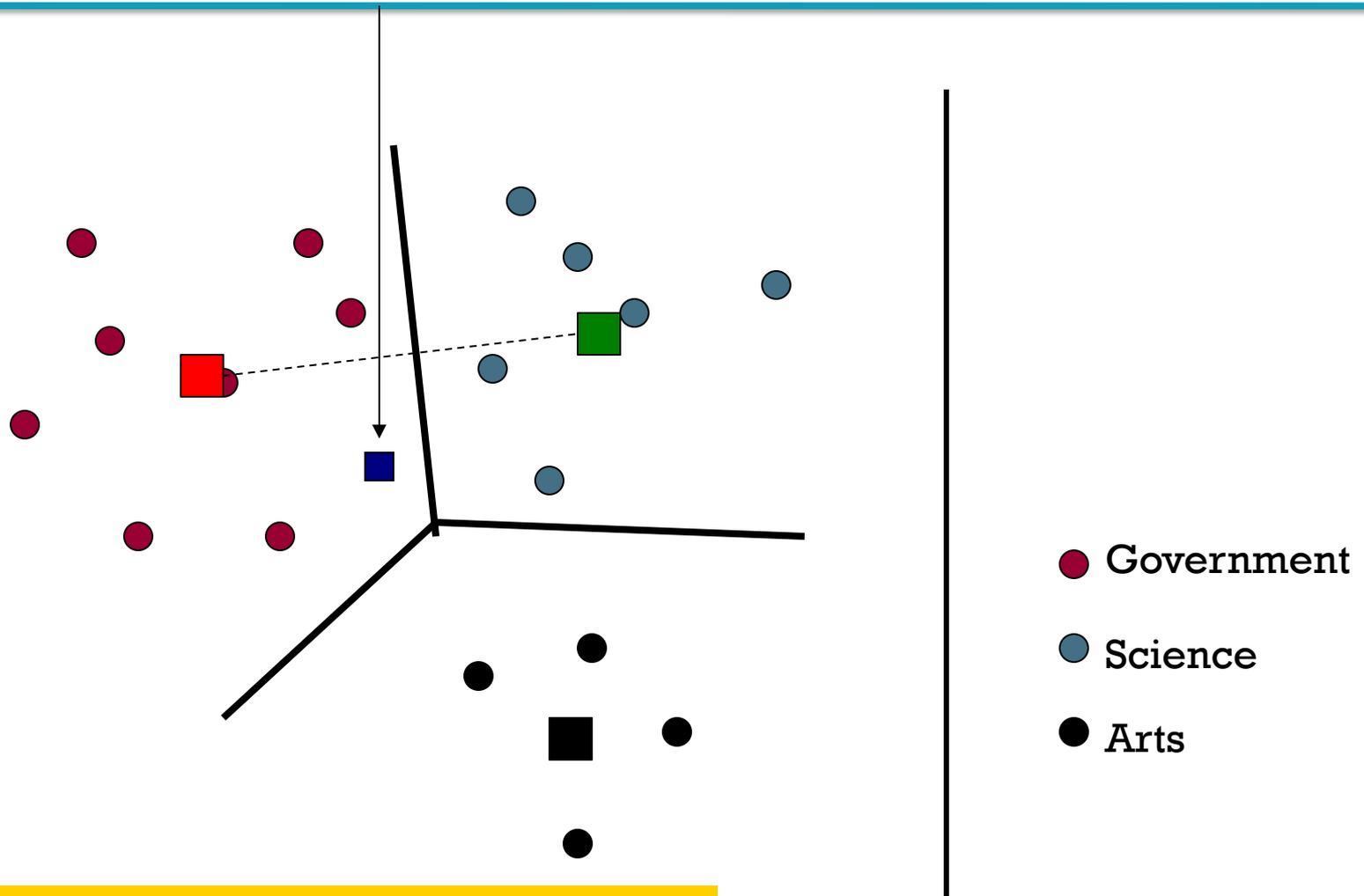
Documents in a Vector Space



Test Document of what class?



Test Document = Government



Our focus: how to find good separators

Definition of centroid

$$\vec{\mu}(c) = \frac{1}{|D_c|} \sum_{d \in D_c} \vec{v}(d)$$

- Where D_c is the set of all documents that belong to class c and $v(d)$ is the vector space representation of d .
- *Note that centroid will in general not be a unit vector even when the inputs are unit vectors.*

Rocchio classification

- Rocchio forms a simple representative for each class: the centroid/prototype
- Classification: nearest prototype/centroid
- It does not guarantee that classifications are consistent with the given training data

Why not?

Two-class Rocchio as a linear classifier

- Line or hyperplane defined by:

$$\sum_{i=1}^M w_i d_i = \theta$$

- For Rocchio, set:

$$\vec{w} = \vec{\mu}(c_1) - \vec{\mu}(c_2)$$

$$\theta = 0.5 \times (|\vec{\mu}(c_1)|^2 - |\vec{\mu}(c_2)|^2)$$

Linear classifier: Example

- Class: “interest” (as in interest rate)
- Example features of a linear classifier
 - 0.70 w_i prime t_i
 - 0.67 rate
 - 0.63 interest
 - 0.60 rates
 - 0.46 discount
 - 0.43 bundesbank
 - -0.71 w_i dlrs t_i
 - -0.35 world
 - -0.33 sees
 - -0.25 year
 - -0.24 group
 - -0.24 dlr
- To classify, find dot product of feature vector and weights

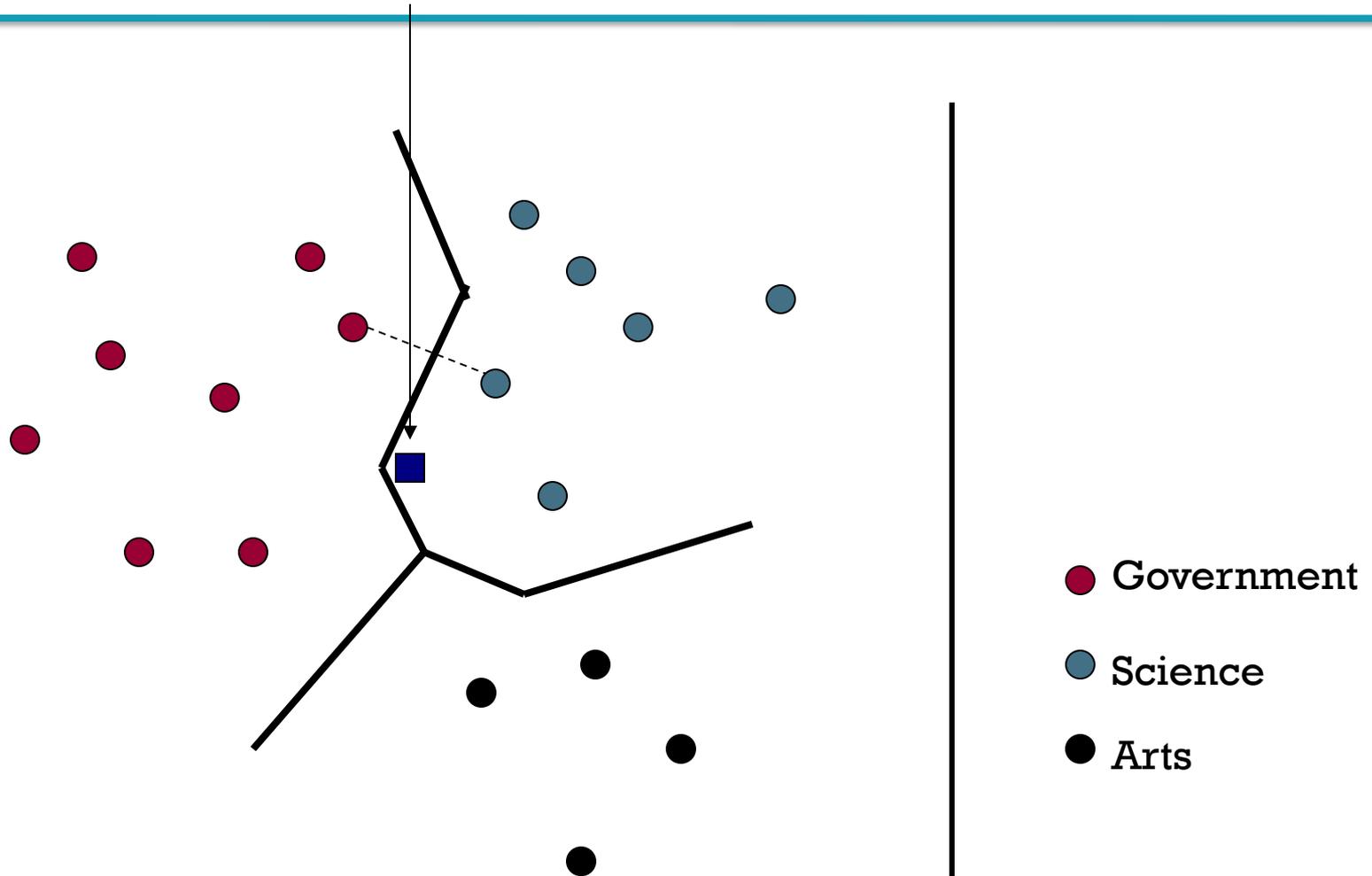
Rocchio classification

- A simple form of Fisher's linear discriminant
- Little used outside text classification
 - It has been used quite effectively for text classification
 - But in general worse than Naïve Bayes
- Again, cheap to train and test documents

k Nearest Neighbor Classification

- k NN = k Nearest Neighbor
- To classify a document d :
- Define k -neighborhood as the k nearest neighbors of d
- Pick the majority class label in the k -neighborhood
- For larger k can roughly estimate $P(c|d)$ as $\#(c)/k$

Test Document = Science



Voronoi diagram

Nearest-Neighbor Learning

- Learning: just store the labeled training examples D
- Testing instance x (*under 1NN*):
 - Compute similarity between x and all examples in D .
 - Assign x the category of the most similar example in D .
- Does not compute anything beyond storing the examples
- Also called:
 - Case-based learning
 - Memory-based learning
 - Lazy learning
- Rationale of kNN: contiguity hypothesis

k Nearest Neighbor

- Using only the closest example (1NN) is subject to errors due to:
 - A single atypical example.
 - Noise (i.e., an error) in the category label of a single training example.
- More robust: find the k examples and return the majority category of these k
- k is typically odd to avoid ties; 3 and 5 are most common

Nearest Neighbor with Inverted Index

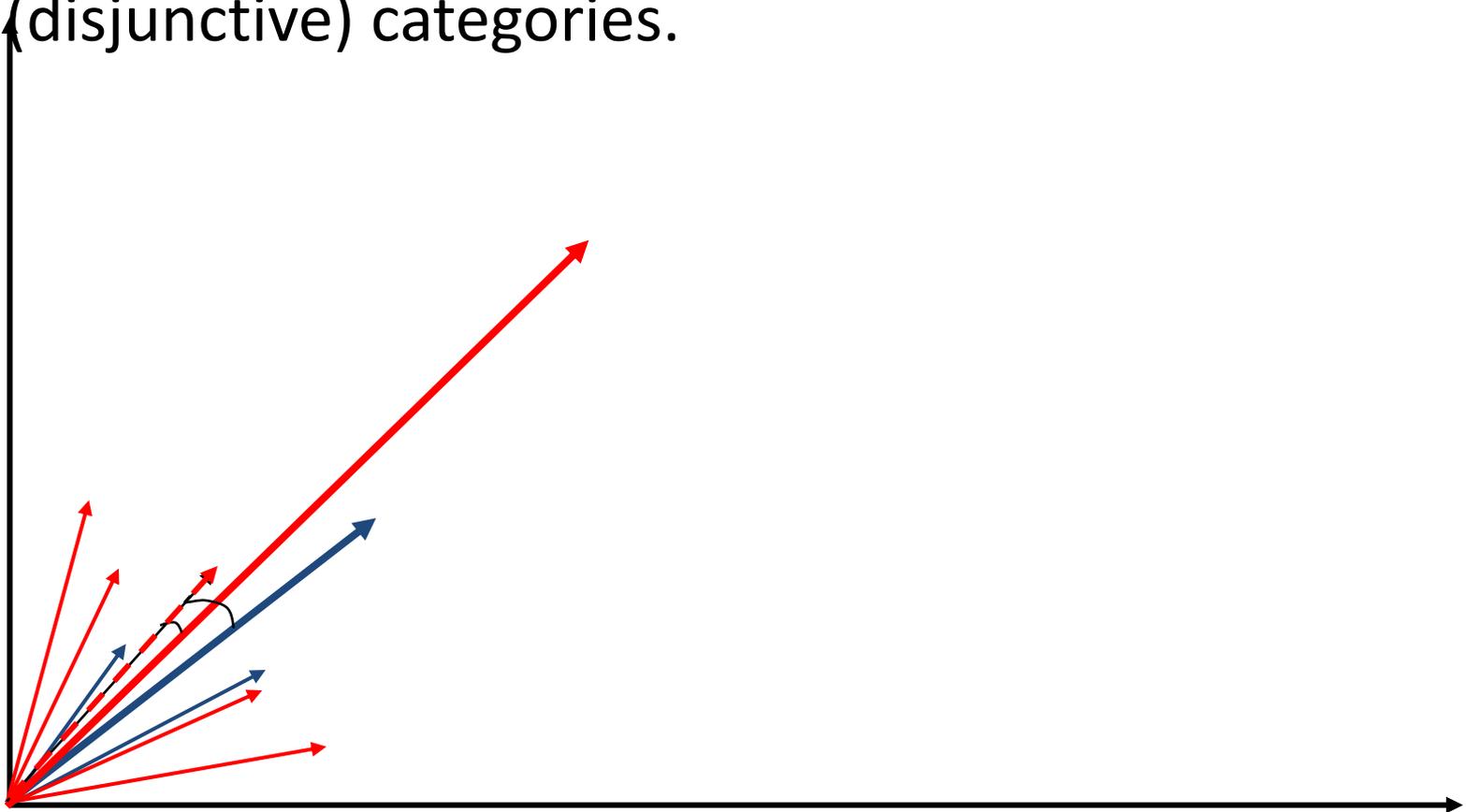
- Naively finding nearest neighbors requires a linear search through $|D|$ documents in collection
- But determining k nearest neighbors is the same as determining the k best retrievals using the test document as a query to a database of training documents.
- Use standard vector space inverted index methods to find the k nearest neighbors.
- **Testing Time:** $O(B|V_t|)$ where B is the average number of training documents in which a test-document word appears.
 - Typically $B \ll |D|$

kNN: Discussion

- No feature selection necessary
- No training necessary
- Scales well with large number of classes
 - Don't need to train n classifiers for n classes
- Classes can influence each other
 - Small changes to one class can have ripple effect
- Done naively, very expensive at test time
- In most cases it's more accurate than NB or Rocchio
 - As the amount of data goes to infinity, it has to be a great classifier! – it's “Bayes optimal”

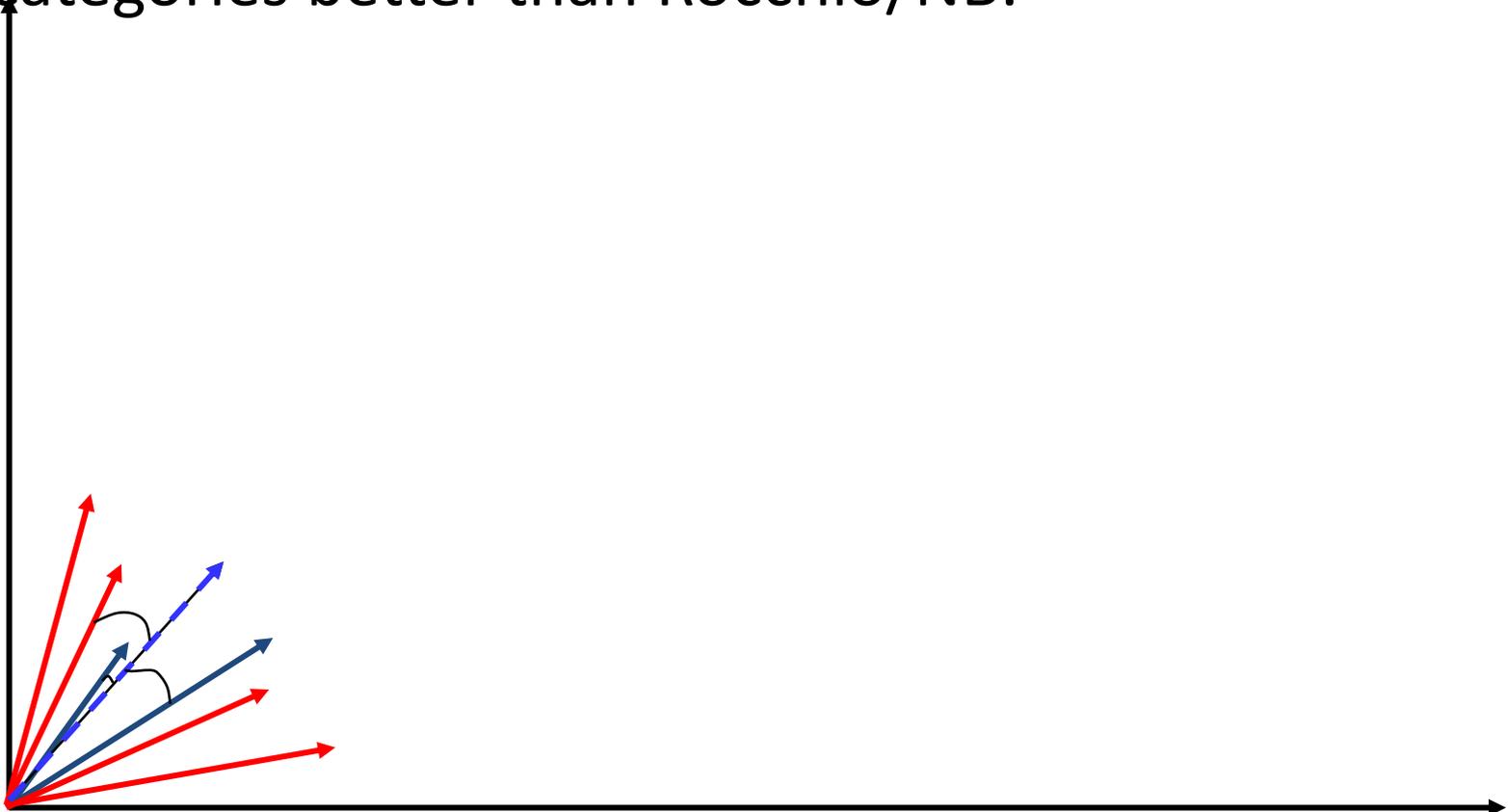
Rocchio Anomaly

- Prototype models have problems with polymorphic (disjunctive) categories.



3 Nearest Neighbor vs. Rocchio

- Nearest Neighbor tends to handle polymorphic categories better than Rocchio/NB.



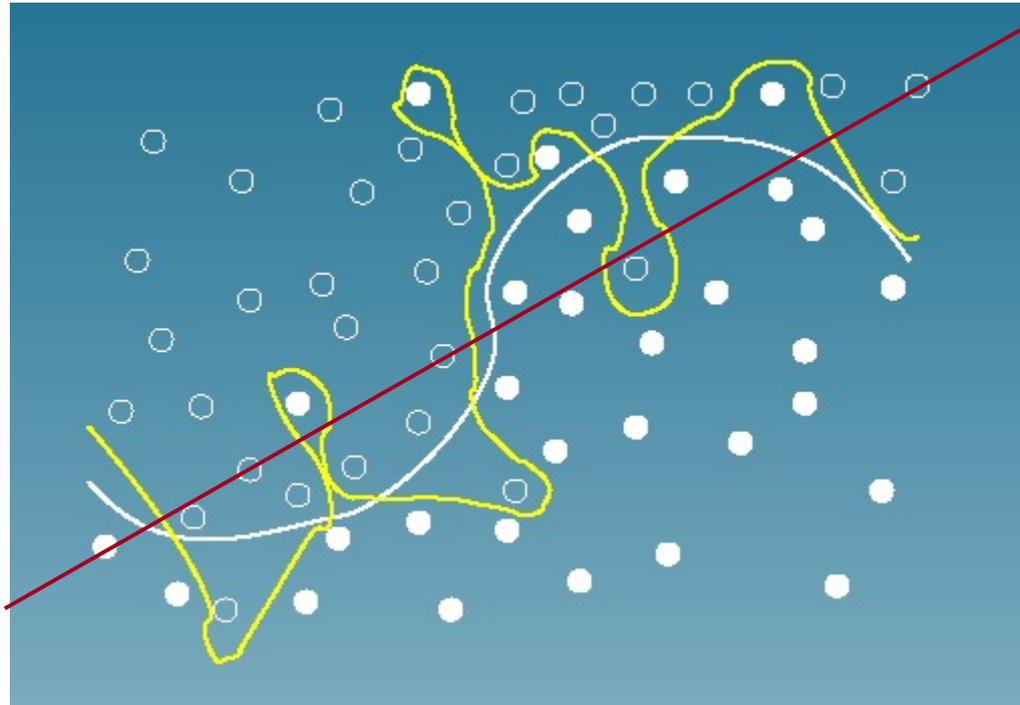
Bias vs. capacity – notions and terminology

- Consider asking a botanist: **Is an object a tree?**
 - Too much *capacity*, low *bias*
 - Botanist who memorizes
 - Will always say “no” to new object (e.g., different # of leaves)
 - Not enough capacity, high bias
 - Lazy botanist
 - Says “yes” if the object is green
 - You want the middle ground

kNN vs. Naive Bayes

- Bias/Variance tradeoff
 - Variance \approx Capacity
- kNN has **high variance** and **low bias**.
 - Infinite memory
- Rocchio/NB has **low variance** and **high bias**.
 - Linear decision surface between classes

Bias vs. variance: Choosing the correct model capacity



Summary: Representation of Text Categorization Attributes

- Representations of text are usually very high dimensional
 - “The curse of dimensionality”
- High-bias algorithms should generally work best in high-dimensional space
 - They prevent overfitting
 - They generalize more
- For most text categorization tasks, there are many relevant features & many irrelevant ones

Which classifier do I use for a given text classification problem?

- Is there a learning method that is optimal for all text classification problems?
- No, because there is a tradeoff between bias and variance.
- Factors to take into account:
 - How much training data is available?
 - How simple/complex is the problem? (linear vs. nonlinear decision boundary)
 - How noisy is the data?
 - How stable is the problem over time?
 - For an unstable problem, it's better to use a simple and robust classifier.

DISTRIBUTED REPRESENTATIONS

How can we more robustly match a user's search intent?

We want to **understand** a query, not just do String equals()

- If user searches for [Dell notebook battery size], we would like to match documents discussing “Dell laptop battery capacity”
- If user searches for [Seattle motel], we would like to match documents containing “Seattle hotel”

A pure keyword-matching IR system does nothing to help....

Simple facilities that we have already discussed do a bit to help

- Spelling correction
- Stemming / case folding

But we'd like to better **understand** when query/document match

How can we more robustly match a user's search intent?

Query expansion:

- **Relevance feedback** could allow us to capture this if we get near enough to matching documents with these words
- We can also use information on **word similarities**:
 - A manual **thesaurus** of synonyms for query expansion
 - A **measure of word similarity**
 - Calculated from a big document collection
 - Calculated by query log mining (common on the web)

Document expansion:

- Use of **anchor text** may solve this by providing human authored synonyms, but not for new or less popular web pages, or non-hyperlinked collections

Example of manual thesaurus

The screenshot displays the PubMed interface. At the top left is the NCBI logo, and at the top center is the PubMed logo. On the top right is the National Library of Medicine (NLM) logo. Below the logos is a navigation bar with tabs for PubMed, Nucleotide, Protein, Genome, Structure, PopSet, and Taxonomy. The search bar contains the text "Search PubMed for cancer" with "PubMed" in a dropdown menu. To the right of the search bar are "Go" and "Clear" buttons. Below the search bar are links for "Limits", "Preview/Index", "History", "Clipboard", and "Details". On the left side, there is a sidebar with links for "About Entrez", "Text Version", "Entrez PubMed", "Overview", "Help | FAQ", "Tutorial", "New/Noteworthy", "E-Utilities", "PubMed Services", "Journals Database", "MeSH Browser", "Single Citation", and "Metabay". The main content area shows the "PubMed Query:" section with the following query:

```
("neoplasms"[MeSH Terms] OR cancer[Text Word])
```

 At the bottom of the query area are "Search" and "URL" buttons.

Search log query expansion

- Context-free query expansion ends up problematic
 - [wet ground] \approx [wet earth]
 - So expand [ground] \Rightarrow [ground earth]
 - But [ground coffee] \neq [earth coffee]
- You can learn query context-specific rewritings from search logs by attempting to identify the same user making a second attempt at the same user need
 - [Hinton word vector]
 - [Hinton word embedding]
- In this context, [vector] \approx [embedding]
 - But not when talking about a *disease vector* or C++!

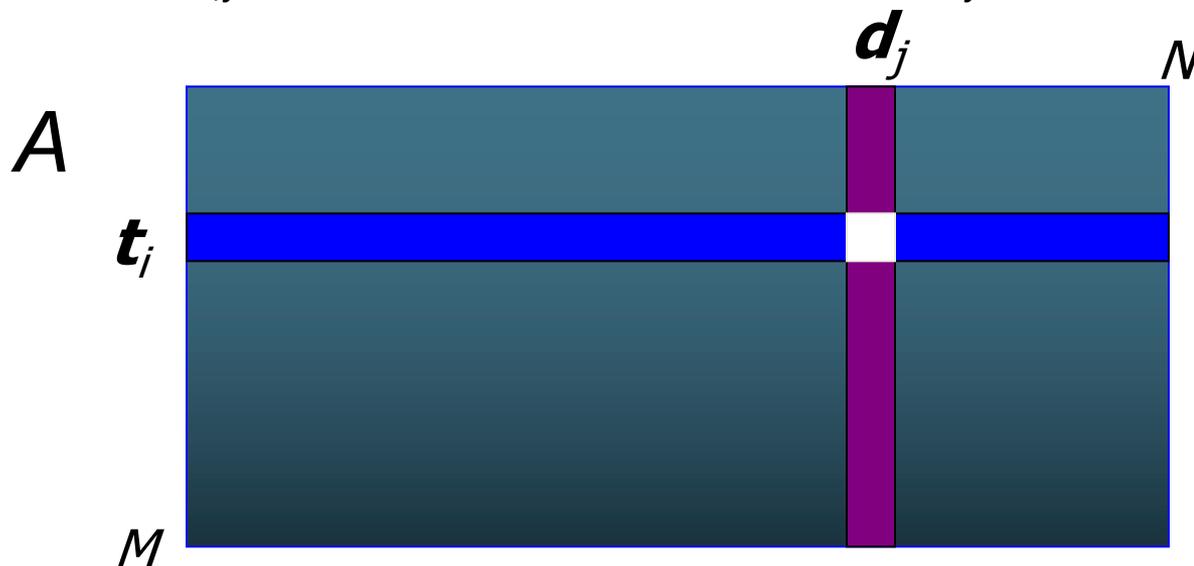
Automatic Thesaurus Generation

- Attempt to generate a thesaurus automatically by analyzing a collection of documents
- Fundamental notion: similarity between two words
- **Definition 1: Two words are similar if they co-occur with similar words.**
- **Definition 2: Two words are similar if they occur in a given grammatical relation with the same words.**
- You can harvest, peel, eat, prepare, etc. apples and pears, so apples and pears must be similar.
- **Co-occurrence based is more robust, grammatical relations are more accurate.**



Simple Co-occurrence Thesaurus

- Simplest way to compute one is based on term-term similarities in $C = AA^T$ where A is term-document matrix.
- $w_{i,j}$ = (normalized) weight for (t_i, d_j)



- For each t_i , pick terms with high values in C

What does C contain if A is a term-doc incidence (0/1) matrix?

Automatic thesaurus generation

example ... sort of works

Word	Nearest neighbors
absolutely	absurd, whatsoever, totally, exactly, nothing
bottomed	dip, copper, drops, topped, slide, trimmed
captivating	shimmer, stunningly, superbly, plucky, witty
doghouse	dog, porch, crawling, beside, downstairs
makeup	repellent, lotion, glossy, sunscreen, skin, gel
mediating	reconciliation, negotiate, cease, conciliation
keeping	hoping, bring, wiping, could, some, would
lithographs	drawings, Picasso, Dali, sculptures, Gauguin
pathogens	toxins, bacteria, organisms, bacterial, parasites
senses	grasp, psyche, truly, clumsy, naïve, innate

Too little data (10s of millions of words) treated by **too sparse method**.
 100,000 words = 10^{10} entries in C .

How can we represent term relations?

- With the standard symbolic encoding of terms, each term is a dimension
- Different terms have no inherent similarity
- $\text{motel} [0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 1 \ 0 \ 0 \ 0 \ 0]^T$
 $\text{hotel} [0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 3 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0]^T = 0$
- If query on *hotel* and document has *motel*, then our query and document vectors are **orthogonal**

Can you directly learn term relations?

- Basic IR is scoring on $q^T d$
- No treatment of synonyms; no machine learning
- Can we learn parameters W to rank via $q^T W d$?

"search ranking" q^T

"information retrieval ranking" d

$$\begin{pmatrix} 1 & 0 & 0 & 1 & 0 \\ \text{se} & \text{in} & \text{re} & \text{ra} & \text{or} \end{pmatrix} \begin{pmatrix} 1 & 0.7 & 0.5 & 0 & 0 \\ 0.3 & 1 & 0.2 & 0 & 0 \\ 0.5 & 0.2 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0.7 \\ 0 & 0 & 0 & 0.7 & 1 \end{pmatrix} \begin{pmatrix} 0 \\ 1 \\ 1 \\ 1 \\ 0 \\ \text{se} \\ \text{in} \\ \text{re} = 2.2 \\ \text{ra} \\ \text{or (dering)} \end{pmatrix}$$

- Cf. Query translation models: Berger and Lafferty (1999)
- Problem is again sparsity – W is huge $> 10^{10}$

Is there a better way?

- Idea:
 - Can we learn a dense low-dimensional representation of a word in \mathbb{R}^d such that dot products $u^T v$ express word similarity?
 - We could still if we want to include a “translation” matrix between vocabularies (e.g., cross-language): $u^T W v$
 - But **now W is small!**
 - Supervised Semantic Indexing (Bai et al. *Journal of Information Retrieval* 2009) shows successful use of learning W for information retrieval
- But we'll develop direct similarity in this class

Distributional similarity based representations

- You can get a lot of value by representing a word by means of its neighbors
- “You shall know a word by the company it keeps”
 - (J. R. Firth 1957: 11)
- One of the most successful ideas of modern statistical NLP



...government debt problems turning into **banking** crises as happened in 2009...
...saying that Europe needs unified **banking** regulation to replace the hodgepodge...
...India has just given its **banking** system a shot in the arm...

↖ These words will represent *banking* ↗

Solution: Low dimensional vectors

- The number of topics that people talk about is small (in some sense)
 - Clothes, movies, politics, ...
- Idea: store “most” of the important information in a fixed, small number of dimensions: a dense vector
- Usually 25 – 1000 dimensions
- How to reduce the dimensionality?
 - Go from big, sparse co-occurrence count vector to low dimensional “word embedding”

Traditional Way:

Latent Semantic Indexing/Analysis

- Use Singular Value Decomposition (SVD) – kind of like Principal Components Analysis (PCA) for an arbitrary rectangular matrix – or just random projection to find a low-dimensional basis or orthogonal vectors
- Theory is that similarity is preserved as much as possible
- You can actually gain in IR (slightly) by doing LSA, as “noise” of term variation gets replaced by semantic “concepts”
- Somewhat popular in the 1990s [Deerwester et al. 1990, etc.]
 - But **results were always somewhat iffy** (... it worked sometimes)
 - Hard to implement efficiently in an IR system (dense vectors!)
- Discussed in *IIR* chapter 18, but not discussed further here
 - Not on the exam (!!!)

“NEURAL EMBEDDINGS”

Word meaning is defined in terms of vectors

- We will build a dense vector for each word type, chosen so that it is good at predicting other words appearing in its context
... those other words also being represented by vectors ... it all gets a bit recursive

$$\mathit{banking} = \begin{pmatrix} 0.286 \\ 0.792 \\ -0.177 \\ -0.107 \\ 0.109 \\ -0.542 \\ 0.349 \\ 0.271 \end{pmatrix}$$

Neural word embeddings - visualization



Basic idea of learning neural network word embeddings

- We define a model that aims to predict between a center word w_t and context words in terms of word vectors
- $p(\text{context} | w_t) = \dots$
- which has a loss function, e.g.,
- $J = 1 - p(w_{-t} | w_t)$
- We look at many positions t in a big language corpus
- We keep adjusting the vector representations of words to minimize this loss

Idea: Directly learn low-dimensional word vectors based on ability to predict

- Old idea: Learning representations by back-propagating errors. (Rumelhart et al., 1986)
- A neural probabilistic language model (Bengio et al., 2003)
- NLP (almost) from Scratch (Collobert & Weston, 2008)
- A recent, even simpler and faster model: word2vec (Mikolov et al. 2013) → intro now
- The GloVe model from Stanford (Pennington, Socher, and Manning 2014) connects back to matrix factorization
- Per-token representations: Deep contextual word representations: ELMo, ULMfit, **BERT**

Non-linear
and slow

Fast
bilinear
models

Current
state of
the art

Word2vec is a family of algorithms

[Mikolov et al. 2013]

Predict between every word and its context words!

Two algorithms

- 1. Skip-grams (SG)**

Predict context words given target (position independent)

- 2. Continuous Bag of Words (CBOW)**

Predict target word from bag-of-words context

Two (moderately efficient) training methods

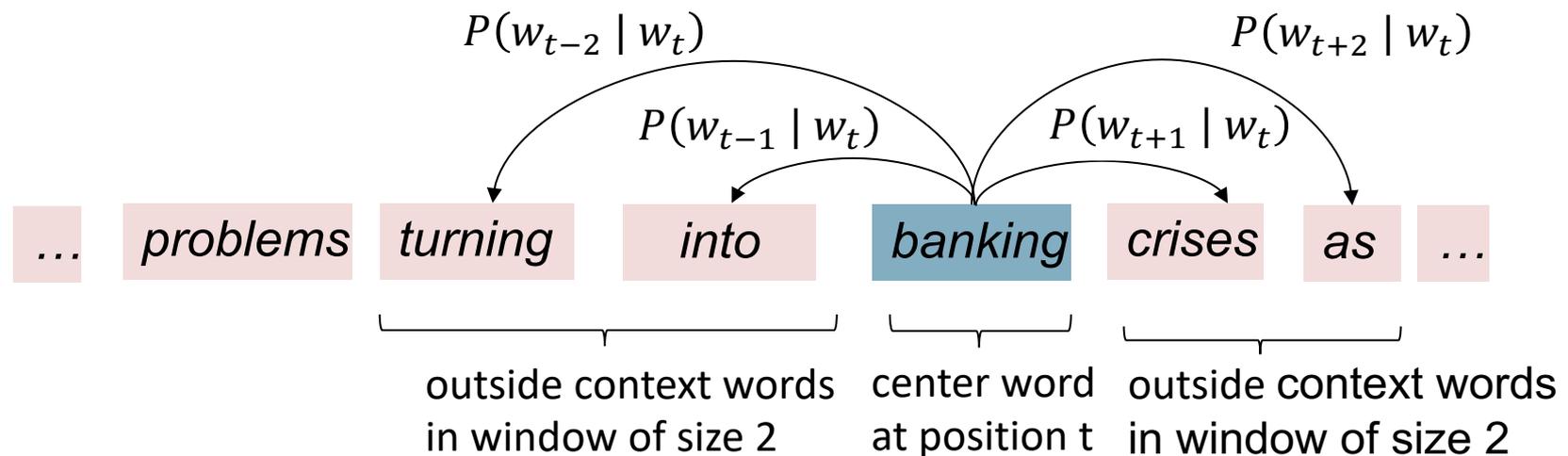
- 1. Hierarchical softmax**

- 2. Negative sampling**

- 3. Naïve softmax**

Word2Vec Skip-gram Overview

- Example windows and process for computing $P(w_{t+j} | w_t)$



Word2vec: objective function

For each position $t = 1, \dots, T$, predict context words within a window of fixed size m , given center word w_j .

Likelihood =

$$L(\theta) = \prod_{t=1}^T \prod_{\substack{-m \leq j \leq m \\ j \neq 0}} P(w_{t+j} | w_t; \theta)$$

θ is all variables
to be optimized

sometimes called *cost* or *loss* function

The **objective function** $J(\theta)$ is the (average) negative log likelihood:

$$J(\theta) = -\frac{1}{T} \log L(\theta) = -\frac{1}{T} \sum_{t=1}^T \sum_{\substack{-m \leq j \leq m \\ j \neq 0}} \log P(w_{t+j} | w_t; \theta)$$

Minimizing objective function \Leftrightarrow Maximizing predictive accuracy

Word2vec: objective function

- We want to minimize the objective function:

$$J(\theta) = -\frac{1}{T} \sum_{t=1}^T \sum_{\substack{-m \leq j \leq m \\ j \neq 0}} \log P(w_{t+j} | w_t; \theta)$$

- Question: How to calculate $P(w_{t+j} | w_t; \theta)$?
- Answer: We will use two vectors per word w :
 - v_w when w is a center word
 - u_w when w is a context word
- Then for a center word c and a context word o :

$$P(o|c) = \frac{\exp(u_o^T v_c)}{\sum_{w \in V} \exp(u_w^T v_c)}$$

Word2vec: prediction function

Exponentiation makes anything positive

Dot product compares similarity of o and c .

$$u^T v = u \cdot v = \sum_{i=1}^n u_i v_i$$

Larger dot product = larger probability

$$P(o|c) = \frac{\exp(u_o^T v_c)}{\sum_{w \in V} \exp(u_w^T v_c)}$$

Normalize over entire vocabulary to give probability distribution

- This is an example of the **softmax function** $\mathbb{R}^n \rightarrow (0,1)^n$

$$\text{softmax}(x_i) = \frac{\exp(x_i)}{\sum_{j=1}^n \exp(x_j)} = p_i$$

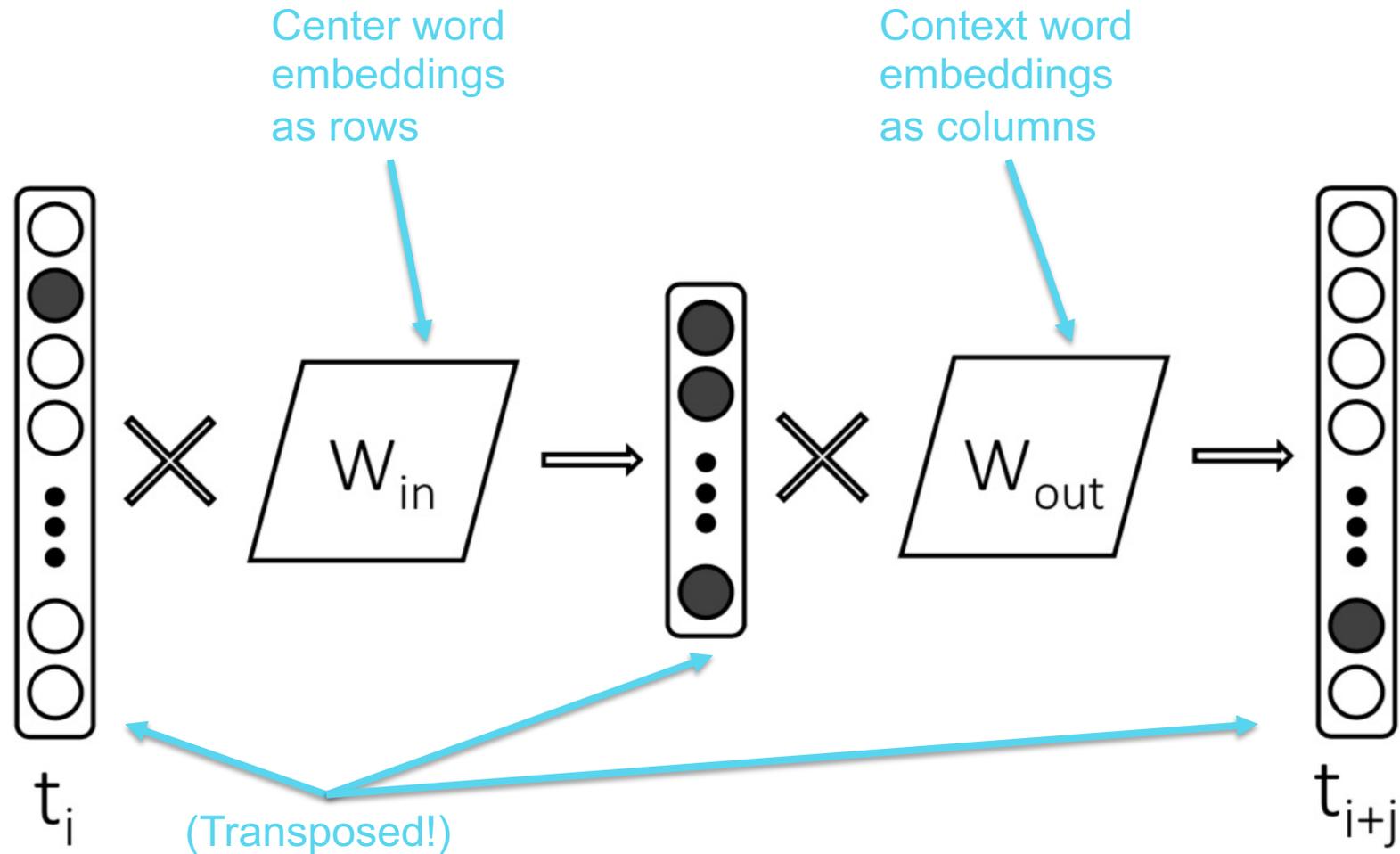
Open region

- The softmax function maps arbitrary values x_i to a probability distribution

p_i

- “max” because amplifies probability of largest x_i
- “soft” because still assigns some probability to smaller x_i
- Frequently used in neural networks/Deep Learning

Word2vec: 2 matrices of parameters



To learn good word vectors: Compute **all** vector gradients!

- We often define the set of **all** parameters in a model in terms of one long vector θ

- In our case with d -dimensional vectors and V many words:

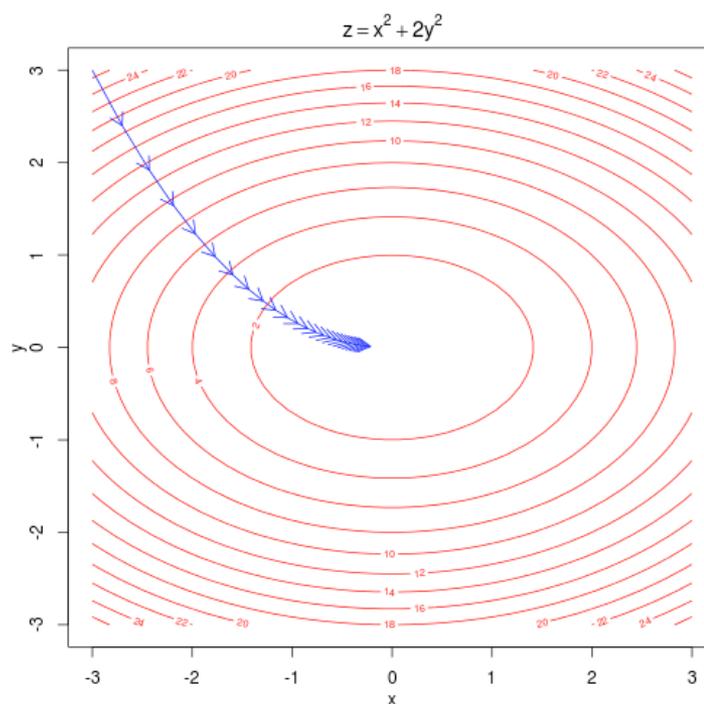
$$\theta = \begin{bmatrix} v_{aardvark} \\ v_a \\ \vdots \\ v_{zebra} \\ u_{aardvark} \\ u_a \\ \vdots \\ u_{zebra} \end{bmatrix} \in \mathbb{R}^{2dV}$$

- We then optimize these parameters

Note: Every word has two vectors! Makes it simpler!

Intuition of how to minimize loss for a simple function over two parameters

We start at a random point and walk in the steepest direction, which is given by the derivative of the function



Contour lines show points of equal value of objective function

Descending by using derivatives

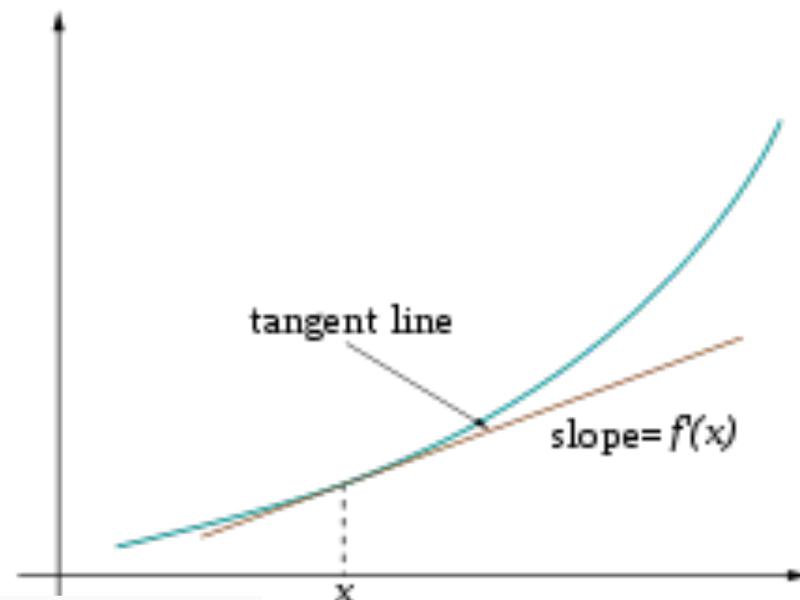
We will minimize a cost function by gradient descent

Trivial example: (from Wikipedia)

Find a local minimum of the function

$$f(x) = x^4 - 3x^3 + 2,$$

with derivative $f'(x) = 4x^3 - 9x^2$



```
x_old = 0
x_new = 6 # The algorithm starts at x=6
eps = 0.01 # step size
precision = 0.00001

def f_derivative(x):
    return 4 * x**3 - 9 * x**2

while abs(x_new - x_old) > precision:
    x_old = x_new
    x_new = x_old - eps * f_derivative(x_old)

print("Local minimum occurs at", x_new)
```

Subtracting a fraction of the gradient moves you towards the minimum!

Vanilla Gradient Descent Code

$$\theta^{new} = \theta^{old} - \alpha \nabla_{\theta} J(\theta)$$

```
while True:  
    theta_grad = evaluate_gradient(J, corpus, theta)  
    theta = theta - alpha * theta_grad
```

Stochastic Gradient Descent

- But Corpus may have 40B tokens and windows
- You would wait a very long time before making a single update!
- **Very** bad idea for pretty much all neural nets!
- Instead: We update parameters after each window t
→ Stochastic gradient descent (SGD)

$$\theta^{new} = \theta^{old} - \alpha \nabla_{\theta} J_t(\theta)$$

```
while True:  
    window = sample_window(corpus)  
    theta_grad = evaluate_gradient(J, window, theta)  
    theta = theta - alpha * theta_grad
```

Working out how to optimize a neural network is really all the chain rule!

Chain rule! If $y = f(u)$ and $u = g(x)$, i.e. $y = f(g(x))$, then:

$$\frac{dy}{dx} = \frac{dy}{du} \frac{du}{dx} = \frac{df(u)}{du} \frac{dg(x)}{dx}$$

Simple example: $\frac{dy}{dx} = \frac{d}{dx} 5(x^3 + 7)^4$

$$y = f(u) = 5u^4$$

$$u = g(x) = x^3 + 7$$

$$\frac{dy}{du} = 20u^3$$

$$\frac{du}{dx} = 3x^2$$

$$\frac{dy}{dx} = 20(x^3 + 7)^3 \cdot 3x^2$$

Objective Function

$$\text{Maximize } J'(\theta) = \prod_{t=1}^T \prod_{\substack{-m \leq j \leq m \\ j \neq 0}} p(w'_{t+j} | w_t; \theta)$$

Or minimize
neg. log
likelihood

$$J(\theta) = -\frac{1}{T} \sum_{t=1}^T \sum_{\substack{-m \leq j \leq m \\ j \neq 0}} \log p(w'_{t+j} | w_t)$$

[negate to minimize;
log is monotone]

↑
text
length

↑
window
size

where

$$p(o|c) = \frac{\exp(u_o^T v_c)}{\sum_{w=1}^V \exp(u_w^T v_c)}$$

word IDs ↗

We now take derivatives to work out minimum

Each word type
(vocab entry)
has two word
representations:
as center word
and context word

$$\frac{\partial}{\partial v_c} \log \frac{\exp(u_0^T v_c)}{\sum_{w=1}^V \exp(u_w^T v_c)}$$

$$= \underbrace{\frac{\partial}{\partial v_c} \log \exp(u_0^T v_c)}_{\textcircled{1}} - \underbrace{\frac{\partial}{\partial v_c} \log \sum_{w=1}^V \exp(u_w^T v_c)}_{\textcircled{2}}$$

$$\textcircled{1} \quad \frac{\partial}{\partial v_c} \underbrace{\log \exp(u_0^T v_c)}_{\text{inverses}} = \frac{\partial}{\partial v_c} u_0^T v_c = u_0$$

Vector!
Not high school
single variable
calculus

You can do things one variable at a time, and this may be helpful when things get gnarly.

$$\forall j \quad \frac{\partial}{\partial (v_c)_j} u_0^T v_c = \frac{\partial}{\partial (v_c)_j} \sum_{i=1}^d (u_0)_i (v_c)_i = (u_0)_j$$

Each term is zero except when $i=j$

$$\textcircled{2} \frac{\partial}{\partial v_c} \log \underbrace{\sum_{w=1}^v \exp(u_w^T v_c)}_f$$

$z = g(v_c)$

$$= \frac{1}{\sum_{w=1}^v \exp(u_w^T v_c)}$$

$$\frac{\partial}{\partial v_c} f(g(v_c)) = \frac{\partial f}{\partial z} \cdot$$

$$= \frac{1}{\sum_{w=1}^v \exp(u_w^T v_c)}$$

$$\frac{\partial}{\partial v_c} \sum_{x=1}^v \exp(u_x^T v_c)$$

Important to change index

$$\frac{\partial z}{\partial v_c} \quad \text{Use chain rule}$$

$$\left(\sum_{x=1}^v \frac{\partial}{\partial v_c} \underbrace{\exp(u_x^T v_c)}_f \right)$$

$z = g(v_c)$

Move deriv inside sum

$$\left(\sum_{x=1}^v \exp(u_x^T v_c) \frac{\partial}{\partial v_c} u_x^T v_c \right)$$

Chain rule

$$\left(\sum_{x=1}^v \exp(u_x^T v_c) u_x \right)$$

$$\frac{\partial}{\partial v_c} \log(p(o|c)) = u_o - \frac{1}{\sum_{w=1}^V \exp(u_w^T v_c)} \cdot \left(\sum_{x=1}^V \exp(u_x^T v_c) u_x \right)$$

$$= u_o - \sum_{x=1}^V \frac{\exp(u_x^T v_c)}{\sum_{w=1}^V \exp(u_w^T v_c)} u_x$$

Distribute
term
across sum

$$= u_o - \underbrace{\sum_{x=1}^V p(x|c)}_{\text{this an expectation: average over all context vectors weighted by their probability}} u_x$$

= observed - expected

This is just the derivatives for the center vector parameters
Also need derivatives for output vector parameters
(they're similar)

Then we have derivative w.r.t. all parameters and can minimize

Linear Relationships in word2vec

These representations are *very good* at encoding **similarity** and **dimensions of similarity**!

- Analogies testing dimensions of similarity can be solved quite well just by doing vector subtraction in the embedding space

Syntactically

- $x_{apple} - x_{apples} \approx x_{car} - x_{cars} \approx x_{family} - x_{families}$
- Similarly for verb and adjective morphological forms

Semantically (Semeval 2012 task 2)

- $x_{shirt} - x_{clothing} \approx x_{chair} - x_{furniture}$
- $x_{king} - x_{man} \approx x_{queen} - x_{woman}$

Word Analogies

Test for linear relationships, examined by Mikolov et al.

a:b :: c:?



$$d = \arg \max_x \frac{(w_b - w_a + w_c)^T w_x}{\|w_b - w_a + w_c\|}$$

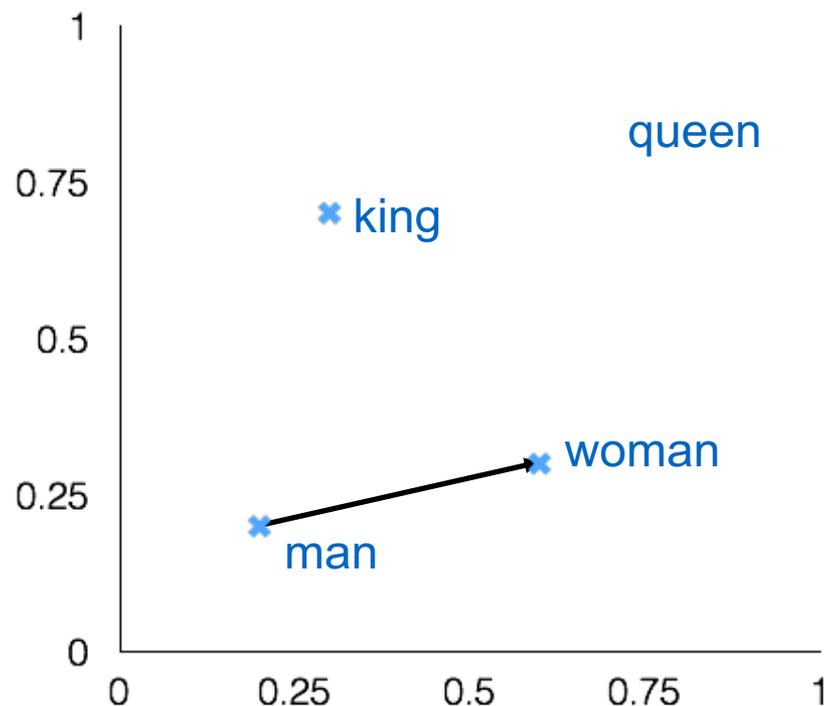
man:woman :: king:?

+ king [0.30 0.70]

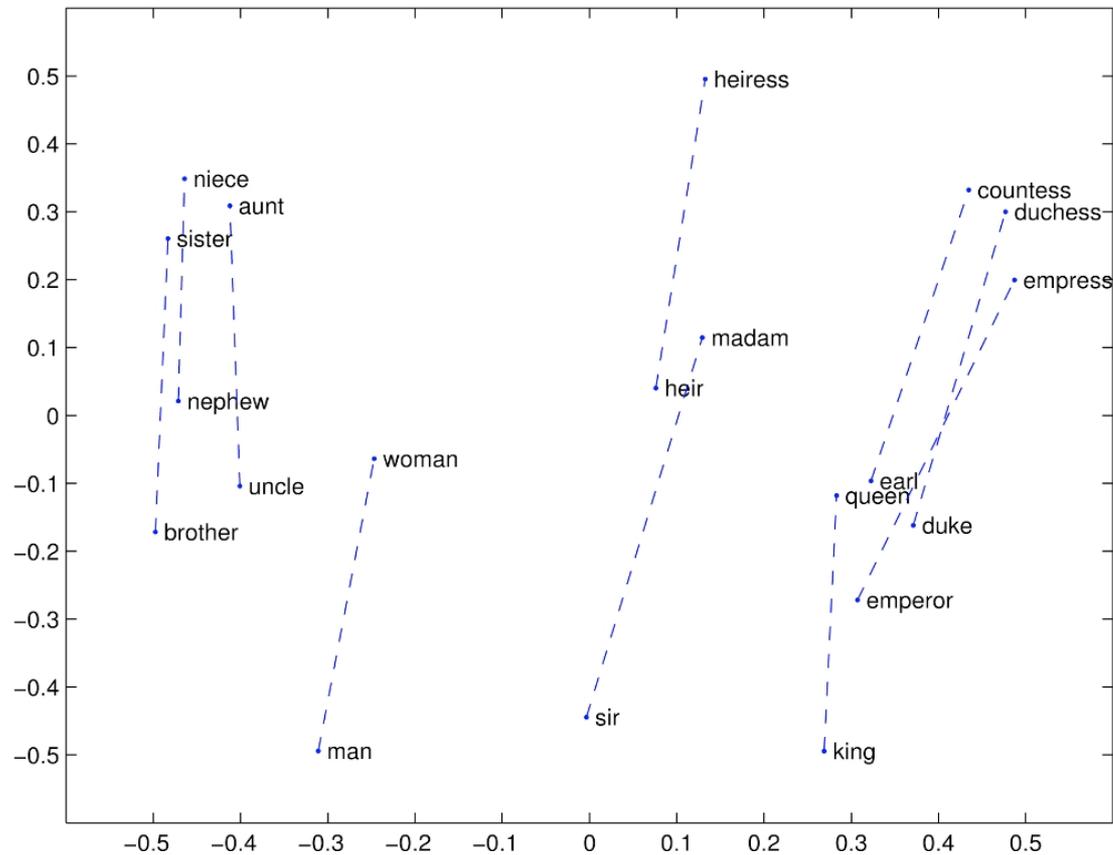
- man [0.20 0.20]

+ woman [0.60 0.30]

queen [0.70 0.80]

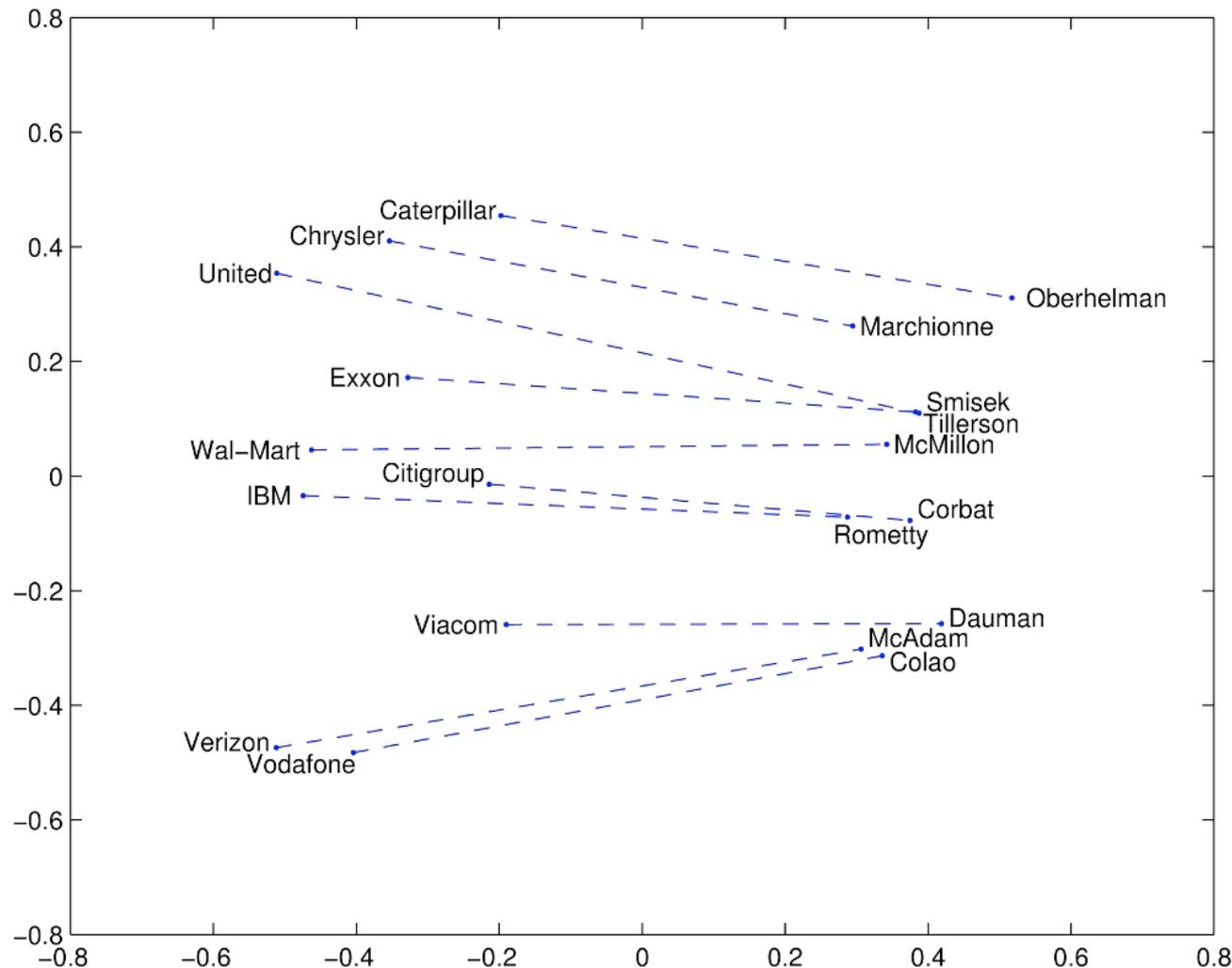


GloVe Visualizations

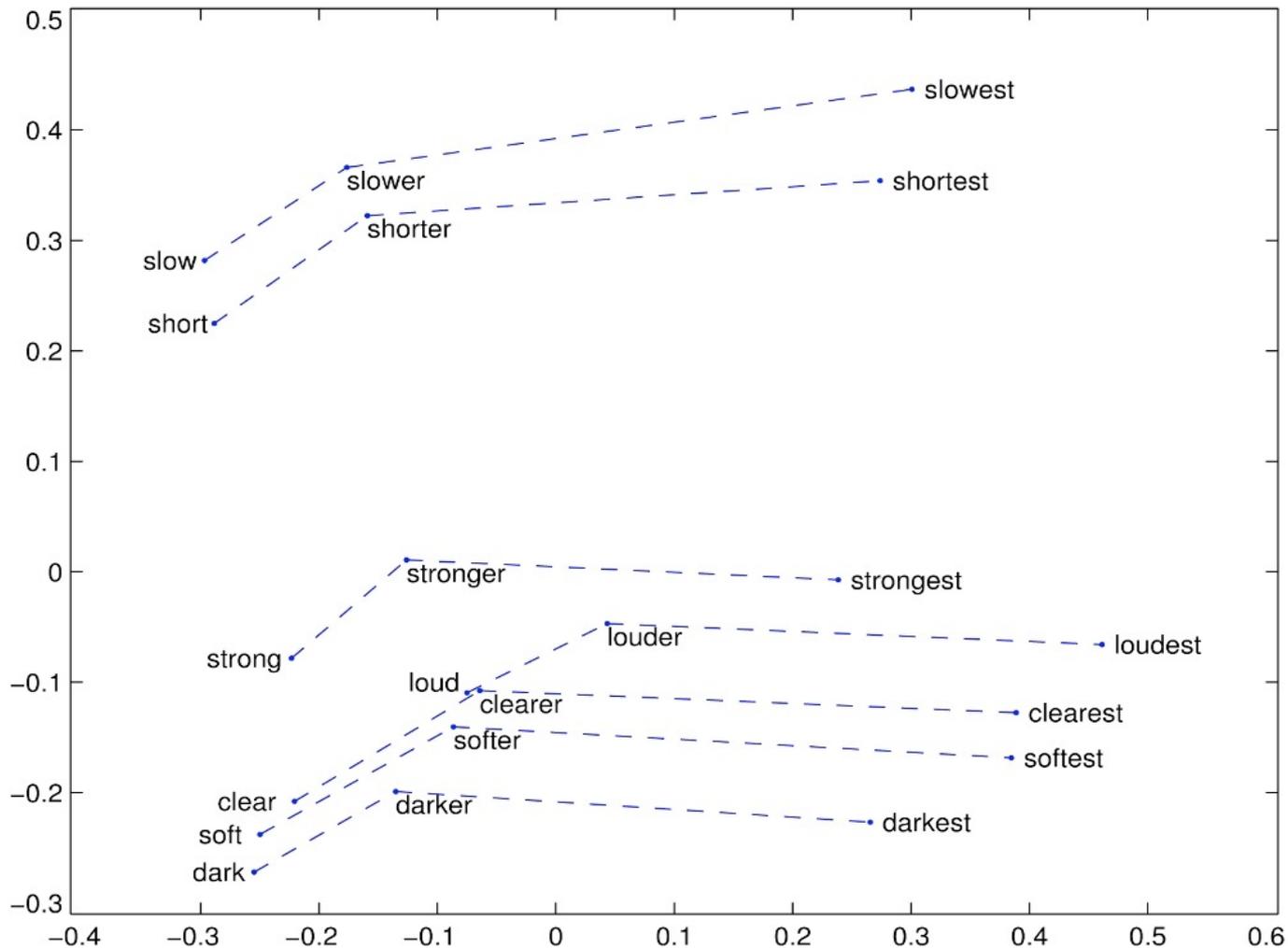


<http://nlp.stanford.edu/projects/glove/>

Glove Visualizations: Company - CEO



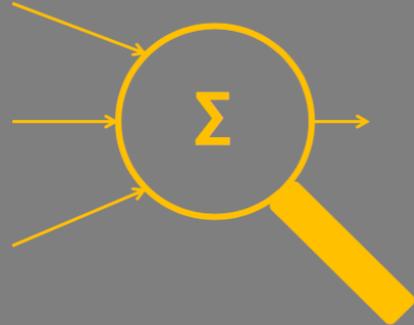
Glove Visualizations: Superlatives



Application to Information Retrieval

Application is just beginning – we’re “at the end of the early years”

- Google’s RankBrain – little is publicly known
 - Bloomberg article by Jack Clark (Oct 26, 2015):
<http://www.bloomberg.com/news/articles/2015-10-26/google-turning-its-lucrative-web-search-over-to-ai-machines>
 - A result reranking system. “3rd most valuable ranking signal”
 - But note: more of the potential value is in the tail?
- New SIGIR Neu-IR workshop series (2016 on)



Neu-IR (2016)
The Neural Information Retrieval Workshop @ SIGIR
Pisa, Tuscany, Italy on 21st July, 2016
research.microsoft.com/neuir2016

An application to information retrieval

Nalisnick, Mitra, Craswell & Caruana. 2016. Improving Document Ranking with Dual Word Embeddings. *WWW 2016 Companion*.

<http://research.microsoft.com/pubs/260867/pp1291-Nalisnick.pdf>

Mitra, Nalisnick, Craswell & Caruana. 2016. A Dual Embedding Space Model for Document Ranking. [arXiv:1602.01137](https://arxiv.org/abs/1602.01137) [cs.IR]

Builds on BM25 model idea of “aboutness”

- Not just term repetition indicating aboutness
- Relationship between query terms and *all* terms in the document indicates aboutness (BM25 uses only query terms)

Makes clever argument for different use of word and context vectors in word2vec’s CBOW/SGNS or GloVe

Modeling document aboutness:

Results from a search for Albuquerque

 d_1

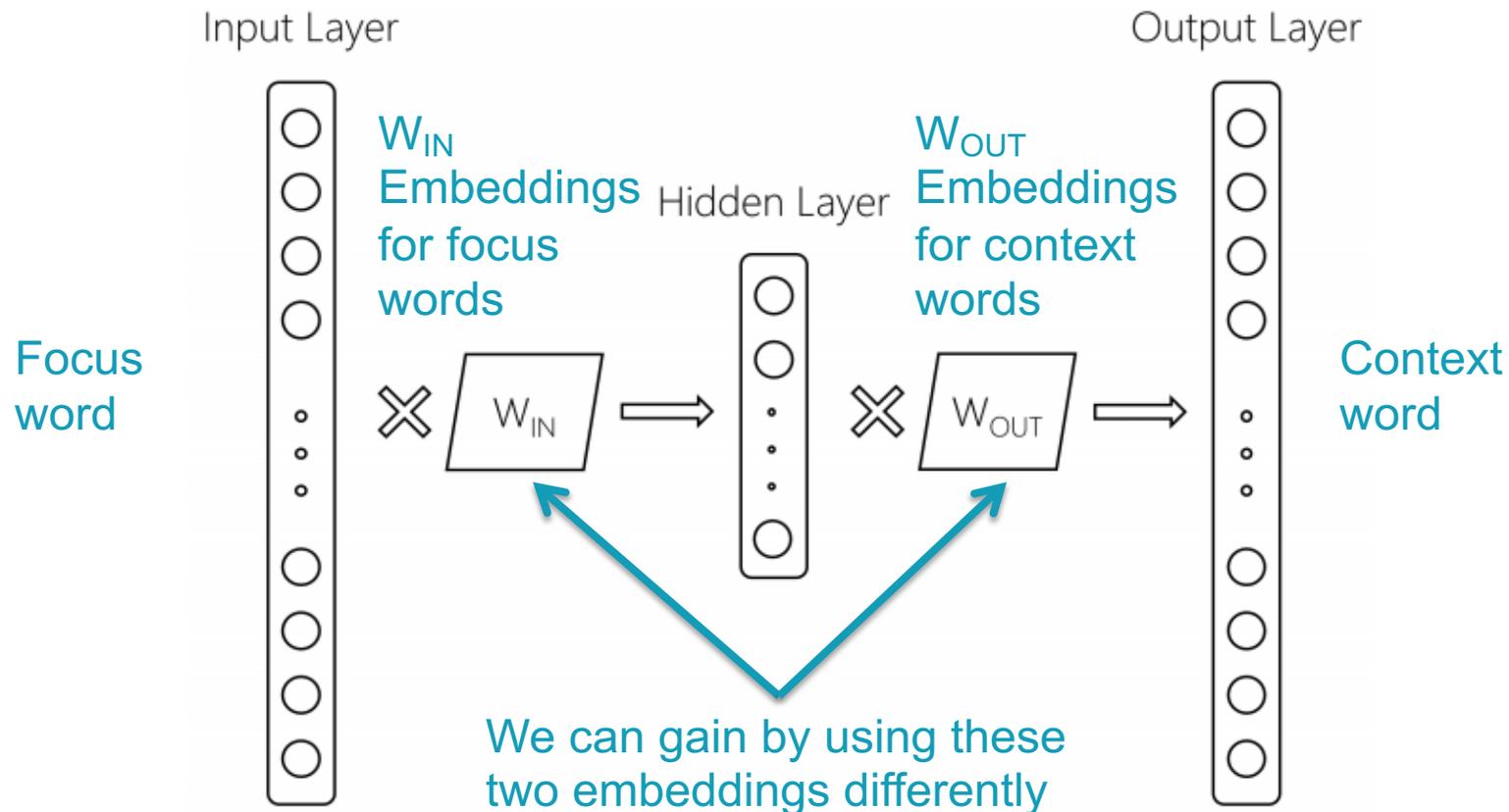
Allen suggested that they could program a BASIC interpreter for the device; after a call from Gates claiming to have a working interpreter, MITS requested a demonstration. Since they didn't actually have one, Allen worked on a simulator for the Altair while Gates developed the interpreter. Although they developed the interpreter on a simulator and not the actual device, the interpreter worked flawlessly when they demonstrated the interpreter to MITS in Albuquerque, New Mexico in March 1975; MITS agreed to distribute it, marketing it as Altair BASIC.

 d_2

Albuquerque is the most populous city in the U.S. state of New Mexico. The high-altitude city serves as the county seat of Bernalillo County, and it is situated in the central part of the state, straddling the Rio Grande. The city population is 557,169 as of the July 1, 2014, population estimate from the United States Census Bureau, and ranks as the 32nd-largest city in the U.S. The Metropolitan Statistical Area (or MSA) has a population of 902,797 according to the United States Census Bureau's most recently available estimate for July 1, 2013.

Using 2 word embeddings

word2vec model with 1 word of context



Using 2 word embeddings

yale		seahawks	
IN-IN	IN-OUT	IN-IN	IN-OUT
yale	yale	seahawks	seahawks
harvard	faculty	49ers	highlights
nyu	alumni	broncos	jerseys
cornell	orientation	packers	tshirts
tulane	haven	nfl	seattle
tufts	graduate	steelers	hats

Dual Embedding Space Model (DESM)

- Simple model
- A document is represented by the centroid of its word vectors

$$\bar{\mathbf{D}} = \frac{1}{|D|} \sum_{\mathbf{d}_j \in D} \frac{\mathbf{d}_j}{\|\mathbf{d}_j\|}$$

- Query-document similarity is average over query words of cosine similarity

$$DESM(Q, D) = \frac{1}{|Q|} \sum_{q_i \in Q} \frac{\mathbf{q}_i^T \bar{\mathbf{D}}}{\|\mathbf{q}_i\| \|\bar{\mathbf{D}}\|}$$

Dual Embedding Space Model (DESM)

- What works best is to use the OUT vectors for the document and the IN vectors for the query

$$DESM_{IN-OUT}(Q, D) = \frac{1}{|Q|} \sum_{q_i \in Q} \frac{q_{IN,i}^T \overline{D_{OUT}}}{\|q_{IN,i}\| \|\overline{D_{OUT}}\|}$$

- This way similarity measures *aboutness* – words that appear with this word – which is more useful in this context than *(distributional) semantic similarity*

Experiments

- Train word2vec from either
 - 600 million Bing queries
 - 342 million web document sentences
- Test on 7,741 randomly sampled Bing queries
 - 5 level eval (Perfect, Excellent, Good, Fair, Bad)
- Two approaches
 1. Use DESM model to rerank top results from BM25
 2. Use DESM alone or a mixture model of it and BM25

$$MM(Q, D) = \alpha DESM(Q, D) + (1 - \alpha) BM25(Q, D)$$

$$\alpha \in \mathbb{R}, 0 \leq \alpha \leq 1$$

Results – reranking k -best list

	Explicitly Judged Test Set		
	NDCG@1	NDCG@3	NDCG@10
BM25	23.69	29.14	44.77
LSA	22.41*	28.25*	44.24*
DESM (IN-IN, trained on body text)	23.59	29.59	45.51*
DESM (IN-IN, trained on queries)	23.75	29.72	46.36*
DESM (IN-OUT, trained on body text)	24.06	30.32*	46.57*
DESM (IN-OUT, trained on queries)	25.02*	31.14*	47.89*

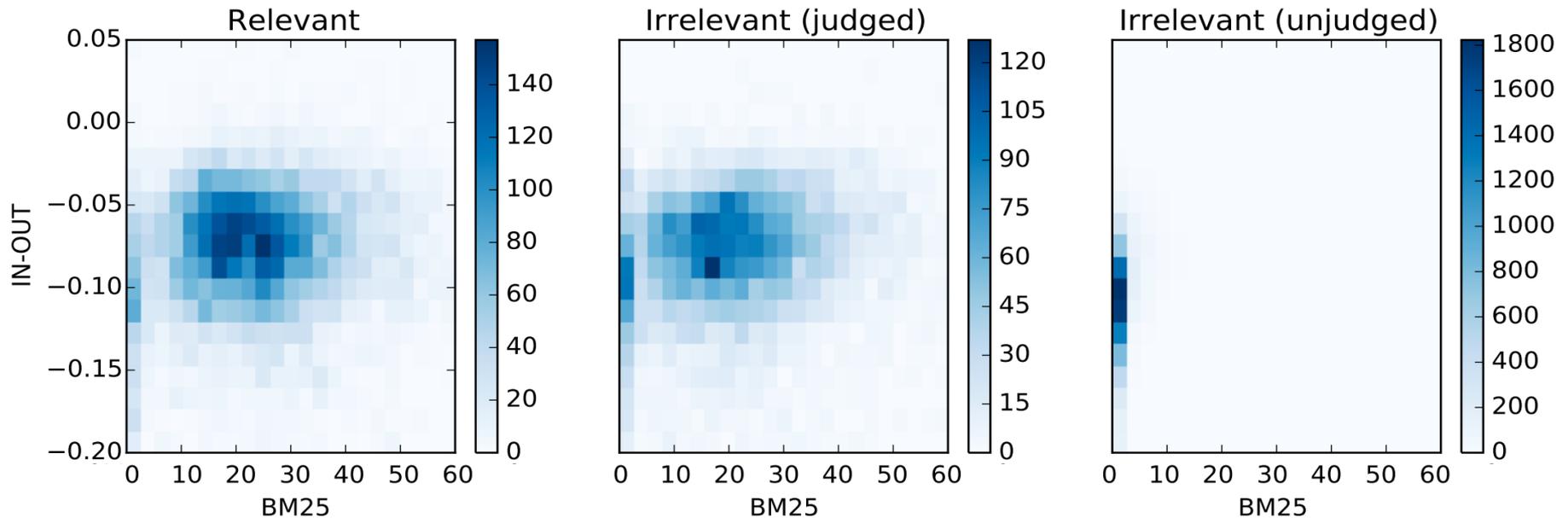
Pretty decent gains – e.g., 2% for NDCG@3

Gains are bigger for model trained on queries than docs

Results – whole ranking system

	Explicitly Judged Test Set		
	NDCG@1	NDCG@3	NDCG@10
BM25	21.44	26.09	37.53
LSA	04.61*	04.63*	04.83*
DESM (IN-IN, trained on body text)	06.69*	06.80*	07.39*
DESM (IN-IN, trained on queries)	05.56*	05.59*	06.03*
DESM (IN-OUT, trained on body text)	01.01*	01.16*	01.58*
DESM (IN-OUT, trained on queries)	00.62*	00.58*	00.81*
BM25 + DESM (IN-IN, trained on body text)	21.53	26.16	37.48
BM25 + DESM (IN-IN, trained on queries)	21.58	26.20	37.62
BM25 + DESM (IN-OUT, trained on body text)	21.47	26.18	37.55
BM25 + DESM (IN-OUT, trained on queries)	21.54	26.42*	37.86*

A possible explanation



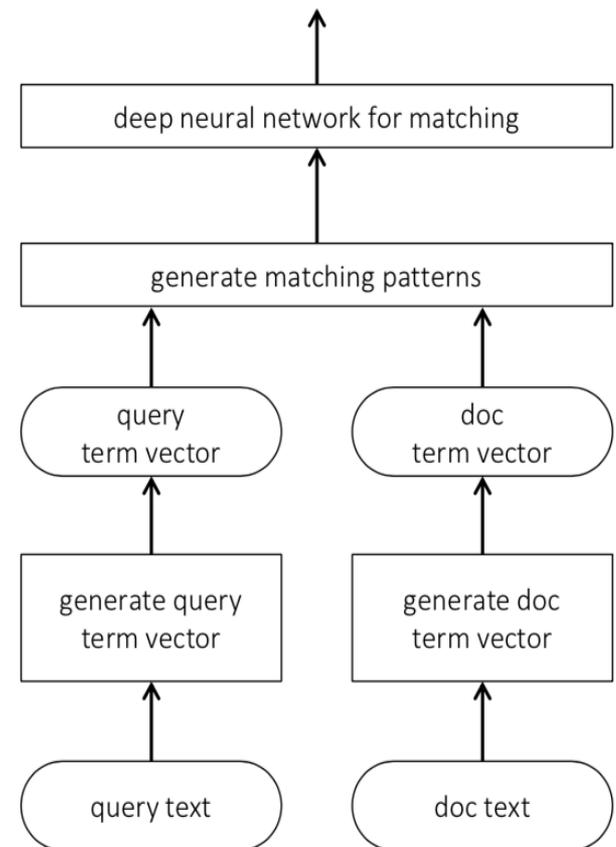
IN-OUT has some ability to prefer Relevant to close-by (judged) non-relevant, but it's scores induce too much noise vs. BM25 to be usable alone

DESM conclusions

- DESM is a weak ranker but effective at finding subtler similarities/aboutness
- It is effective at, but only at, reranking at least somewhat relevant documents
 - For example, DESM can confuse Oxford and Cambridge
 - Bing rarely makes an Oxford/Cambridge mistake!

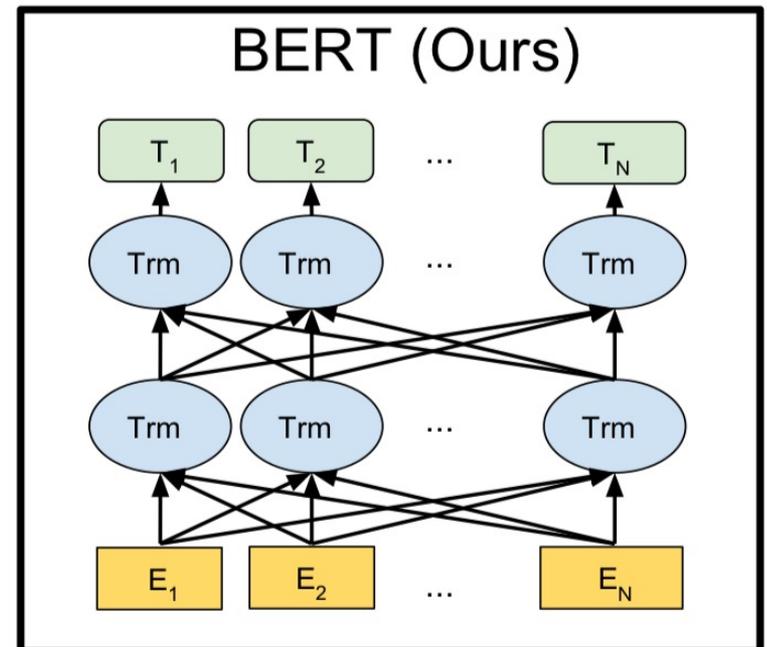
What else can neural nets do in IR?

- Use a neural network as a supervised reranker
- Assume a query and document embedding network (as we have discussed)
- Assume you have (q,d,rel) relevance data
- Learn a neural network (with supervised learning) to predict relevance of (q,d) pair
- An example of “machine-learned relevance”, which we’ll talk about more next lecture

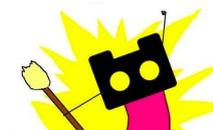


What else can neural nets do in IR?

- BERT: Devlin, Chang, Lee, Toutanova (2018)
- A deep transformer-based neural network
- Builds per-token (in context) representations
- Produces a query/document representation as well
- Or jointly embed query and document and ask for a retrieval score
- Incredibly effective!
- <https://arxiv.org/abs/1810.04805>



Summary: Embed all the things!



Word embeddings are the hot new technology (again!)

Lots of applications wherever knowing word context or similarity helps prediction:

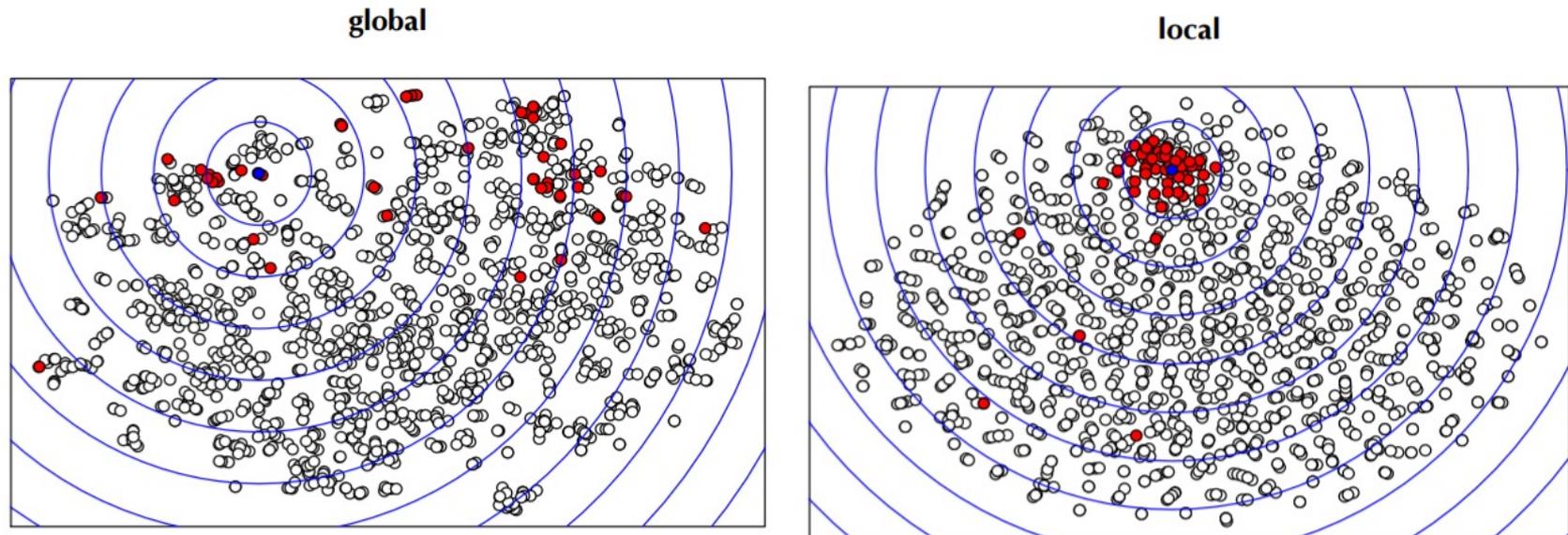
- Synonym handling in search
- Document aboutness
- Ad serving
- Language models: from spelling correction to email response
- Machine translation
- Sentiment analysis
- ...

Global vs. local embedding [Diaz 2016]

global	local
cutting	tax
squeeze	deficit
reduce	vote
slash	budget
reduction	reduction
spend	house
lower	bill
halve	plan
soften	spend
freeze	billion

Figure 3: Terms similar to ‘cut’ for a word2vec model trained on a general news corpus and another trained only on documents related to ‘gasoline tax’.

Global vs. local embedding [Diaz 2016]



Train w2v on documents from first round of retrieval

Fine-grained word sense disambiguation

Figure 5: Global versus local embedding of highly relevant terms. Each point represents a candidate expansion term. Red points have high frequency in the relevant set of documents. White points have low or no frequency in the relevant set of documents. The blue point represents the query. Contours indicate distance from the query.

