

Information Retrieval

CS 547/DS 547

Worcester Polytechnic Institute
Department of Computer Science
Instructor: Prof. Kyumin Lee

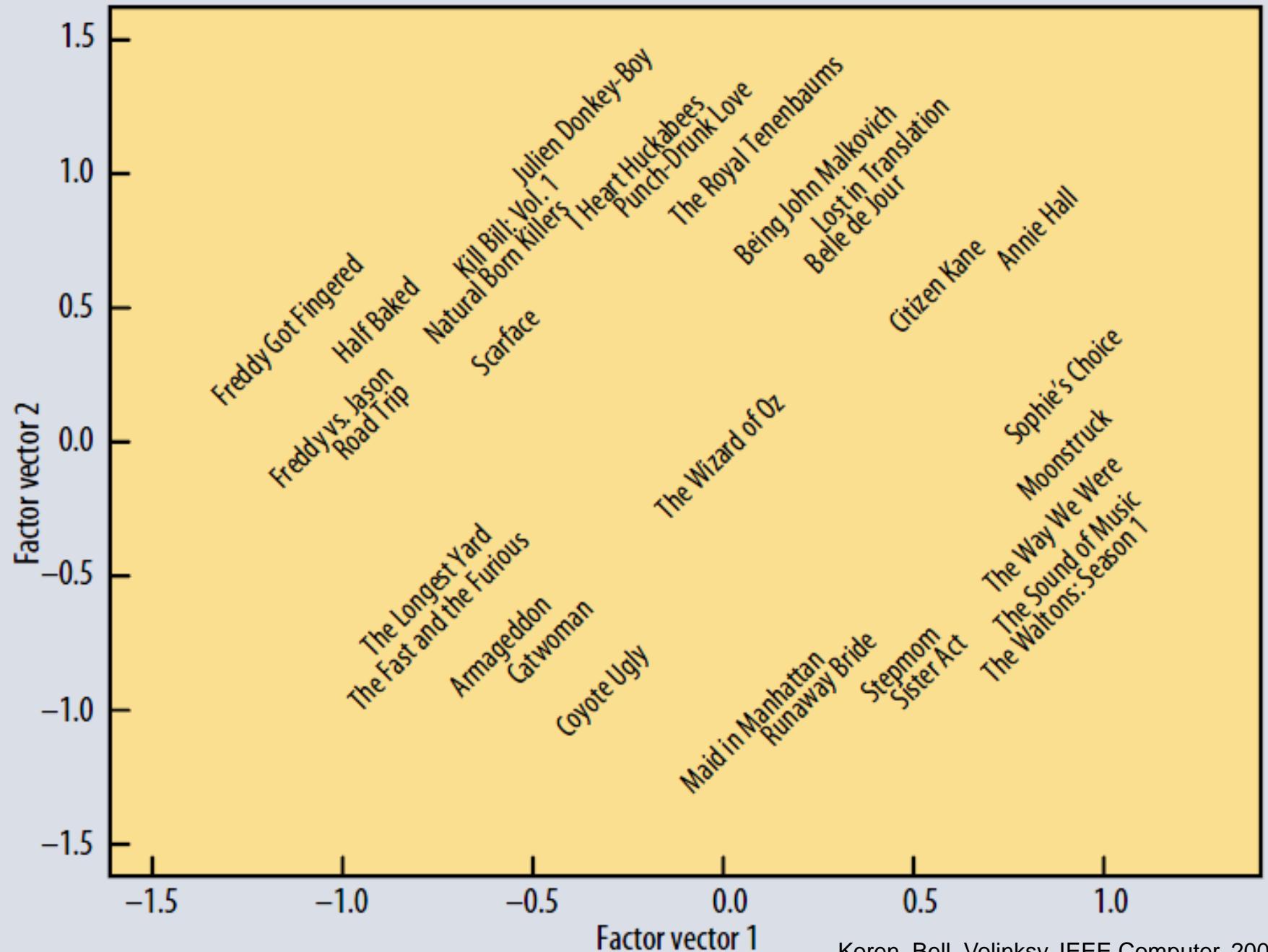
Project Team

1. Bo Yu, Jin Yang, Kangjian Wu
2. Vignesh Sundaram, Amey More, Akanksha Pawar, Padmesh Naik
3. Xiaofan Zhou, Yiming Liu, Yuyuan Liu, Akhil Daphara
4. Sidhant Jain, Parth dhruv, Darshan Swami, Aditya Ramesh
5. Chandra Rachabathuni, Ayush Shinde, Chao Wang, Anthony Chen, Adhiraj Budukh
6. Joe Scheufele, Alex Alvarez, Matt Suyer, Jason Dykstra
7. Adityavikram Gurao, Snehith Varma Datla, Sree Likhith Dasari, Supreeth Bandi
8. Samar, Bao, Michael, Megan

HW4

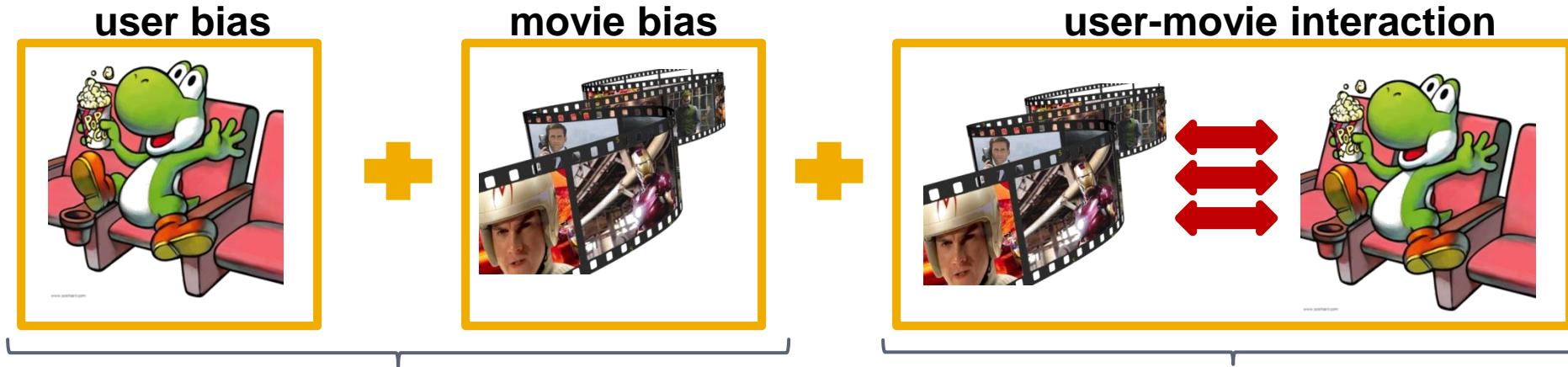
- https://canvas.wpi.edu/courses/46542/assignments/292272?module_item_id=917255
- Due date is March 31

Recommenders



Extending Latent Factor Model to Include Biases

Modeling Biases and Interactions



Baseline predictor

- Separates users and movies
- Benefits from insights into user's behavior
- Among the main practical contributions of the competition

User-Movie interaction

- Characterizes the matching between users and movies
- Attracts most research in the field
- Benefits from algorithmic and mathematical innovations

- μ = overall mean rating
- b_x = bias of user x
- b_i = bias of movie i

Putting It All Together

$$r_{xi} = \mu + b_x + b_i + q_i \cdot p_x$$

Overall mean rating Bias for user x Bias for movie i User-Movie interaction

■ Example:

- Mean rating: $\mu = 3.7$
- You are a critical reviewer: your ratings are 1 star lower than the mean: $b_x = -1$
- Star Wars gets a mean rating of 0.5 higher than average movie: $b_i = +0.5$
- Predicted rating for you on Star Wars:
 $= 3.7 - 1 + 0.5 = 3.2$

Fitting the New Model

- **Solve:**

$$\min_{Q,P} \sum_{(x,i) \in R} \left(r_{xi} - (\mu + b_x + b_i + q_i p_x) \right)^2$$

goodness of fit

$$+ \left(\lambda_1 \sum_i \|q_i\|^2 + \lambda_2 \sum_x \|p_x\|^2 + \lambda_3 \sum_x \|b_x\|^2 + \lambda_4 \sum_i \|b_i\|^2 \right)$$

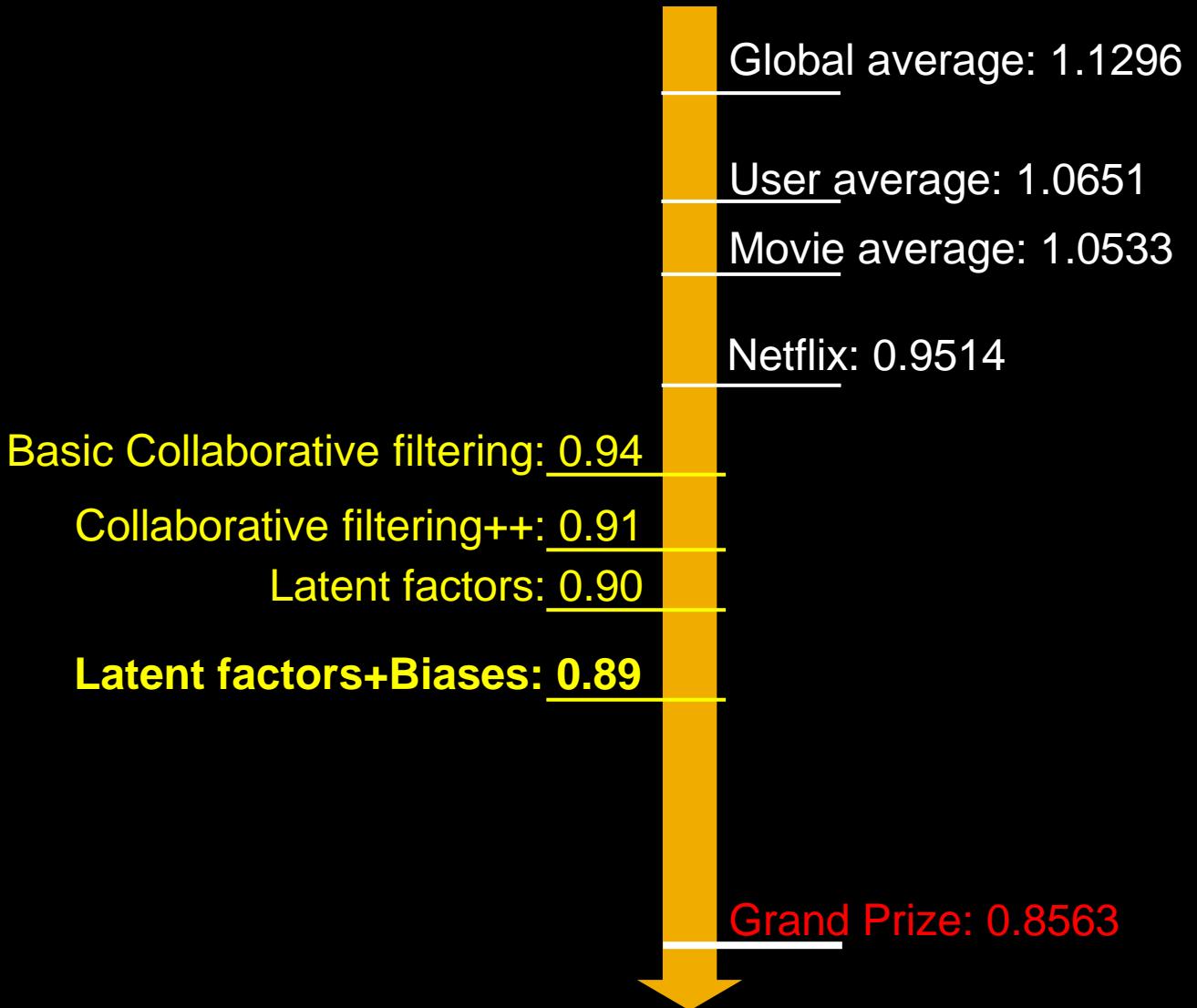
regularization

λ is selected via grid-search on a validation set

- **Stochastic gradient decent to find parameters**

- **Note:** Both biases b_x, b_i as well as interactions q_i, p_x are treated as parameters (and we learn them)

Performance of Various Methods

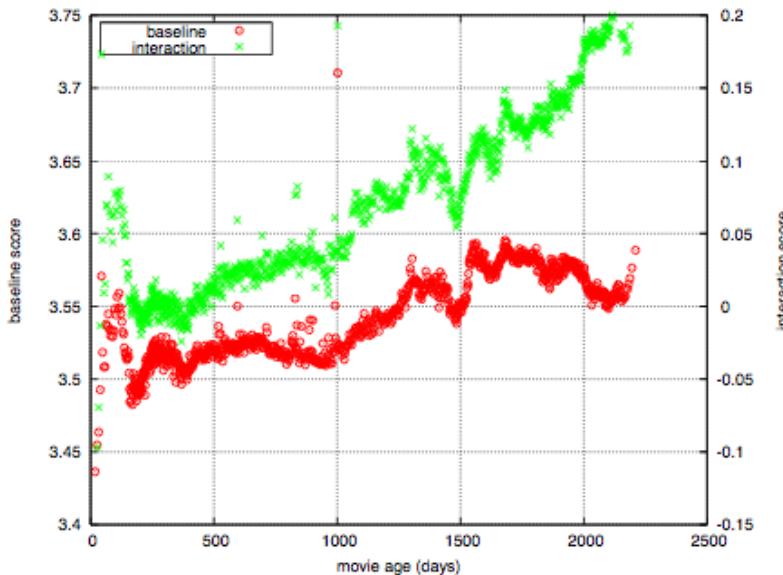
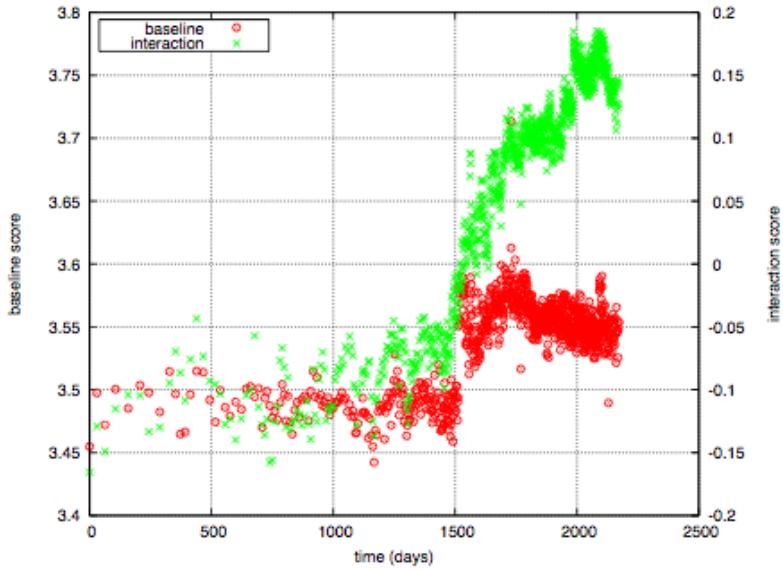


The Netflix Challenge: 2006-09

Temporal Biases Of Users

- **Sudden rise in the average movie rating (early 2004)**
 - Improvements in Netflix
 - GUI improvements
 - Meaning of rating changed
- **Movie age**
 - Users prefer new movies without any reasons
 - Older movies are just inherently better than newer ones

Y. Koren, Collaborative filtering with temporal dynamics, KDD '09



Temporal Biases & Factors

- Original model:

$$r_{xi} = \mu + b_x + b_i + q_i \cdot p_x$$

- Add time dependence to biases:

$$r_{xi} = \mu + b_x(t) + b_i(t) + q_i \cdot p_x$$

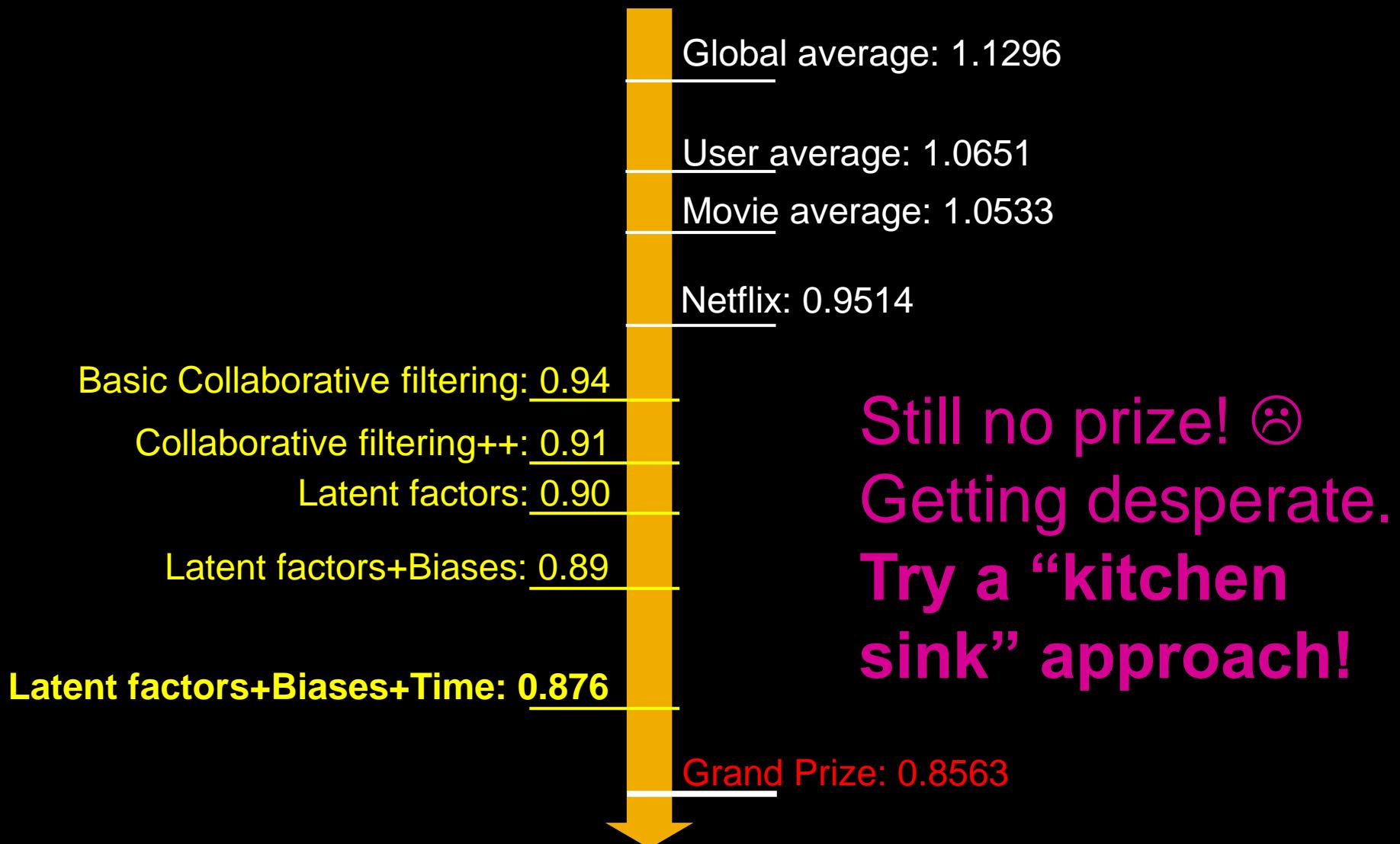
- Make parameters b_x and b_i to depend on time
 - (1) Parameterize time-dependence by linear trends
 - (2) Each bin corresponds to 10 consecutive weeks

$$b_i(t) = b_i + b_{i,\text{Bin}(t)}$$

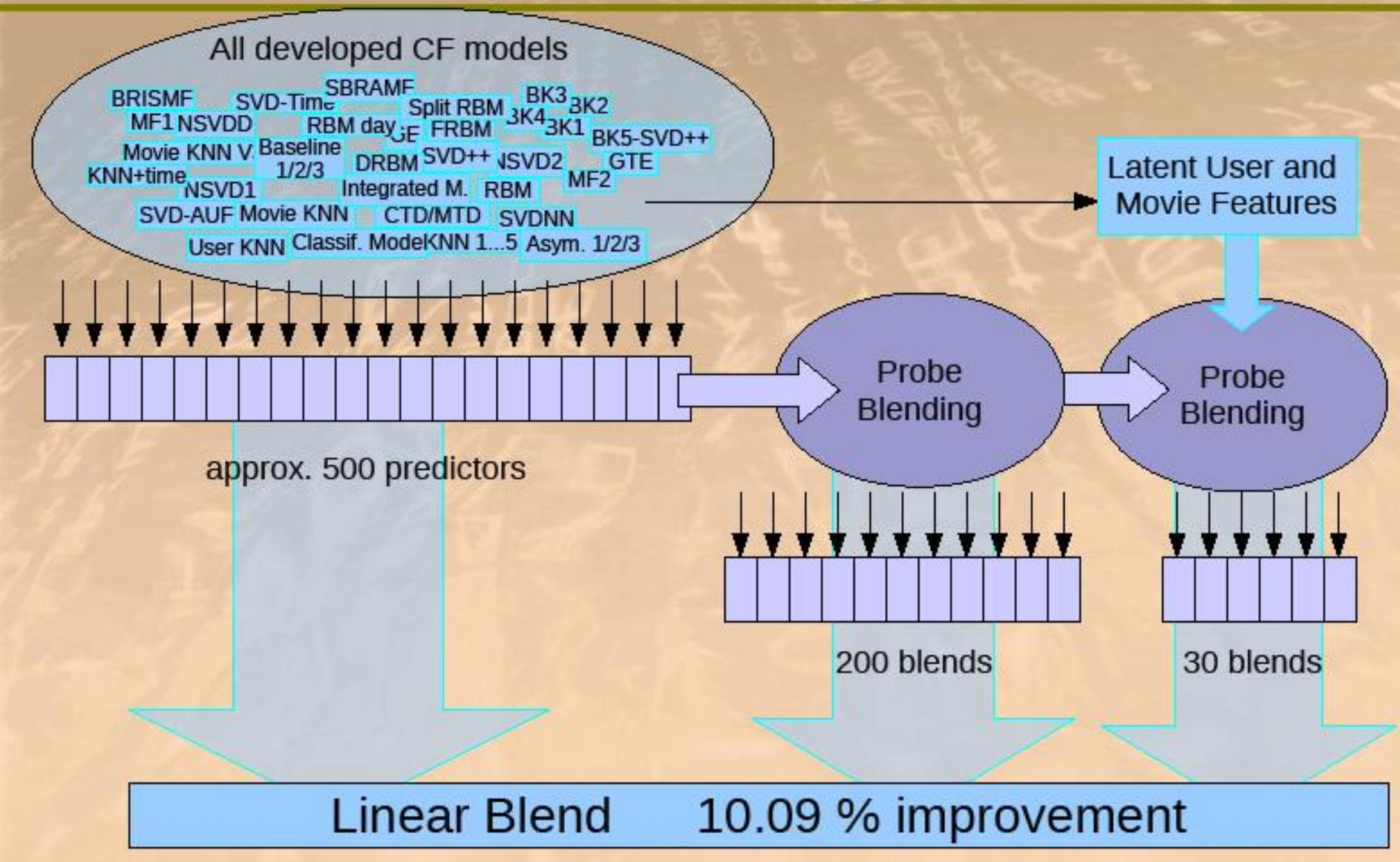
- Add temporal dependence to factors

- $p_x(t)$... user preference vector on day t

Performance of Various Methods



The big picture Solution of BellKor's Pragmatic Chaos



Standing on June 26th 2009

NETFLIX

Netflix Prize

Home Rules Leaderboard Register Update Submit Download

Leaderboard

Display top 20 leaders.

Rank	Team Name	Best Score	% Improvement	Last Submit Time
1	BellKor's Pragmatic Chaos	0.8558	10.05	2009-06-26 18:42:37
Grand Prize - RMSE <= 0.8563				
2	PragmaticTheory	0.8582	9.80	2009-06-25 22:15:51
3	BellKor in BigChaos	0.8590	9.71	2009-05-13 08:14:09
4	Grand Prize Team	0.8593	9.68	2009-06-12 08:20:24
5	Dace	0.8604	9.56	2009-04-22 05:57:03
6	BigChaos	0.8613	9.47	2009-06-23 23:06:52
Progress Prize 2008 - RMSE = 0.8616 - Winning Team: BellKor in BigChaos				
7	BellKor	0.8620	9.40	2009-06-24 07:16:02
8	Gravity	0.8634	9.25	2009-04-22 18:31:32
9	Opera Solutions	0.8638	9.21	2009-06-26 23:18:13
10	BruceDengDiaoCiYiYou	0.8638	9.21	2009-06-27 00:55:55
11	pengpengzhou	0.8638	9.21	2009-06-27 01:06:43
12	xvector	0.8639	9.20	2009-06-26 13:49:04
13	xiangliang	0.8639	9.20	2009-06-26 07:47:34

June 26th submission triggers 30-day “last call”

Netflix Prize

COMPLETED[Home](#) | [Rules](#) | [Leaderboard](#) | [Update](#) | [Download](#)

Leaderboard

Showing Test Score. [Click here to show quiz score](#)Display top leaders.

Rank	Team Name	Best Test Score	% Improvement	Best Submit Time
1	BellKor's Pragmatic Chaos	0.8567	10.06	2009-07-26 18:18:28
2	The Ensemble	0.8567	10.06	2009-07-26 18:38:22
3	Grand Prize Team	0.8562	9.90	2009-07-10 21:44:40
4	Opera Solutions and Vandelay United	0.8588	9.84	2009-07-10 01:12:31
5	Vandelay Industries!	0.8591	9.81	2009-07-10 00:32:20
6	PragmaticTheory	0.8594	9.77	2009-06-24 12:06:56
7	BellKor in BigChaos	0.8601	9.70	2009-05-13 08:14:09
8	Dace	0.8612	9.59	2009-07-24 17:18:43
9	Feeds2	0.8622	9.48	2009-07-12 13:11:51
10	BigChaos	0.8623	9.47	2009-04-07 12:33:59
11	Opera Solutions	0.8623	9.47	2009-07-24 00:34:07
12	BellKor	0.8624	9.46	2009-07-26 17:19:11

Grand Prize - RMSE = 0.8567 - Winning Team: BellKor's Pragmatic Chaos

1	BellKor's Pragmatic Chaos	0.8567	10.06	2009-07-26 18:18:28
2	The Ensemble	0.8567	10.06	2009-07-26 18:38:22
3	Grand Prize Team	0.8562	9.90	2009-07-10 21:44:40
4	Opera Solutions and Vandelay United	0.8588	9.84	2009-07-10 01:12:31
5	Vandelay Industries!	0.8591	9.81	2009-07-10 00:32:20
6	PragmaticTheory	0.8594	9.77	2009-06-24 12:06:56
7	BellKor in BigChaos	0.8601	9.70	2009-05-13 08:14:09
8	Dace	0.8612	9.59	2009-07-24 17:18:43
9	Feeds2	0.8622	9.48	2009-07-12 13:11:51
10	BigChaos	0.8623	9.47	2009-04-07 12:33:59
11	Opera Solutions	0.8623	9.47	2009-07-24 00:34:07
12	BellKor	0.8624	9.46	2009-07-26 17:19:11

Progress Prize 2008 - RMSE = 0.8627 - Winning Team: BellKor in BigChaos

13	xiangliang	0.8642	9.27	2009-07-15 14:53:22
14	Gravity	0.8643	9.26	2009-04-22 18:31:32
15	Ces	0.8651	9.18	2009-06-21 19:24:53
16	Invisible Ideas	0.8653	9.15	2009-07-15 15:53:04
17	Just a guy in a garage	0.8662	9.06	2009-05-24 10:02:54
18	J Dennis Su	0.8666	9.02	2009-03-07 17:16:17
19	Craig Carmichael	0.8666	9.02	2009-07-25 16:00:54
20	acmehill	0.8668	9.00	2009-03-21 16:20:50

Progress Prize 2007 - RMSE = 0.8723 - Winning Team: KorBell

Million \$ Awarded Sept 21st 2009



Acknowledgments

- Some slides and plots borrowed from
Yehuda Koren, Robert Bell and Padhraic Smyth, Jure Leskovec
- **Further reading:**
 - Y. Koren, Collaborative filtering with temporal dynamics, KDD '09
 - [Matrix Factorization Techniques for Recommender Systems](#)
 - [How the Netflix Prize was won](#)

Today

- Text Classification: Definition and Overview
- Vector Space Classification
 - Rocchio
 - kNN
 - Naïve Bayes



Earthquuuuuuuuakessss!!!



VIDEO

POLITICS

SPORTS

SCIENCE/TECH

LOCAL

ENTERTAINMENT

Grandmother Classifies 79% Of Everything A Shame

NEWS • Family • Local • ISSUE 47•47 ISSUE 45•29 • Jul 18, 2009



3.0K



382



20

SANDUSKY, OH—According to those close to Gertrude Wharton, the grandmother of nine declares 79 percent of everything she witnesses, experiences, or hears about from friends to be "a shame."



Wharton, 83, says it's a shame her husband isn't alive to see what a shame their grandchildren have become.

"No matter what happens, her response is always, 'That's a shame,'" said Wharton's son Kevin, 46. "From the recent passing of her friend Lillian to the fact that her coupon for chicken bouillon cubes expired last week, I can't have a conversation with her without being told something is a shame. Is this really how she sees the world now?"

Though Wharton, 83, has proclaimed things to be a shame in the past, her current usage of the word has been a growing cause for concern. Witnesses report that in the past 24 hours alone she has

Standing queries

- The path from IR to text classification:
 - You have an information need to monitor, say:
 - [presidential polls](#)
 - 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

<http://www.google.com/alerts>

A text classification task: Email spam filtering

From: "" <takworlld@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>

=====

How would you write a program that would automatically detect and delete this type of message?

Formal definition of Text Classification: Training

Given:

- A **document space** X
 - Documents are represented in this space – typically some type of high-dimensional space.
- A fixed set of **classes** $C = \{c_1, c_2, \dots, c_J\}$
 - The classes are human-defined for the needs of an application (e.g., relevant vs. nonrelevant).
- A **training set** D of labeled documents with each labeled document $\langle d, c \rangle \in X \times C$

Using a learning method or **learning algorithm**, we then wish to learn a **classifier** γ that maps documents to classes:

$$\gamma : X \rightarrow C$$

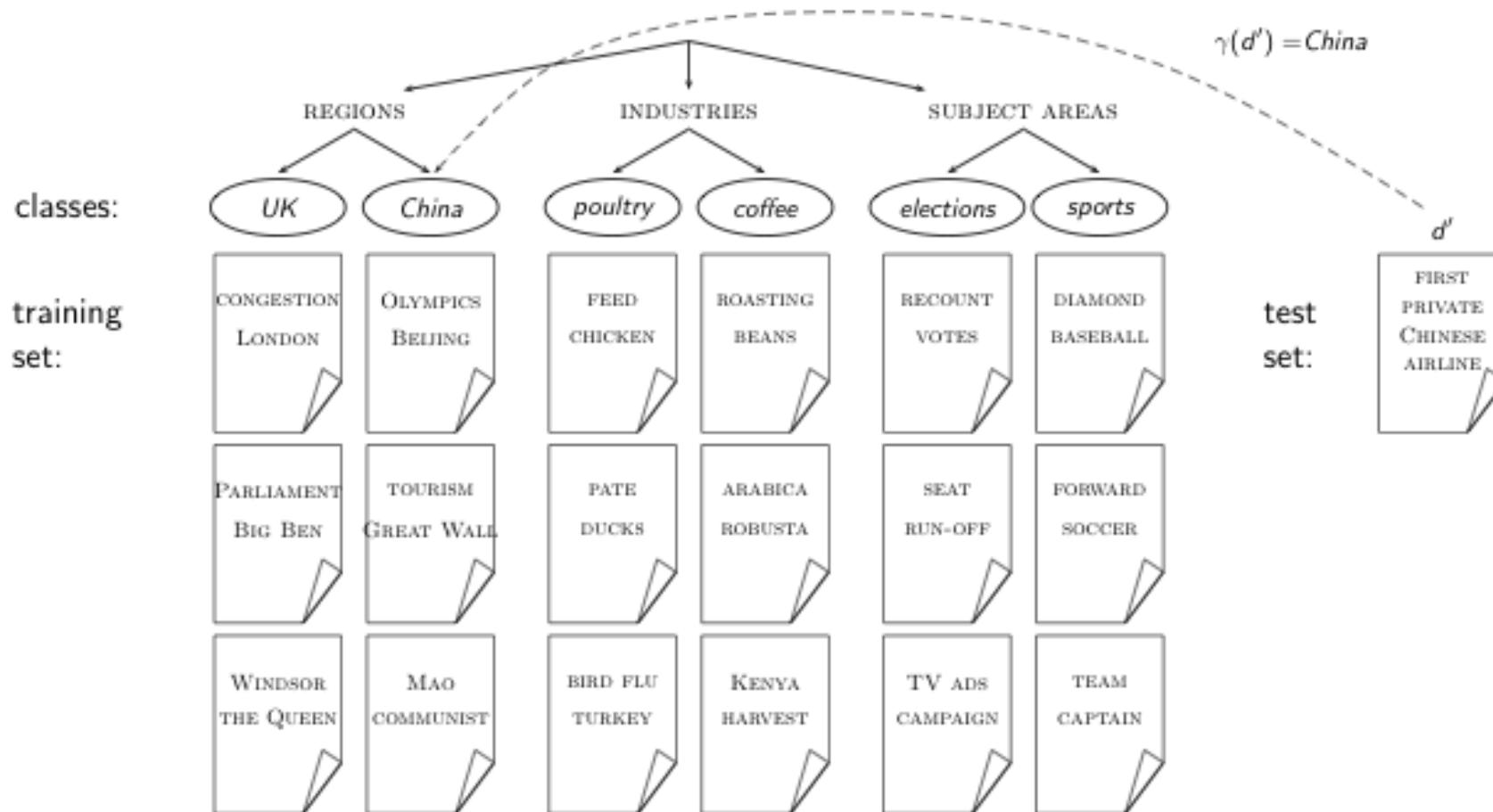
Formal definition of Text Classification : Application/Testing

Given: a description $d \in X$ of a document

Determine: $\gamma(d) \in C$,

that is, the class that is most appropriate for d

Topic classification



Exercise

- Find examples of uses of text classification in information retrieval

Examples of how search engines use classification

- Language identification (classes: English vs. French etc.)
- The automatic detection of spam pages (spam vs. nonspam)
- The automatic detection of sexually explicit content (sexually explicit vs. not)
- Topic-specific or *vertical* search – restrict search to a “vertical” like “related to health” (relevant to vertical vs. not)
- Standing queries (e.g., Google Alerts)
- Sentiment detection: is a movie or product review positive or negative (positive vs. negative)

How can we classify?
Any classification method?

Classification in Machine Learning

- Text classification as a learning problem
- (i) Supervised learning of a the classification function Υ and
(ii) its application to classifying new documents
- We will look at a couple of methods for doing this:Rocchio,
kNN, Naive Bayes
- No free lunch: requires hand-classified training data
- But this manual classification can be done by non-experts.

Vector Space Classification

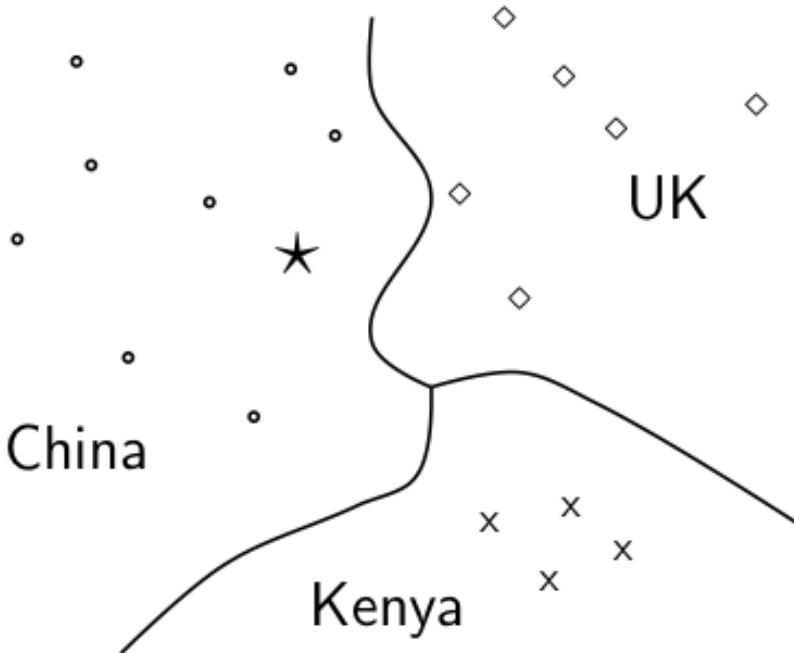
Recall: Vector Space Representation

- Each document is a vector, one component for each term.
- Terms are axes.
- High-dimensional vector space: 100,000s dimensions
- Normalize vectors (documents) to unit length.
- How can we do classification in this space?

Vector Space Classification

- As before, the training set is a set of documents, each labeled with its class (e.g., topic)
 - In vector space classification, this set corresponds to a labeled set of points (or, equivalently, vectors) in the vector space
-
- Premise 1: Documents in the same class form a contiguous region
 - Premise 2: Documents from different classes don't overlap.
-
- We define lines, surfaces, hypersurfaces to divide regions.

Classes in the vector space



Should the document \star be assigned to China, UK or Kenya?

Find separators between the classes

Based on these separators:

\star should be assigned to China

How do we find separators that do a good job at classifying new documents like

\star ? – Main topic of today

Rocchio

Rocchio classification: Basic idea

- Compute a centroid for each class
 - The centroid is the average of all documents in the class.
- Assign each test document to the class of its closest centroid.

Recall 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 $\vec{v}(d)$ is the vector space representation of d .

Rocchio algorithm

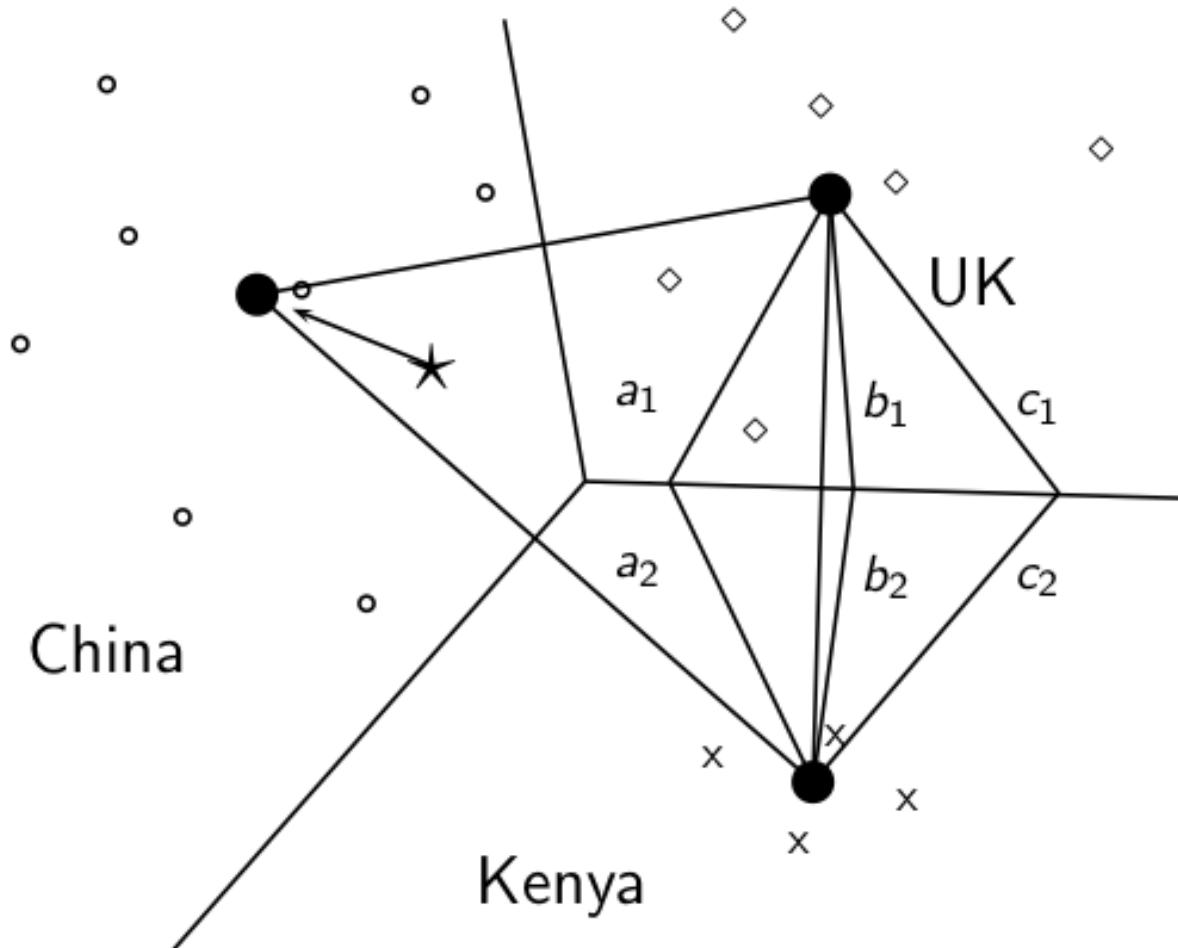
TRAINROCCHIO(\mathbb{C}, \mathbb{D})

- 1 **for each** $c_j \in \mathbb{C}$
- 2 **do** $D_j \leftarrow \{d : \langle d, c_j \rangle \in \mathbb{D}\}$
- 3 $\vec{\mu}_j \leftarrow \frac{1}{|D_j|} \sum_{d \in D_j} \vec{v}(d)$
- 4 **return** $\{\vec{\mu}_1, \dots, \vec{\mu}_J\}$

APPLYROCCHIO($\{\vec{\mu}_1, \dots, \vec{\mu}_J\}, d$)

- 1 **return** $\arg \min_j |\vec{\mu}_j - \vec{v}(d)|$

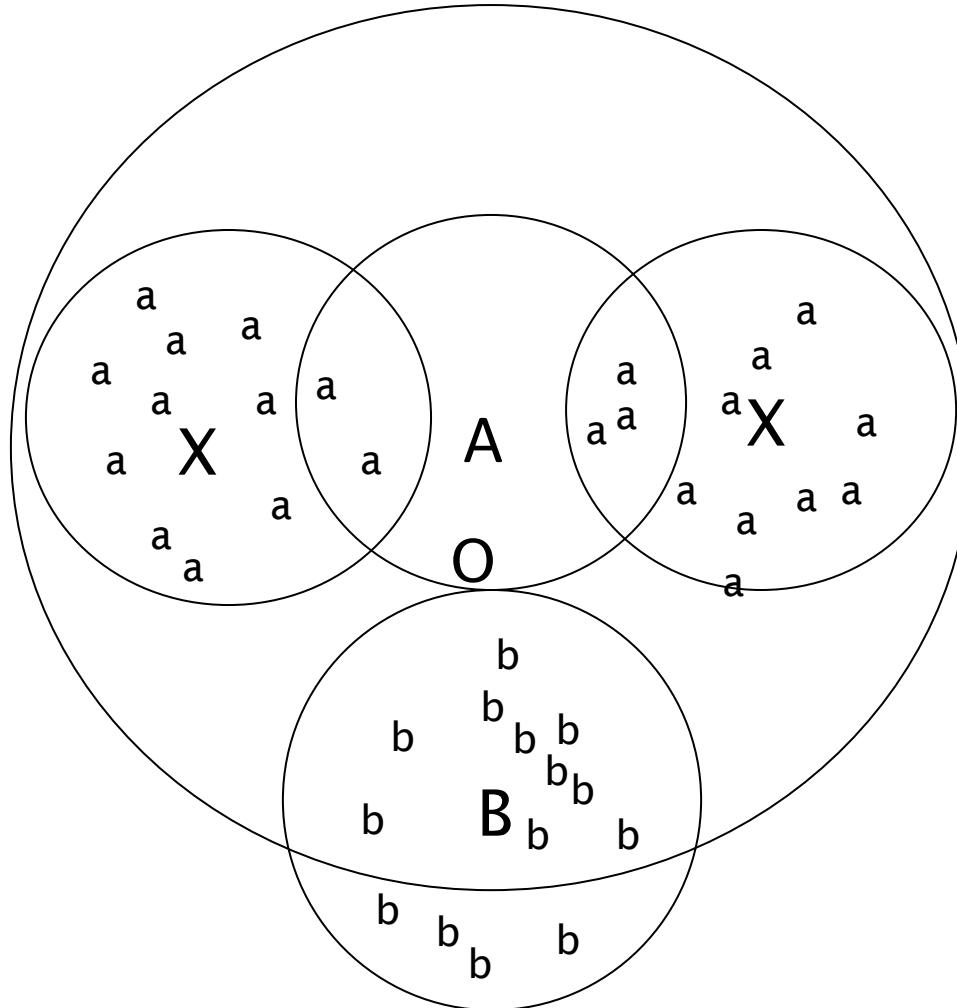
Rocchio illustrated : $a_1 = a_2$, $b_1 = b_2$, $c_1 = c_2$



Rocchio properties

- Rocchio forms a simple representation for each class: the **centroid**
 - We can interpret the centroid as the **prototype** of the class.
- Classification is based on similarity to / distance from centroid/prototype.
- However, does not guarantee that classifications are consistent with the training data!

Rocchio cannot handle nonconvex, multimodal classes



Exercise: Why is Rocchio not expected to do well for the classification task a vs. b here?

- A is centroid of the a's, B is centroid of the b's.
- The point o is closer to A than to B.
- But o is a better fit for the b class.
- A is a multimodal class with two prototypes.
- But in Rocchio we only have one prototype.

kNN (k Nearest Neighbors)

kNN classification

- kNN classification is another vector space classification method.
- It also is very simple and easy to implement.
- kNN is more accurate (in most cases) than Naive Bayes and Rocchio.
- If you need to get a pretty accurate classifier up and running in a short time . . .
- . . . and you don't care about efficiency that much . . .
- . . . use kNN.

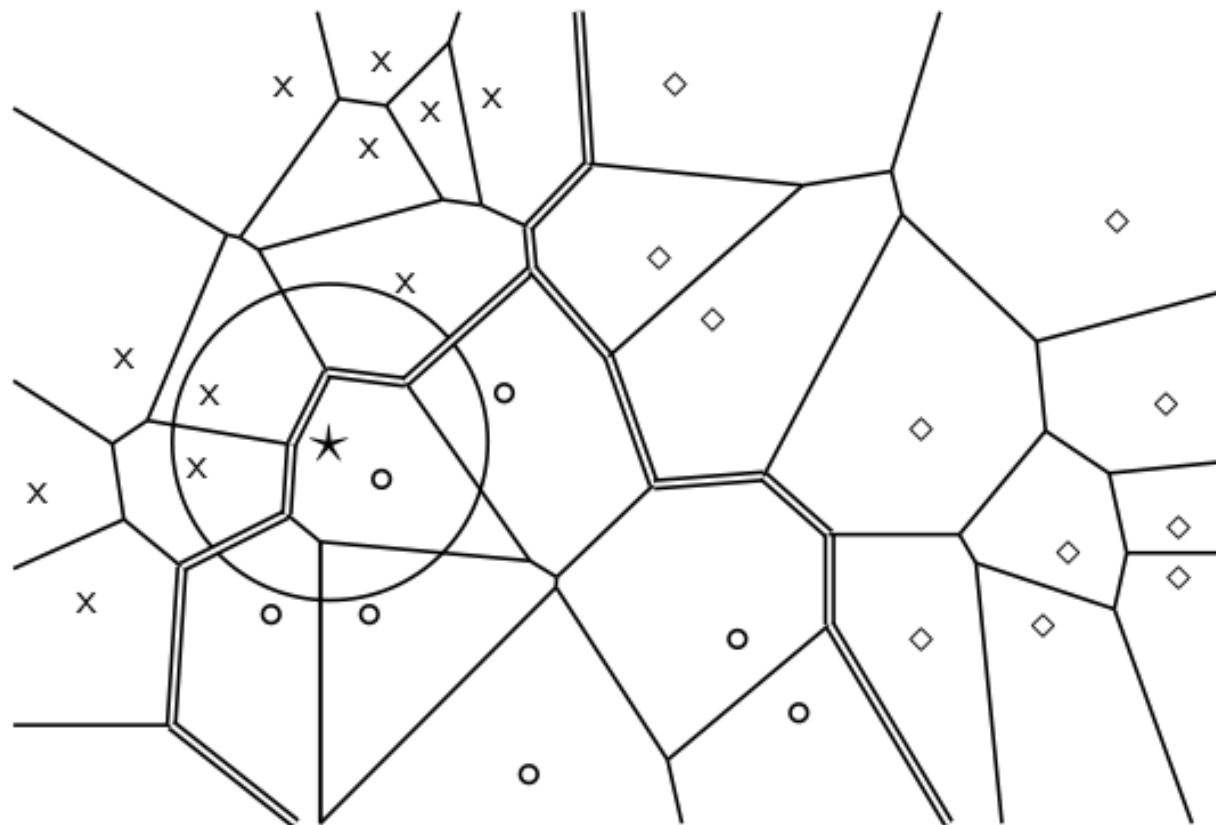
kNN classification

- kNN = k nearest neighbors
- **kNN classification rule for $k = 1$ (1NN):** Assign each test document to the class of its nearest neighbor in the training set.
- 1NN is not very robust – one document can be mislabeled or atypical.
- **kNN classification rule for $k > 1$ (kNN):** Assign each test document to the **majority class of its k nearest neighbors** in the training set.
- Rationale of kNN: contiguity hypothesis
 - We expect a test document d to have the same label as the training documents located in the local region surrounding d .

Probabilistic kNN

- Probabilistic version of kNN: $P(c|d)$ = fraction of k neighbors of d that are in c
- **kNN classification rule for probabilistic kNN:** Assign d to class c with highest $P(c|d)$

Probabilistic kNN



1NN, 3NN
classification
decision
for star?

kNN algorithm

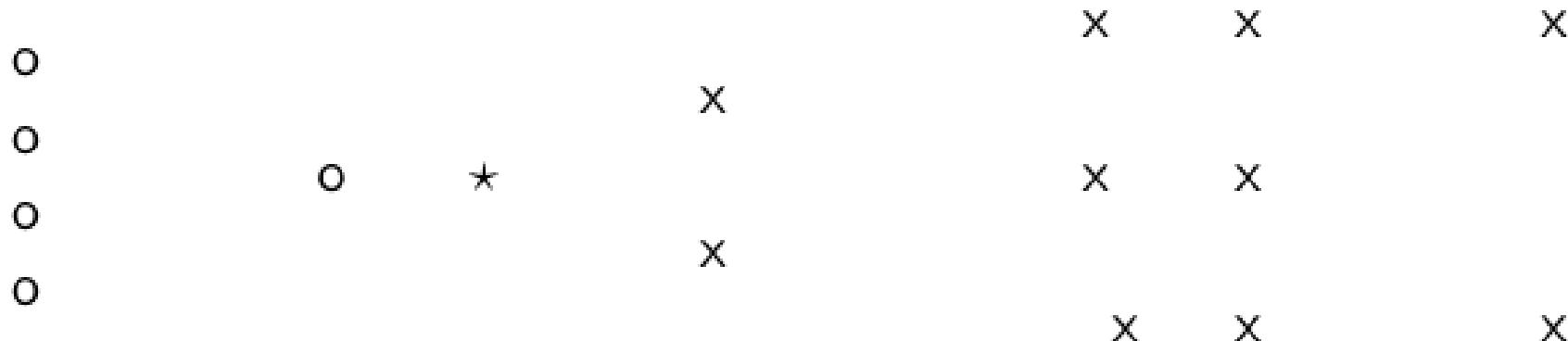
TRAIN-KNN(\mathbb{C}, \mathbb{D})

- 1 $\mathbb{D}' \leftarrow \text{PREPROCESS}(\mathbb{D})$
- 2 $k \leftarrow \text{SELECT-K}(\mathbb{C}, \mathbb{D}')$
- 3 **return** \mathbb{D}', k

APPLY-KNN(\mathbb{D}', k, d)

- 1 $S_k \leftarrow \text{COMPUTENEARESTNEIGHBORS}(\mathbb{D}', k, d)$
- 2 **for each** $c_j \in \mathbb{C}(\mathbb{D}')$
- 3 **do** $p_j \leftarrow |S_k \cap c_j|/k$
- 4 **return** $\arg \max_j p_j$

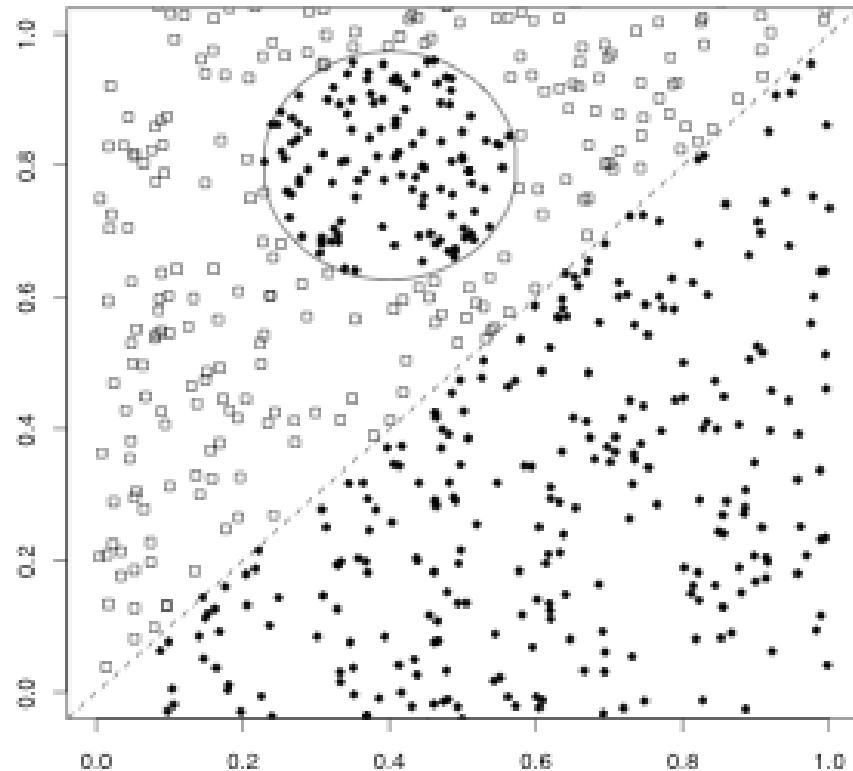
Exercise



How is star classified by:

- (i) 1-NN (ii) 3-NN (iii) 9-NN (iv) 15-NN (v) Rocchio?

A nonlinear problem (Rocchio vs kNN)



- Linear classifier like Rocchio does badly on this task.
- kNN will do well (assuming enough training data)

kNN: Discussion

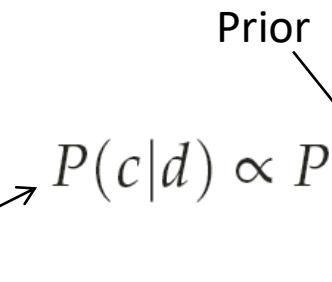
- No training necessary
 - But linear preprocessing of documents is as expensive as training Naive Bayes.
 - We always preprocess the training set, so in reality training time of kNN is linear.
- kNN is very accurate if training set is large.
- But kNN can be very inaccurate if training set is small.

Naive Bayes Classifier

The Naive Bayes Classifier

- The Naive Bayes classifier is a probabilistic classifier
- We compute the probability of a document d being in a class c as follows:

$$P(c|d) \propto P(c) \prod_{1 \leq k \leq n_d} P(t_k|c)$$

Posterior 

- $P(c)$ is the prior probability of c .
- n_d is the length of the document. (number of tokens)
- $P(t_k | c)$ is the conditional probability of term t_k occurring in a document of class c
- $P(t_k | c)$ as a measure of **how much evidence** t_k contributes that c is the correct class.
- If a document's terms do not provide clear evidence for one class vs. another, we choose the c with highest $P(c)$ probability.

Maximum a posteriori class

- Our goal in Naive Bayes classification is to find the “best” class.
- The best class is the most likely or **maximum a posteriori (MAP) class**

c_{map} :

$$c_{\text{map}} = \arg \max_{c \in \mathbb{C}} \hat{P}(c|d) = \arg \max_{c \in \mathbb{C}} \hat{P}(c) \prod_{1 \leq k \leq n_d} \hat{P}(t_k|c)$$

Taking the log

- Multiplying lots of small probabilities can result in [floating point underflow](#).
- Since $\log(xy) = \log(x) + \log(y)$, we can sum log probabilities instead of multiplying probabilities.
- Since log is a monotonic function, the class with the highest score does not change.
- So what we usually compute in practice is:

$$c_{\text{map}} = \arg \max_{c \in \mathbb{C}} [\log \hat{P}(c) + \sum_{1 \leq k \leq n_d} \log \hat{P}(t_k | c)]$$

Naive Bayes classifier

- Classification rule:

$$c_{\text{map}} = \arg \max_{c \in \mathbb{C}} [\log \hat{P}(c) + \sum_{1 \leq k \leq n_d} \log \hat{P}(t_k|c)]$$

- Simple interpretation:

- Each conditional parameter $\log \hat{P}(t_k|c)$ is a weight that indicates how good an indicator t_k is for c .
- The prior $\log \hat{P}(c)$ is a weight that indicates the relative frequency of c .
- The sum of log prior and term weights is then a measure of how much evidence there is for the document being in the class.
- We select the class with the most evidence.

Parameter estimation take 1: Maximum likelihood

- Estimate parameters $\hat{P}(c)$ and $\hat{P}(t_k|c)$ from train data: How?
- Prior:

$$\hat{P}(c) = \frac{N_c}{N}$$

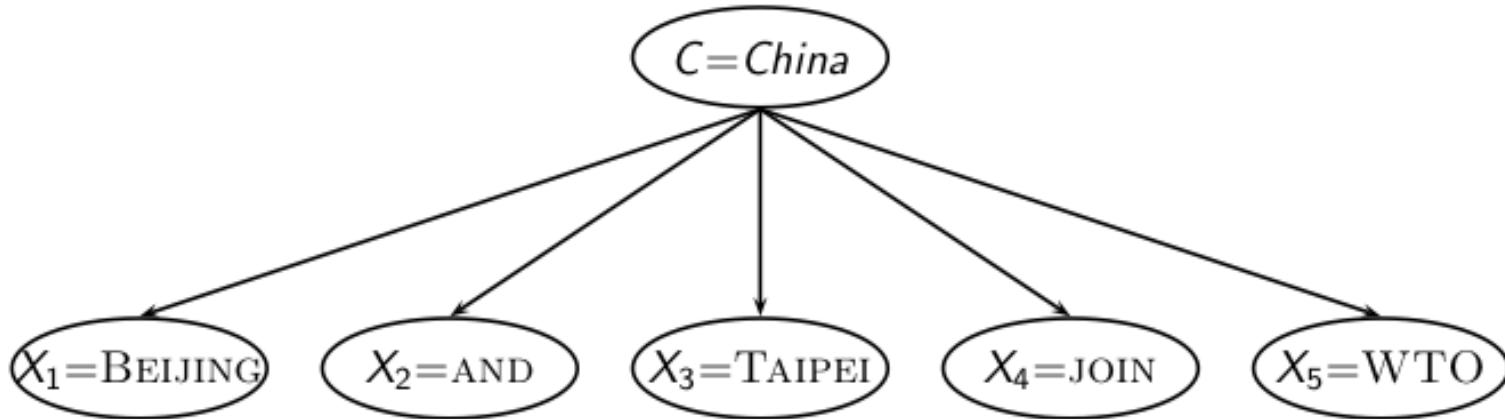
- N_c : number of docs in class c ; N : total number of docs
- Conditional probabilities:

$$\hat{P}(t|c) = \frac{T_{ct}}{\sum_{t' \in V} T_{ct'}}$$

- T_{ct} is the number of tokens of t in training documents from class c (includes multiple occurrences)
- We've made a **Naive Bayes independence assumption** here:

$$\hat{P}(t_{k_1}|c) = \hat{P}(t_{k_2}|c)$$

The problem with maximum likelihood estimates: Zeros



$$\begin{aligned} P(China | d) \propto & P(China) \cdot P(BEIJING|China) \cdot P(AND|China) \\ & \cdot P(TAIPEI|China) \cdot P(JOIN|China) \cdot \\ & P(WTO|China) \end{aligned}$$

$$\hat{P}(WTO|China) = \frac{T_{China,WTO}}{\sum_{t' \in V} T_{China,t'}} = \frac{0}{\sum_{t' \in V} T_{China,t'}} = 0$$

The problem with maximum likelihood estimates: Zeros (cont)

- If there were no occurrences of WTO in documents in class China, we'd get a zero estimate:

$$\hat{P}(\text{WTO}|\text{China}) = \frac{T_{\text{China}, \text{WTO}}}{\sum_{t' \in V} T_{\text{China}, t'}} = 0$$

- → We will get $P(\text{China} | d) = 0$ for any document that contains WTO!
- Zero probabilities cannot be conditioned away.

To avoid zeros: Add-one smoothing

- Before:

$$\hat{P}(t|c) = \frac{T_{ct}}{\sum_{t' \in V} T_{ct'}}$$

- Now: Add one to each count to avoid zeros:

$$\hat{P}(t|c) = \frac{T_{ct} + 1}{\sum_{t' \in V} (T_{ct'} + 1)} = \frac{T_{ct} + 1}{(\sum_{t' \in V} T_{ct'}) + B}$$

- B is the number of different words (in this case the size of the vocabulary: $|V| = M$)

To avoid zeros: Add-one smoothing

- Estimate parameters from the training corpus using add-one smoothing
- For a new document, for each class, compute sum of (i) log of prior and (ii) logs of conditional probabilities of the terms
- Assign the document to the class with the largest score

Naive Bayes: Training

TRAINMULTINOMIALNB(\mathbb{C}, \mathbb{D})

- 1 $V \leftarrow \text{EXTRACTVOCABULARY}(\mathbb{D})$
- 2 $N \leftarrow \text{COUNTDOCS}(\mathbb{D})$
- 3 **for each** $c \in \mathbb{C}$
- 4 **do** $N_c \leftarrow \text{COUNTDOCSINCLASS}(\mathbb{D}, c)$
- 5 $prior[c] \leftarrow N_c/N$
- 6 $text_c \leftarrow \text{CONCATENATETEXTOفالDOCSINCLASS}(\mathbb{D}, c)$
- 7 **for each** $t \in V$
- 8 **do** $T_{ct} \leftarrow \text{COUNTTOKENSOFTERM}(text_c, t)$
- 9 **for each** $t \in V$
- 10 **do** $condprob[t][c] \leftarrow \frac{T_{ct}+1}{\sum_{t'}(T_{ct'}+1)}$
- 11 **return** $V, prior, condprob$

Naive Bayes: Testing

```
APPLYMULTINOMIALNB( $\mathbb{C}, V, prior, condprob, d$ )
1    $W \leftarrow \text{EXTRACTTOKENSFROMDOC}(V, d)$ 
2   for each  $c \in \mathbb{C}$ 
3     do  $score[c] \leftarrow \log prior[c]$ 
4     for each  $t \in W$ 
5       do  $score[c] += \log condprob[t][c]$ 
6   return  $\arg \max_{c \in \mathbb{C}} score[c]$ 
```

Example

	docID	words in document	in $c = China$?
training set	1	Chinese Beijing Chinese	yes
	2	Chinese Chinese Shanghai	yes
	3	Chinese Macao	yes
	4	Tokyo Japan Chinese	no
test set	5	Japan Chinese Chinese Chinese Tokyo	?

- What do we need?
 - Class priors: $P(c)$, $P(\text{not } c)$
 - Conditional probabilities: $P(t|c)$, $P(t|\text{not } c)$

Example

	docID	words in document	in $c = China$?
training set	1	Chinese Beijing Chinese	yes
	2	Chinese Chinese Shanghai	yes
	3	Chinese Macao	yes
	4	Tokyo Japan Chinese	no
test set	5	Japan Chinese Chinese Tokyo	?

$$c_{\text{map}} = \arg \max_{c \in \mathbb{C}} [\log \hat{P}(c) + \sum_{1 \leq k \leq n_d} \log \hat{P}(t_k|c)]$$

$$\hat{P}(c) = \boxed{} \text{ and } \hat{P}(\bar{c}) = \boxed{}$$

$$\hat{P}(\text{Chinese}|c) = \boxed{}$$

$$\hat{P}(\text{Tokyo}|c) = \hat{P}(\text{Japan}|c) = \boxed{}$$

$$\hat{P}(\text{Chinese}|\bar{c}) = \boxed{}$$

$$\hat{P}(\text{Tokyo}|\bar{c}) = \hat{P}(\text{Japan}|\bar{c}) = \boxed{}$$

Example

	docID	words in document	in $c = China$?
training set	1	Chinese Beijing Chinese	yes
	2	Chinese Chinese Shanghai	yes
	3	Chinese Macao	yes
	4	Tokyo Japan Chinese	no
test set	5	Japan Chinese Chinese Chinese Tokyo	?

$$c_{\text{map}} = \arg \max_{c \in \mathbb{C}} [\log \hat{P}(c) + \sum_{1 \leq k \leq n_d} \log \hat{P}(t_k | c)]$$

$$\hat{P}(c) = 3/4 \text{ and } \hat{P}(\bar{c}) = 1/4$$

$$\hat{P}(\text{Chinese}|c) = (5+1)/(8+6) = 6/14 = 3/7$$

$$\hat{P}(\text{Tokyo}|c) = \hat{P}(\text{Japan}|c) = (0+1)/(8+6) = 1/14$$

$$\hat{P}(\text{Chinese}|\bar{c}) = (1+1)/(3+6) = 2/9$$

$$\hat{P}(\text{Tokyo}|\bar{c}) = \hat{P}(\text{Japan}|\bar{c}) = (1+1)/(3+6) = 2/9$$

Example

	docID	words in document	in $c = China$?
training set	1	Chinese Beijing Chinese	yes
	2	Chinese Chinese Shanghai	yes
	3	Chinese Macao	yes
	4	Tokyo Japan Chinese	no
test set	5	Japan Chinese Chinese Chinese Tokyo	?

$$c_{\text{map}} = \arg \max_{c \in \mathbb{C}} [\log \hat{P}(c) + \sum_{1 \leq k \leq n_d} \log \hat{P}(t_k | c)]$$

$$\hat{P}(c|d) \propto \log 3/4 + \log 1/14 + 3 \log 3/7 + \log 1/14 \approx -3.52$$

$$\hat{P}(\bar{c}|d) \propto \log 1/4 + \log 2/9 + 3 \log 2/9 + \log 2/9 \approx -3.86$$

Thus, the classifier assigns the test document to $c = China$.

Example for Rocchio Classification

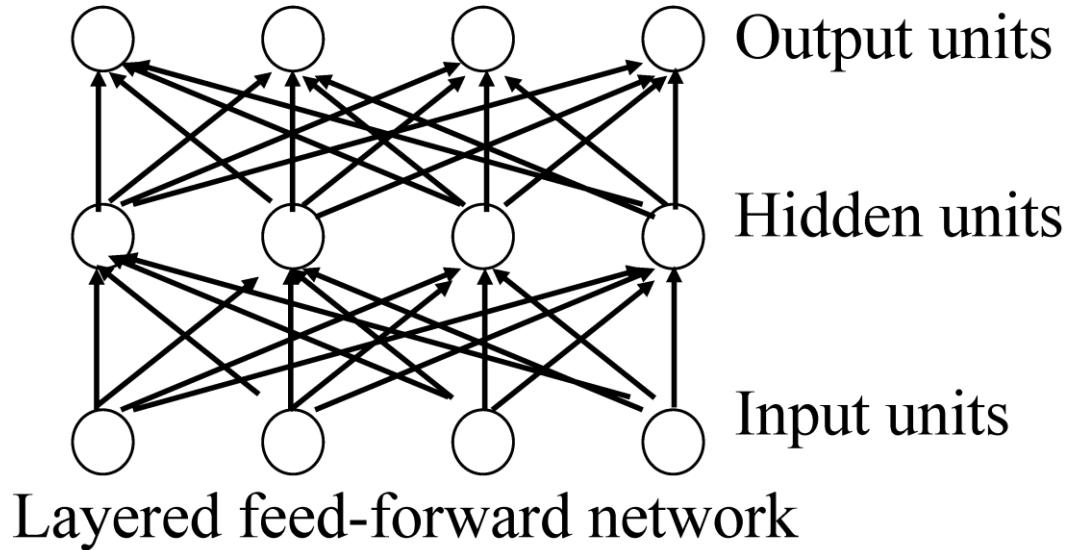
	docID	words in document	in $c = China$?
training set	1	Chinese Beijing Chinese	yes
	2	Chinese Chinese Shanghai	yes
	3	Chinese Macao	yes
	4	Tokyo Japan Chinese	no
test set	5	Japan Chinese Chinese Chinese Tokyo	?

Other Classification Methods

- Neural Network
- Deep Learning (e.g., recurrent neural network, transformer)

Neural Network

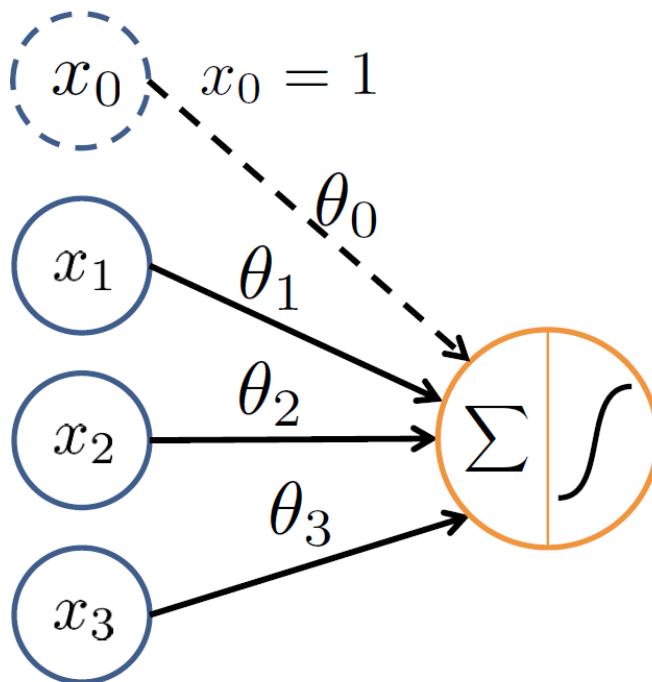
Neural Networks



- Neural networks are made up of **nodes** or **units**, connected by **links**
- Each link has an associated **weight** and **activation level**
- Each node has an **input function** (typically summing over weighted inputs), an **activation function**, and an **output**

Neuron Model: Logistic Unit

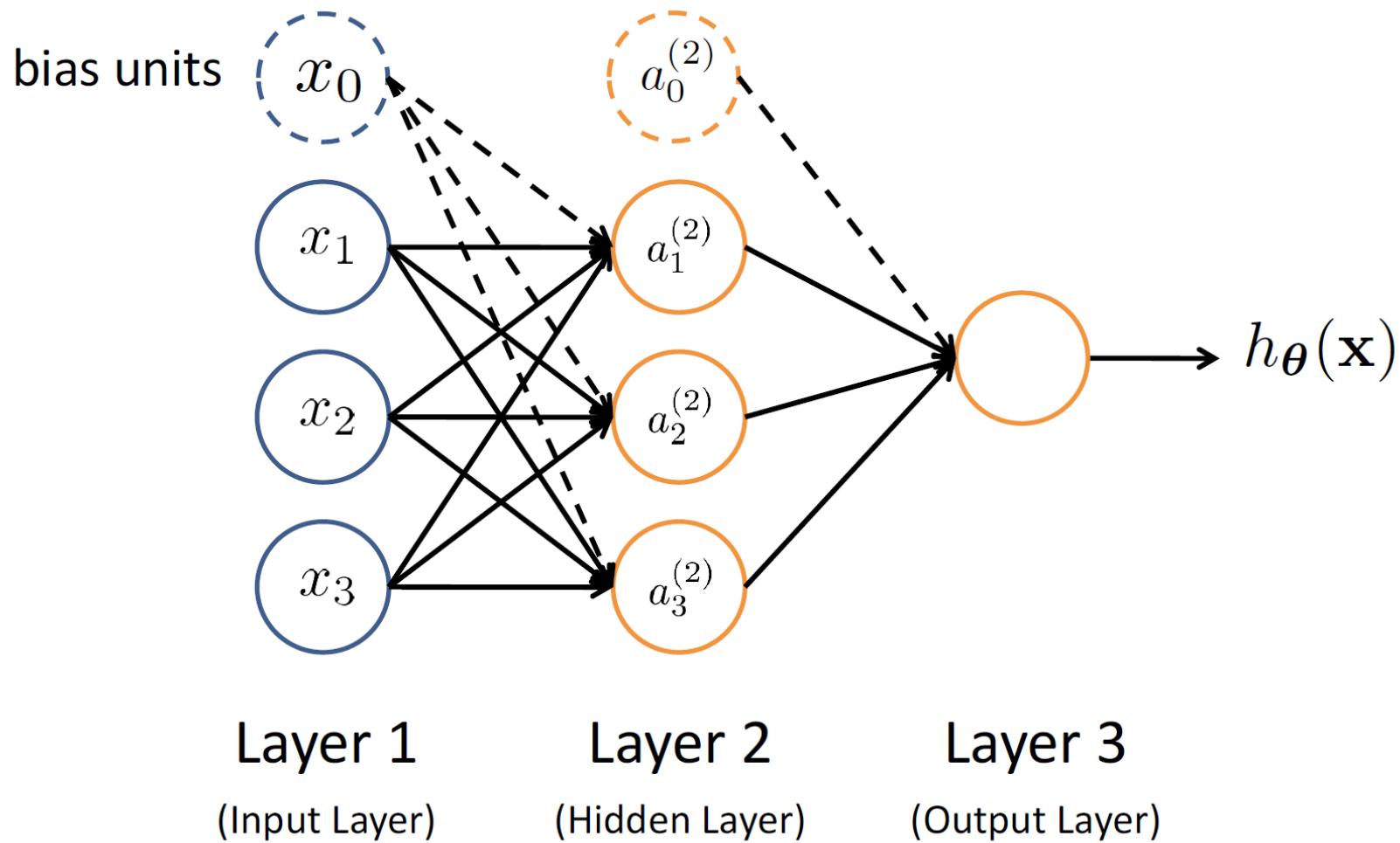
“bias unit”



$$\mathbf{x} = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \end{bmatrix} \quad \boldsymbol{\theta} = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \theta_3 \end{bmatrix}$$

$$h_{\boldsymbol{\theta}}(\mathbf{x}) = g(\boldsymbol{\theta}^T \mathbf{x}) \\ = \frac{1}{1 + e^{-\boldsymbol{\theta}^T \mathbf{x}}}$$

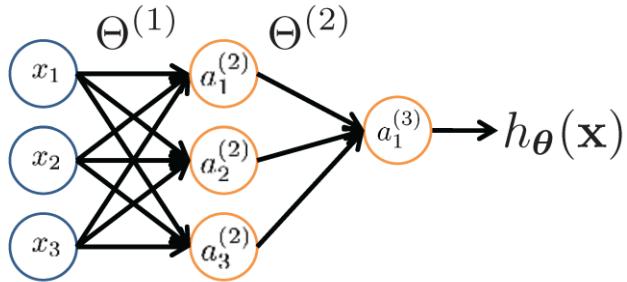
Neural Network



Feed-Forward Process

- Input layer units are set by some exterior function (think of these as **sensors**), which causes their output links to be **activated** at the specified level
- Working forward through the network, the **input function** of each unit is applied to compute the input value
 - Usually this is just the weighted sum of the activation on the links feeding into this node
- The **activation function** transforms this input function into a final value
 - Typically this is a **nonlinear** function, often a **sigmoid** function corresponding to the “threshold” of that node

Neural Network



$a_i^{(j)}$ = “activation” of unit i in layer j
 $\Theta^{(j)}$ = weight matrix controlling function
mapping from layer j to layer $j + 1$

$$a_1^{(2)} = g(\Theta_{10}^{(1)}x_0 + \Theta_{11}^{(1)}x_1 + \Theta_{12}^{(1)}x_2 + \Theta_{13}^{(1)}x_3)$$

$$a_2^{(2)} = g(\Theta_{20}^{(1)}x_0 + \Theta_{21}^{(1)}x_1 + \Theta_{22}^{(1)}x_2 + \Theta_{23}^{(1)}x_3)$$

$$a_3^{(2)} = g(\Theta_{30}^{(1)}x_0 + \Theta_{31}^{(1)}x_1 + \Theta_{32}^{(1)}x_2 + \Theta_{33}^{(1)}x_3)$$

$$h_\Theta(x) = a_1^{(3)} = g(\Theta_{10}^{(2)}a_0^{(2)} + \Theta_{11}^{(2)}a_1^{(2)} + \Theta_{12}^{(2)}a_2^{(2)} + \Theta_{13}^{(2)}a_3^{(2)})$$

If network has s_j units in layer j and s_{j+1} units in layer $j+1$,
then $\Theta^{(j)}$ has dimension $s_{j+1} \times (s_j + 1)$

$$\Theta^{(1)} \in \mathbb{R}^{3 \times 4} \quad \Theta^{(2)} \in \mathbb{R}^{1 \times 4}$$

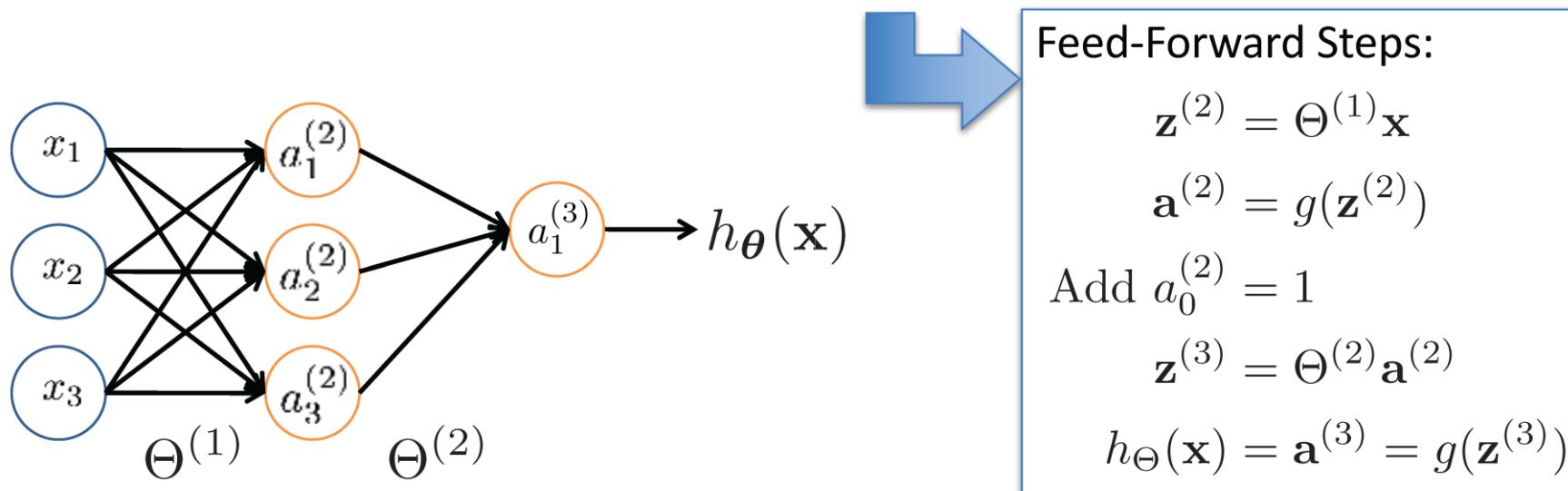
Vectorization

$$a_1^{(2)} = g \left(\Theta_{10}^{(1)} x_0 + \Theta_{11}^{(1)} x_1 + \Theta_{12}^{(1)} x_2 + \Theta_{13}^{(1)} x_3 \right) = g \left(z_1^{(2)} \right)$$

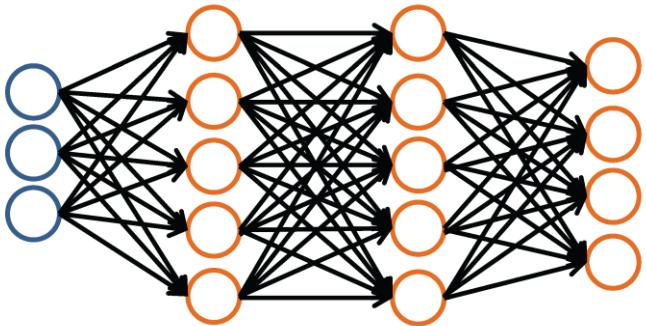
$$a_2^{(2)} = g \left(\Theta_{20}^{(1)} x_0 + \Theta_{21}^{(1)} x_1 + \Theta_{22}^{(1)} x_2 + \Theta_{23}^{(1)} x_3 \right) = g \left(z_2^{(2)} \right)$$

$$a_3^{(2)} = g \left(\Theta_{30}^{(1)} x_0 + \Theta_{31}^{(1)} x_1 + \Theta_{32}^{(1)} x_2 + \Theta_{33}^{(1)} x_3 \right) = g \left(z_3^{(2)} \right)$$

$$h_{\Theta}(\mathbf{x}) = g \left(\Theta_{10}^{(2)} a_0^{(2)} + \Theta_{11}^{(2)} a_1^{(2)} + \Theta_{12}^{(2)} a_2^{(2)} + \Theta_{13}^{(2)} a_3^{(2)} \right) = g \left(z_1^{(3)} \right)$$



Neural Network Classification



Given:

$$\{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_n, y_n)\}$$

$s \in \mathbb{N}^{+L}$ contains # nodes at each layer

- $s_0 = d$ (# features)

Binary classification

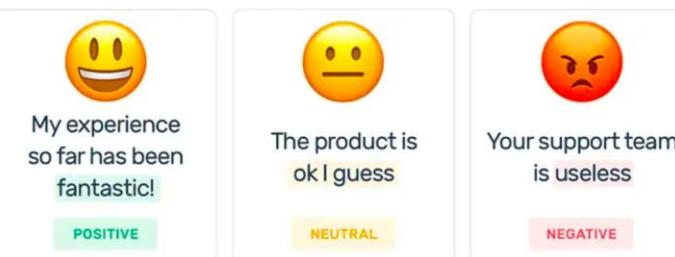
$y = 0$ or 1

1 output unit ($s_{L-1} = 1$)

Multi-class classification (K classes)

$\mathbf{y} \in \mathbb{R}^K$ e.g. $\begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}$ e.g., sentiment analysis
positive neural negative

K output units ($s_{L-1} = K$)



Neural Network Learning

Learning in NN: Backpropagation

- We cycle through our examples
 - If the output of the network is correct, no changes are made
 - If there is an error, weights are adjusted to reduce the error
- The trick is to assess the blame for the error and divide it among the contributing weights

General NN Cost Function

$$J(\Theta) = \frac{1}{n} \sum_{i=1}^n \mathcal{L}(h_\Theta(\mathbf{x}_i), y_i) + \frac{\lambda}{2n} \sum_{l=1}^{L-1} \sum_{i=1}^{s_{l-1}} \sum_{j=1}^{s_l} \left(\Theta_{ji}^{(l)} \right)^2$$

Given $h_\Theta(\mathbf{x}_i) = \hat{y}$

Loss of **Binary classification (1 or 0)** ↪ log loss

$$\mathcal{L}(\hat{y}, y) = -(y \log \hat{y} + (1 - y) \log(1 - \hat{y}))$$

Loss of **Regression** ↪ squared error loss

$$\mathcal{L}(\hat{y}, y) = \frac{1}{2} (\hat{y} - y)^2$$

Loss of **Multi-class classification with softmax regression** ↪ cross-entropy loss

$$\mathcal{L}(\hat{y}, y) = - \sum_{j=1}^k \mathbf{1}\{y = j\} \log \hat{y}_j = - \sum_{j=1}^k y_j \log \hat{y}_j$$

Optimizing the Neural Network

$$J(\Theta) = \frac{1}{n} \sum_{i=1}^n \mathcal{L}(h_\Theta(\mathbf{x}_i), \mathbf{y}_i) + \frac{\lambda}{2n} \sum_{l=1}^{L-1} \sum_{i=1}^{s_{l-1}} \sum_{j=1}^{s_l} \left(\Theta_{ji}^{(l)} \right)^2$$

Solve via: $\min_{\Theta} J(\Theta)$

$J(\Theta)$ is not convex, so GD on a neural net yields a local optimum

- But, tends to work well in practice

Need code to compute:

- $J(\Theta)$
- $\frac{\partial}{\partial \Theta_{ij}^{(l)}} J(\Theta)$

Training a Neural Network via Gradient Descent with Backprop

Backpropagation

Given: training set $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$

Initialize all $\Theta^{(l)}$ randomly (NOT to 0!)

Loop // each iteration is called an epoch

Set $\Delta_{ij}^{(l)} = 0 \quad \forall l, i, j$ (Used to accumulate gradient)

For each training instance (\mathbf{x}_i, y_i) :

Set $\mathbf{a}^{(1)} = \mathbf{x}_i$

Compute $\{\mathbf{a}^{(2)}, \dots, \mathbf{a}^{(L)}\}$ via forward propagation

Compute $\boldsymbol{\delta}^{(L)} = \mathbf{a}^{(L)} - y_i$

Compute errors $\{\boldsymbol{\delta}^{(L-1)}, \dots, \boldsymbol{\delta}^{(2)}\}$

Compute gradients $\Delta_{ij}^{(l)} = \Delta_{ij}^{(l)} + a_j^{(l)} \delta_i^{(l+1)}$

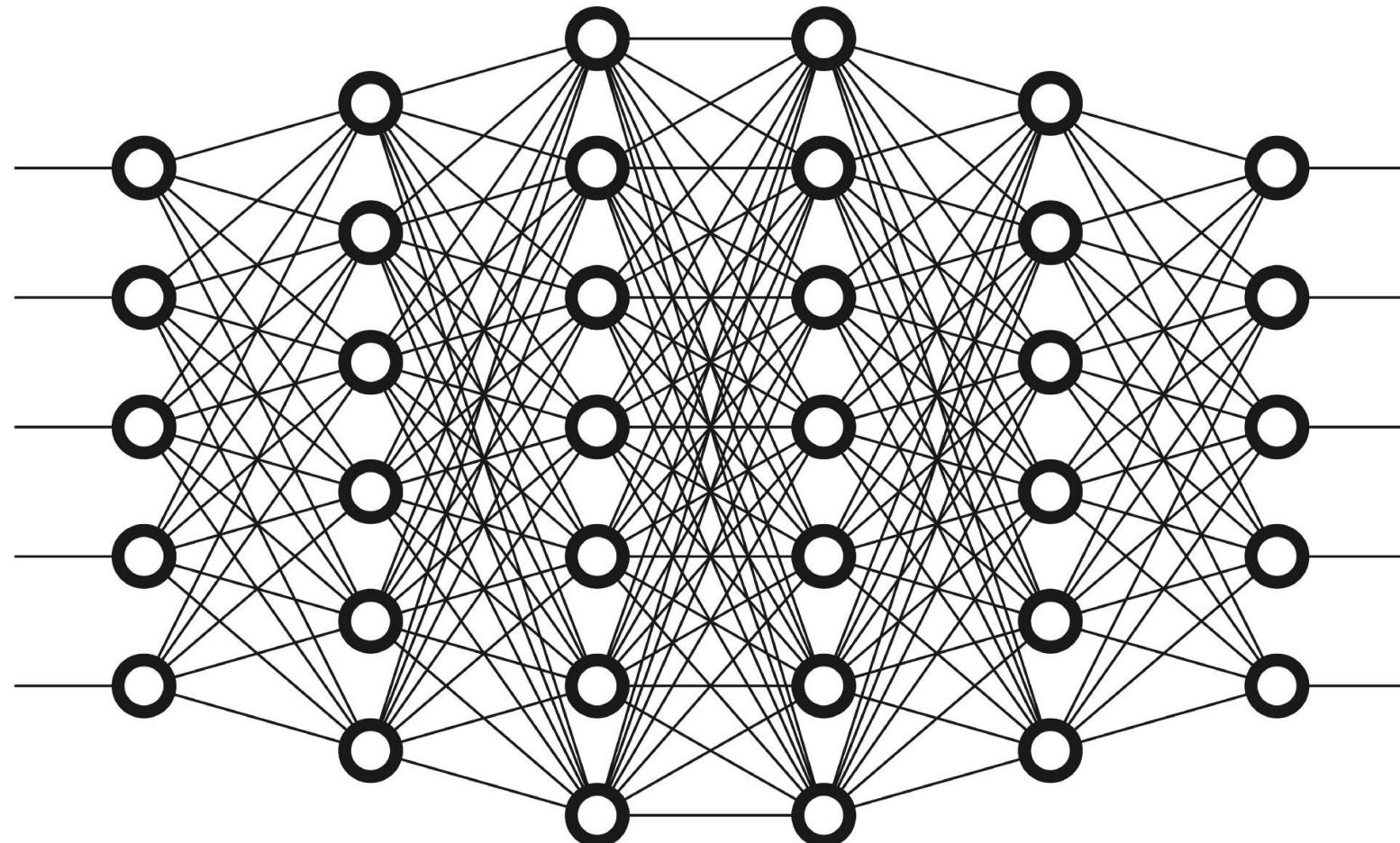
Compute avg regularized gradient $D_{ij}^{(l)} = \begin{cases} \frac{1}{n} \Delta_{ij}^{(l)} + \lambda \Theta_{ij}^{(l)} & \text{if } j \neq 0 \\ \frac{1}{n} \Delta_{ij}^{(l)} & \text{otherwise} \end{cases}$

Update weights via gradient step $\Theta_{ij}^{(l)} = \Theta_{ij}^{(l)} - \alpha D_{ij}^{(l)}$

Until weights converge or max #epochs is reached

Deep Learning

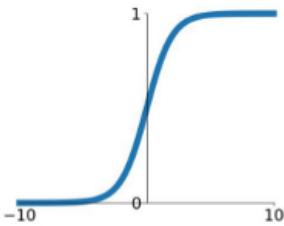
Deep Neural Network



Deep Net Activation Functions

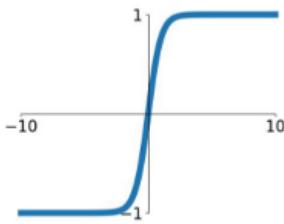
Sigmoid

$$\sigma(x) = \frac{1}{1+e^{-x}}$$



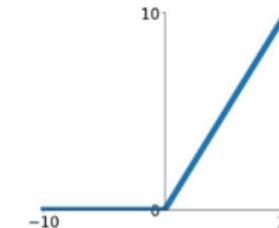
tanh

$$\tanh(x)$$



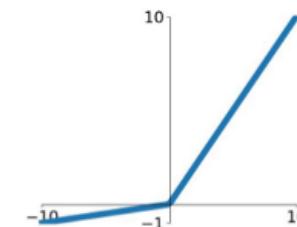
ReLU

$$\max(0, x)$$



Leaky ReLU

$$\max(0.1x, x)$$



Benefits of ReLU

- Cheap to compute as there is no complicated math and hence easier to optimize
- It converges faster.

Best Practices

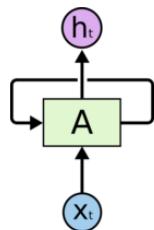
- Use **ReLU** in hidden layer activation, but be careful with the learning rate and monitor the fraction of dead units. If **ReLU** is giving problems. Try Leaky ReLU.

Beyond the Bag-of-words

- Some linguistic phenomena require going beyond the bag-of-words:
 - That's not bad for the first day
 - This is not the worst thing that can happen
 - It would be nice if you acted like you understood
 - This film should be brilliant. The actors are first grade. Stallone plays a happy, wonderful man. His sweet wife is beautiful and adores him. He has a fascinating gift for living life fully. It sounds like a great plot, however, the film is a failure.

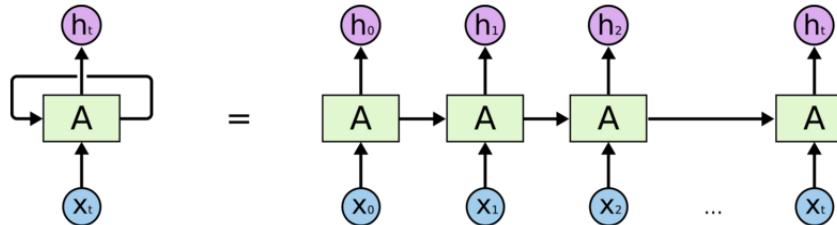
Recurrent Neural Networks

- Recurrent Neural Networks are networks with loops in them, allowing information to persist.



Recurrent Neural Networks have loops.

In the above diagram, a chunk of neural network, A , looks at some input x_t and outputs a value h_t .
A loop allows information to be passed from one step of the network to the next.



An unrolled recurrent neural network.

A recurrent neural network can be thought of as multiple copies of the same network, each passing a message to a successor.
The diagram above shows what happens if we unroll the loop.

Recurrent Neural Networks

- Intuition of Recurrent Neural Networks
 - Human thoughts have persistence; humans don't start their thinking from scratch every second.
 - As you read this sentence, you understand each word based on your understanding of previous words.
 - One of the appeals of RNNs is the idea that they are able to connect previous information to the present task
 - E.g., using previous video frames to inform the understanding of the present frame.
 - E.g., a language model tries to predict the next word based on the previous ones.

```

from tensorflow import keras
from tensorflow.keras import layers
model = keras.Sequential()
model.add(layers.Embedding(10000, 32, input_length=20))
model.add(layers.SimpleRNN(32))
model.add(layers.Dense(1, activation="sigmoid"))
model.summary()

```

Model: "sequential"

Layer (type)	Output Shape	Param #
<hr/>		
embedding (Embedding)	(None, 20, 32)	320000
simple_rnn (SimpleRNN)	(None, 32)	2080
dense (Dense)	(None, 1)	33
<hr/>		
Total params:	322,113	
Trainable params:	322,113	
Non-trainable params:	0	

An embedding is a relatively low-dimensional space into which you can translate high-dimensional vectors.

In Keras, turns positive integers (indexes) into dense vectors of fixed size.

e.g. [[4], [20]] -> [[0.25, 0.1], [0.6, -0.2]]

<https://pytorch.org/docs/stable/generated/torch.nn.Embedding.html>
https://keras.io/api/layers/core_layers/embedding/

```

from tensorflow import keras
from tensorflow.keras import layers
model = keras.Sequential()
model.add(layers.Embedding(10000, 32, input_length=20))
model.add(layers.LSTM(32))
model.add(layers.Dense(1, activation="sigmoid"))
model.summary()

```

Model: "sequential_1"

Layer (type)	Output Shape	Param #
<hr/>		
embedding_1 (Embedding)	(None, 20, 32)	320000

```

from tensorflow import keras
from tensorflow.keras import layers
model = keras.Sequential()
model.add(layers.Embedding(10000, 32, input_length=20))
model.add(layers.GRU(32))
model.add(layers.Dense(1, activation="sigmoid"))
model.summary()

```

Model: "sequential_2"

Layer (type)	Output Shape	Param #
<hr/>		
embedding_2 (Embedding)	(None, 20, 32)	320000

gru (GRU)	(None, 32)	6336
<hr/>		

dense_2 (Dense)	(None, 1)	33
<hr/>		

Total params: 326,369
 Trainable params: 326,369
 Non-trainable params: 0

Transformer Overview

Attention is all you need. 2017. Aswani, Shazeer, Parmar, Uszkoreit, Jones, Gomez, Kaiser, Polosukhin <https://arxiv.org/pdf/1706.03762.pdf>

- Non-recurrent sequence-to-sequence encoder-decoder model
- Task: machine translation with parallel corpus
- Predict each translated word
- Final cost/error function is standard cross-entropy error on top of a softmax classifier

This and related figures from paper ↑

