

# Text Classification Contd + Document Representations

Prof. Sameer Singh

---

CS 295: STATISTICAL NLP

WINTER 2017

January 17, 2017

# Outline

---

Logistic Regression

Brief Intro to Neural Networks

Document Representations

# Outline

---

Logistic Regression

Brief Intro to Neural Networks

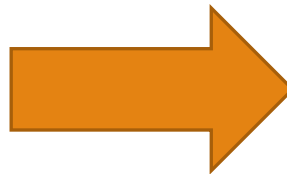
Document Representations

# Text Classification

---

Paper Title

Human machine interface for  
ABC computer applications



CS Area

- ✓ Human Computer Interaction
  - Theory
  - Artificial Intelligence
  - Systems

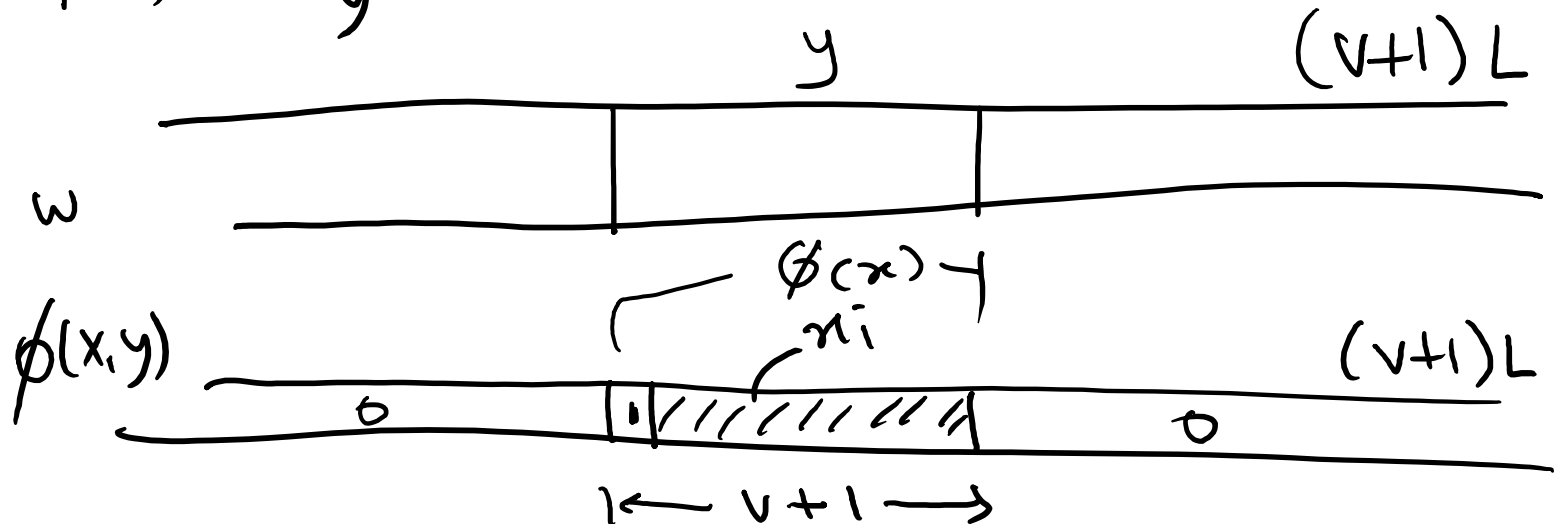
# Linear Models

Human machine interface for  
ABC computer applications

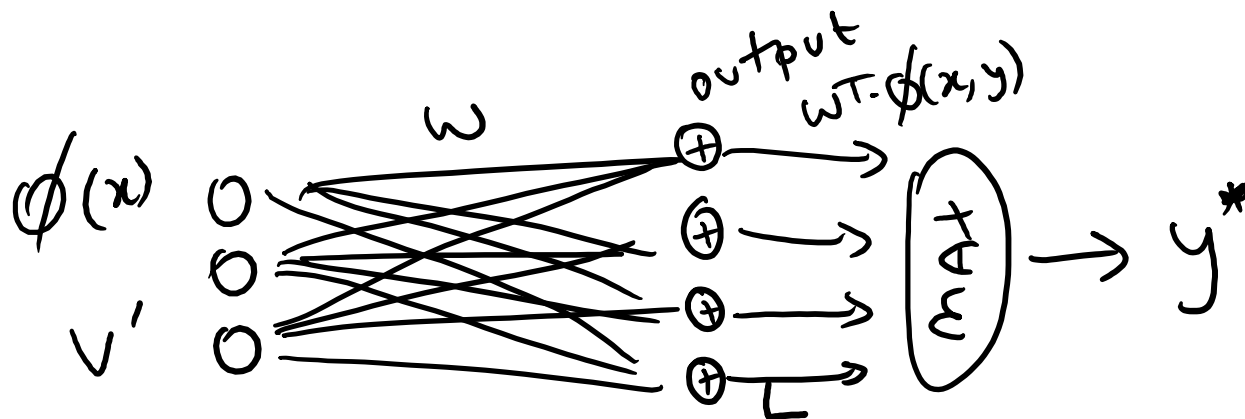
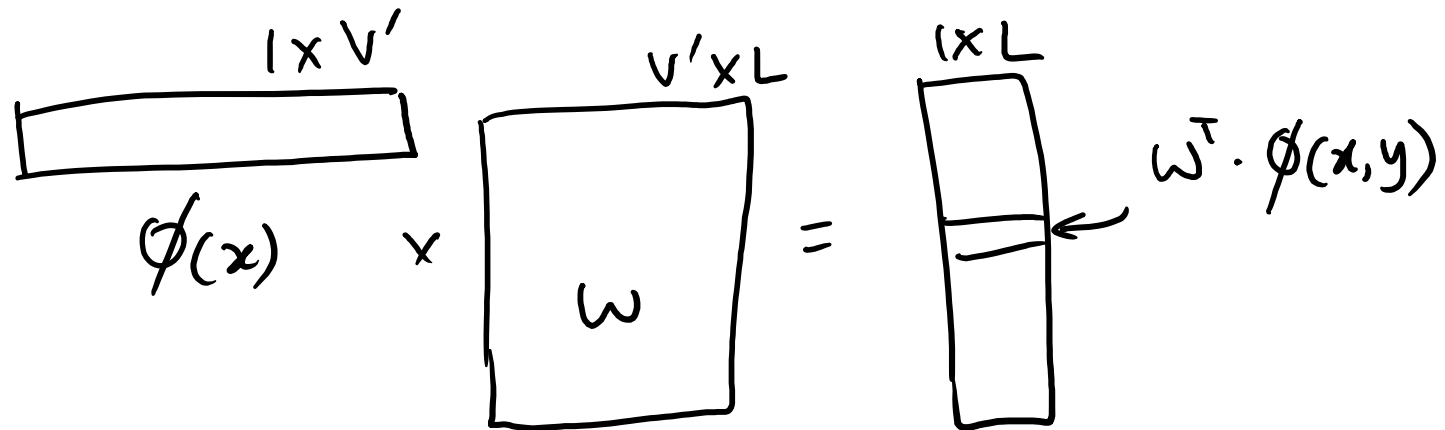
X

→  $y \in \begin{cases} \text{HCI} \\ \text{AI} \\ \text{Th} \\ \text{Sys} \end{cases} \quad L$

$$f(x) = \underset{y}{\text{arg max}} \vec{w} \cdot \phi(x, y)$$

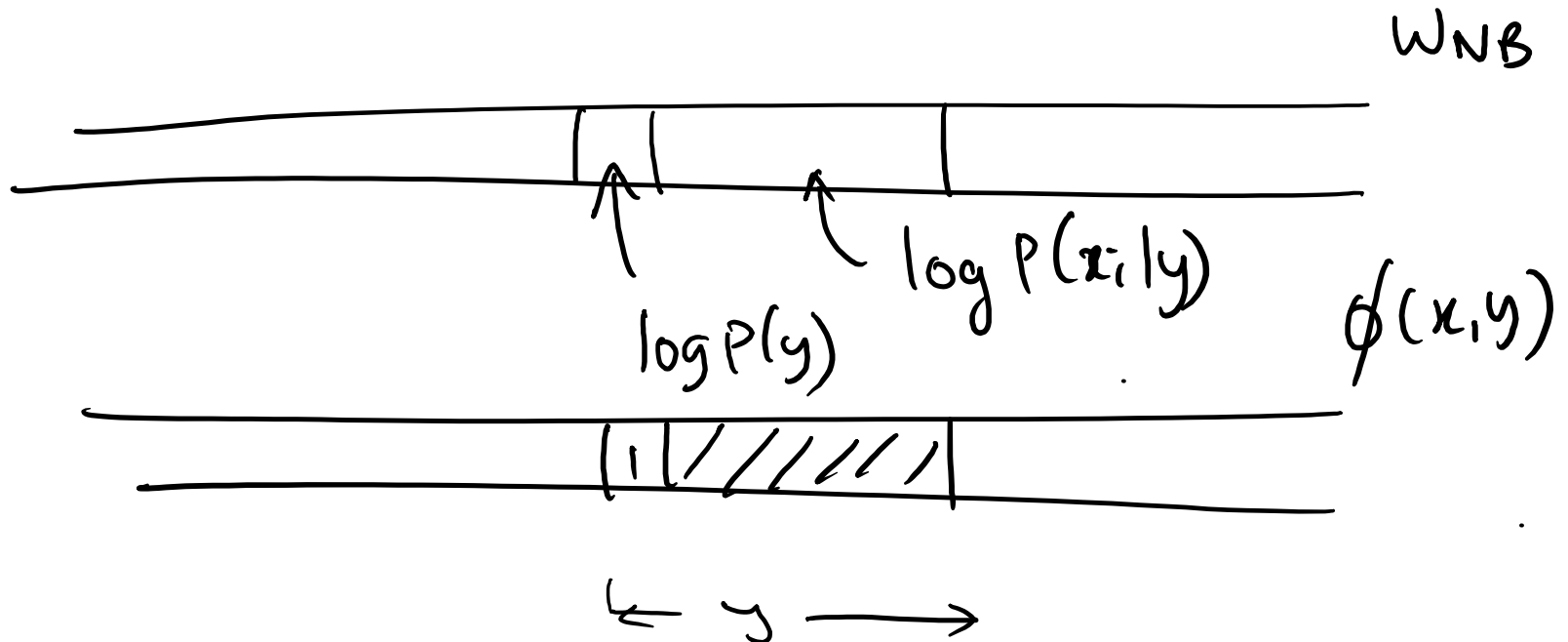


# Matrix/Neural View



# Naïve Bayes as a Linear Model

$$\begin{aligned} f_{NB}(x) &= \operatorname{argmax}_y P(x,y) = \operatorname{argmax}_y P(y) \prod_i P(x_i|y) \\ &= \operatorname{argmax}_y \log P(y) + \sum_i \log P(x_i|y) \end{aligned}$$



# Joint vs Conditional Likelihood

---

$$f_{\text{joint}}(x) = \operatorname{argmax}_y P(x, y) \leftarrow P(x)$$

Conditional

$$f_{\text{cond}}(x) = \operatorname{argmax}_y P(y|x) \leftarrow P(x)$$

$$\theta_{\text{cond}}^* = \operatorname{argmax}_{\theta} \prod_{a \in D} P_{\theta}(y_a | x_a)$$



# Logistic Regression Model

$$f(x) = \operatorname{argmax}_y w \cdot \phi(x, y)$$

$$P(y|x) = \frac{e^{w \cdot \phi(x, y)}}{\sum_y e^{w \cdot \phi(x, y)}}$$

$$f(x) = \operatorname{argmax}_y P(y|x)$$

$$= \operatorname{argmax}_y \frac{e^{w \cdot \phi(x, y)}}{\sum_{y'} e^{w \cdot \phi(x, y)}}$$

$$= \operatorname{argmax}_y w \cdot \phi(x, y)$$

(softmax)

# Logistic Regression: 2 classes

$$p(y=1|x) = \frac{e^{w \cdot \phi(x)}}{1 + e^{w \cdot \phi(x)}} \quad (\text{sigmoid})$$

$$\uparrow e^{0 \cdot \phi(\cdot)}$$

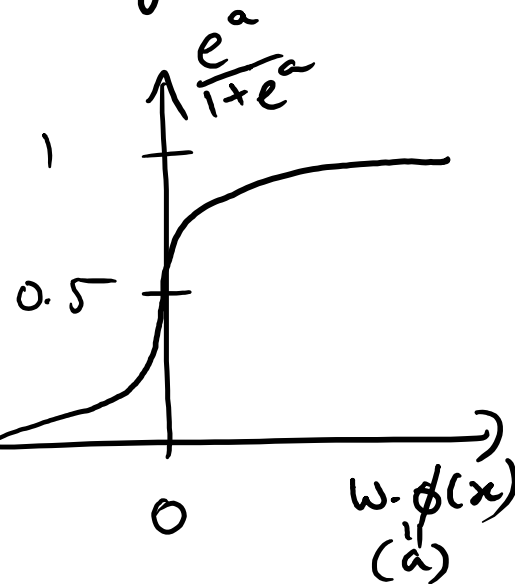
$$\phi_k(x) = \# \text{ word } k$$

→ "interface" ←

$$w_{\text{interface}} > 0$$

"tree"

$$w_{\text{tree}} < 0$$

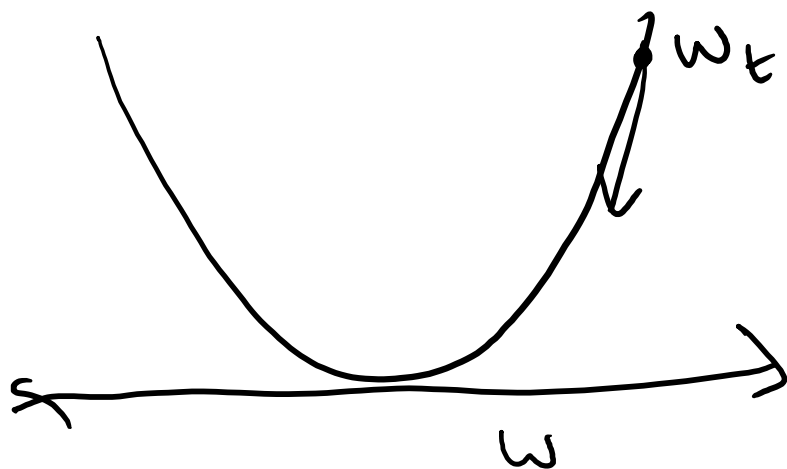


# Estimating the parameters, $\omega$

---

$$\begin{aligned}\omega^* &= \operatorname{argmax}_{\omega} \prod_a P(y_a | x_a) \\ &= \operatorname{argmax}_{\omega} \sum_a \log P(y_a | x_a) \rightarrow \frac{e^{\omega \cdot \phi(x_a, y_a)}}{\sum_y e^{\omega \cdot \phi(x_a, y)}} \\ &= \operatorname{argmax}_{\omega} \sum_a \omega \cdot \phi(x_a, y_a) - \log \sum_y e^{\omega \cdot \phi(x_a, y)} \\ &= \operatorname{argmin}_{\omega} \underbrace{\quad \quad \quad}_{\quad \quad \quad}\end{aligned}$$

# Gradient Descent



$$\operatorname{argmin}_w \underbrace{-\sum_a w \cdot \phi(x_a, y_a) - \log \sum_y \frac{e^{-L}}{\xi}}_L$$

$$\frac{\partial L}{\partial w_i} = \underbrace{\sum_a \phi_i(x_a, y_a)}_{\text{Data}} - \underbrace{E_{p(y|x)} [\phi_i(x_a, y)]}_{\text{Model}}$$

# Tips and Tricks: TF-IDF

## Sparsity of Words

- Remember Zipf's Law? Lots of rare words
- For classification, they can be more informative!

"is"  
"interface"

$N_{id}$  = # word  $i$ , doc  $d$   
(term freq (TF))

$D_i$  = # docs that  $i$  appears in

$D$  = # docs

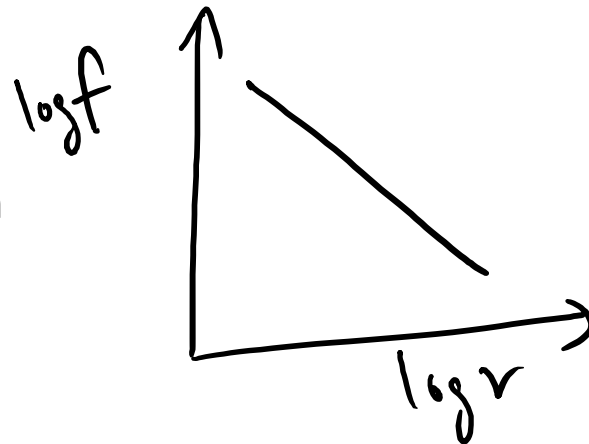
$\sum_d N_{id} > 0$

$$\phi(x_{id}) = tf(i, d) \underbrace{idf(i)}_{\text{inverse doc freq}} = -\log P(i) = -\log \frac{D_i}{D} = \log \frac{D}{D_i}$$

# Tips and Tricks: TF-IDF

## Why use log(proportion)

- It works...
- Importance is not a linear function  
"log"
- IDF is an additive function



$$\begin{aligned}idf(w_1, w_2) &= -\log P(w_1, w_2) = -\log P(w_1)P(w_2) \\ &= -\log P(w_1) - \log P(w_2) \\ &= idf(w_1) + idf(w_2)\end{aligned}$$

# Tips and Tricks: Regularization

## Overfitting

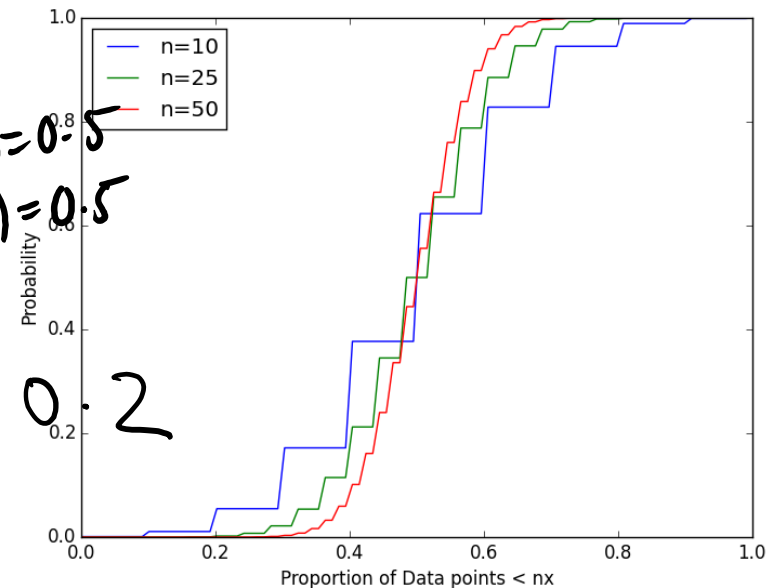
- Training data is finite: thus has spurious correlations
- Rare words that occur with one label!
  - Or don't occur often enough
  - Curse of the Zipf's Law continues

For a word that occurs 10 times...

$$w \quad p(y_1|w) = 0.5$$

$$p(y_2|w) = 0.5$$

$$P(\#_w = y_1 \geq 7) \rightarrow 0.2$$



There are many that occur ~10 times!

# Tips and Tricks: Regularization

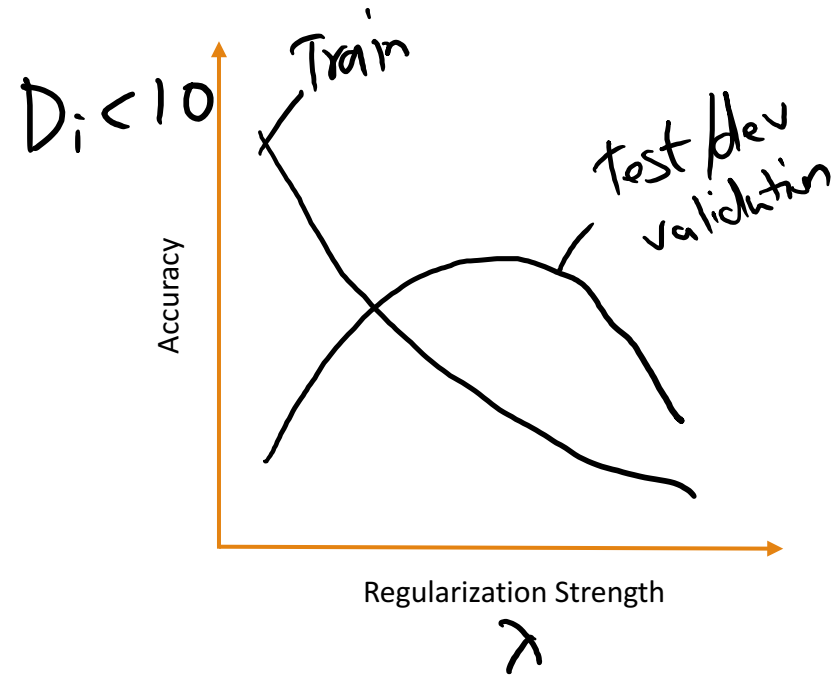
## Fixing Overfitting

- Ignore rare words (opposite of TF-IDF)
- Penalize really high weights...

$$w = \underset{w}{\operatorname{argmin}} \mathcal{L} + \lambda \|w\|_2^2$$

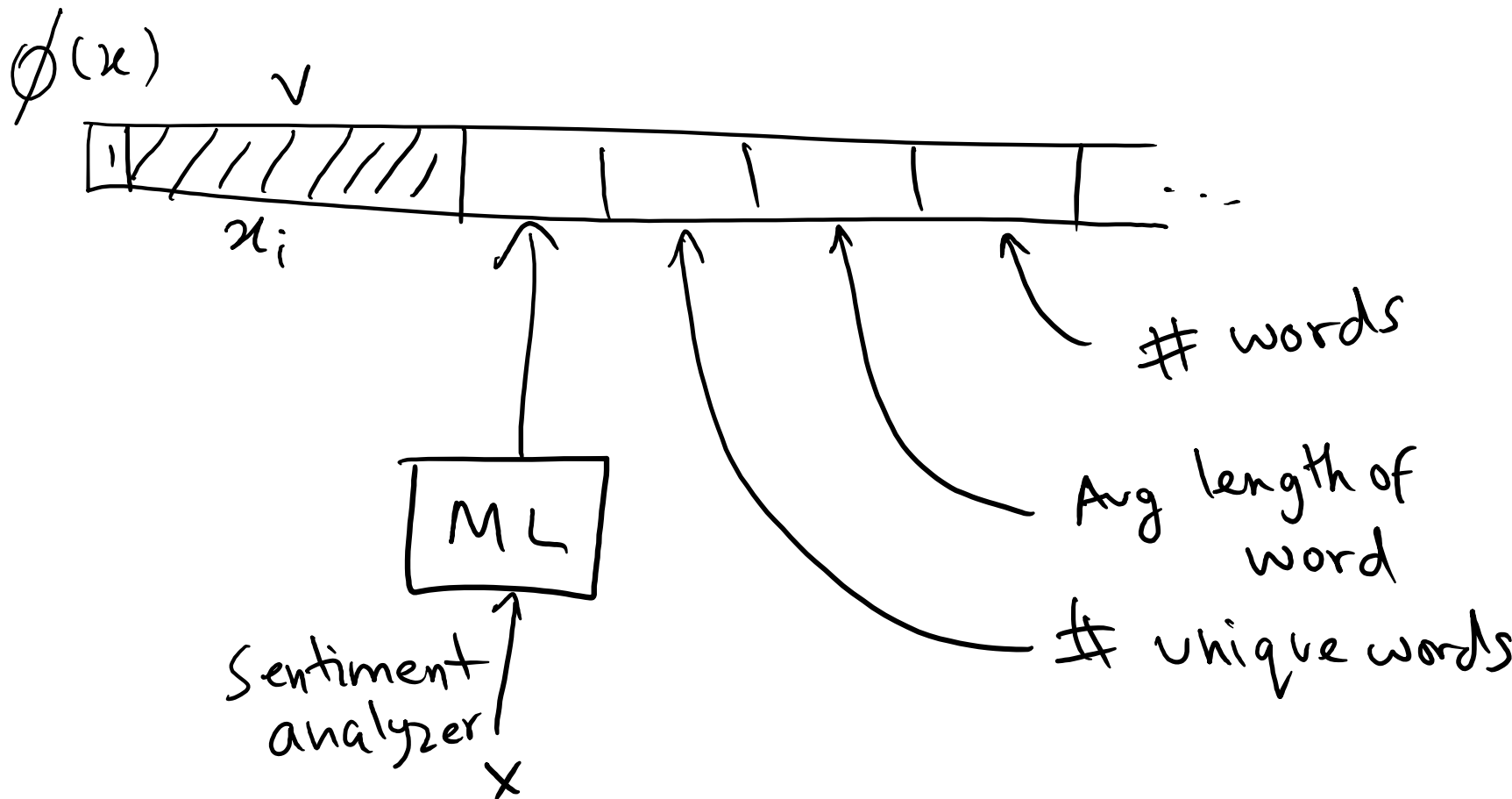
$$\frac{\partial \mathcal{L}}{\partial w_i} + 2\lambda w_i$$

Reg. Strength





# Tips and Tricks: Featurizing



# Outline

---

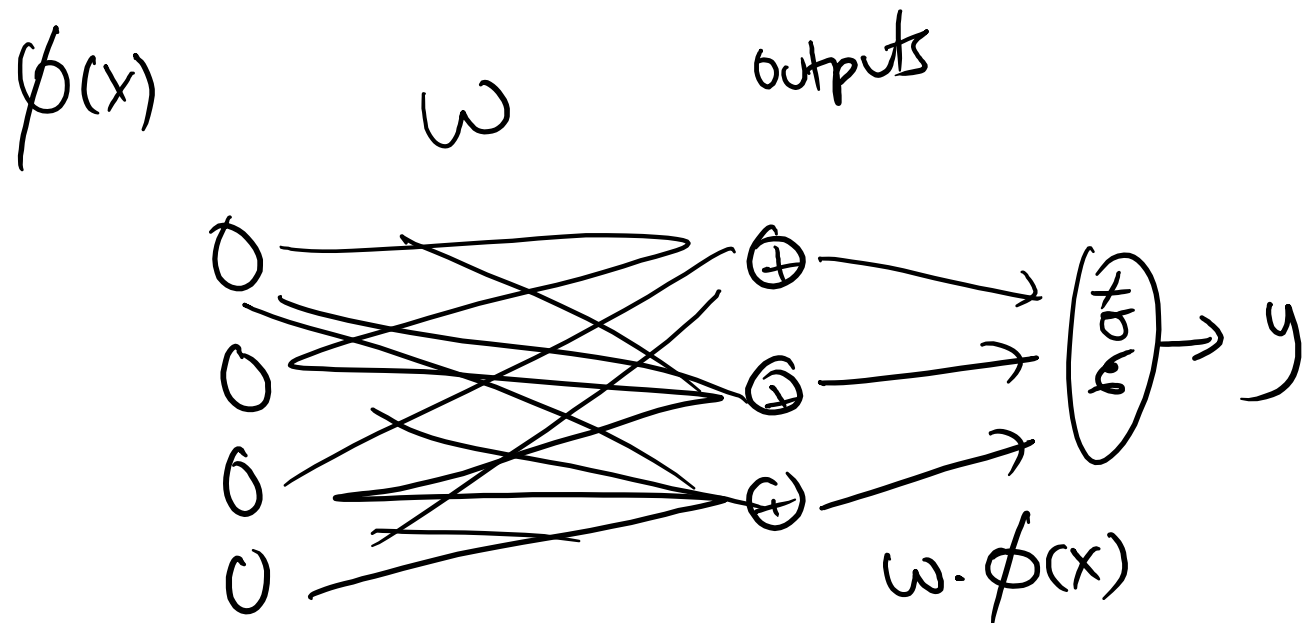
Logistic Regression

**Brief Intro to Neural Networks**

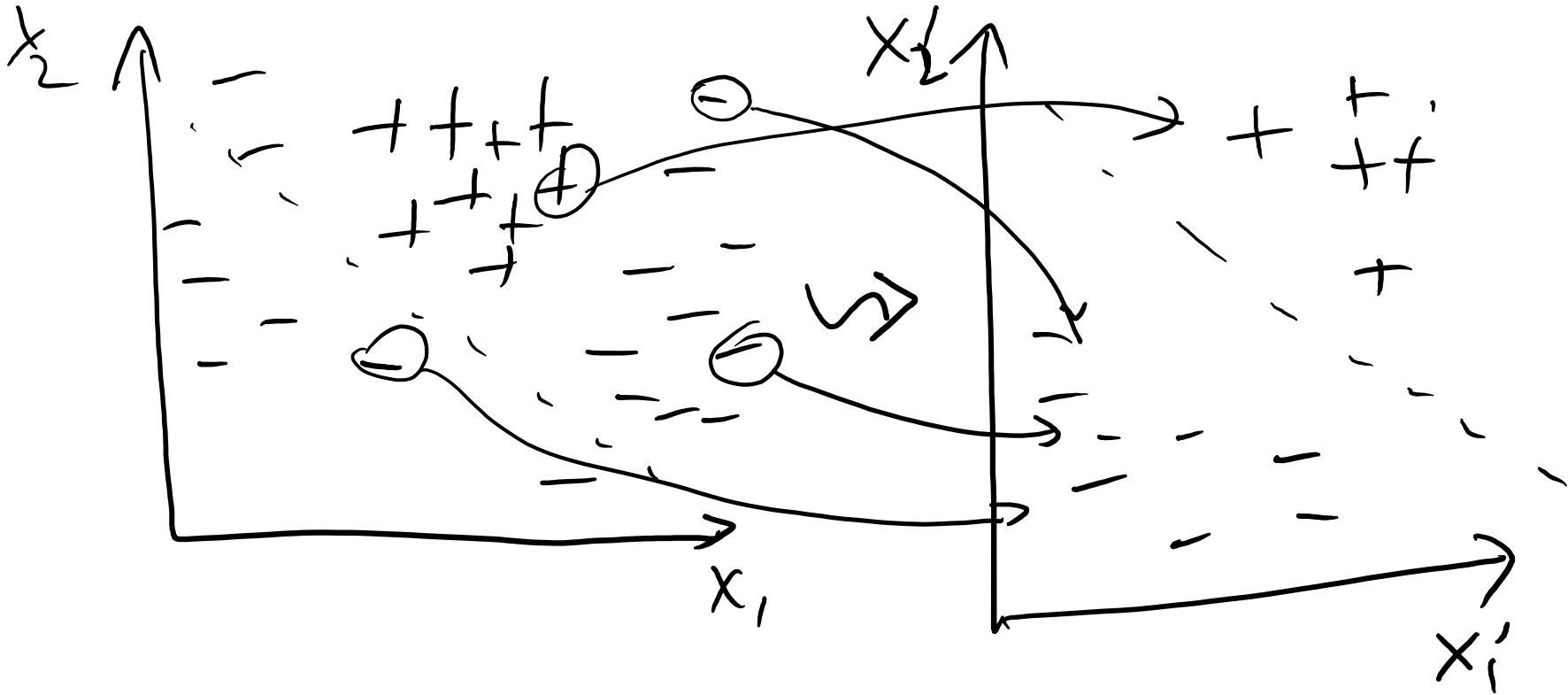
Document Representations

# Neural View of Log. Regression

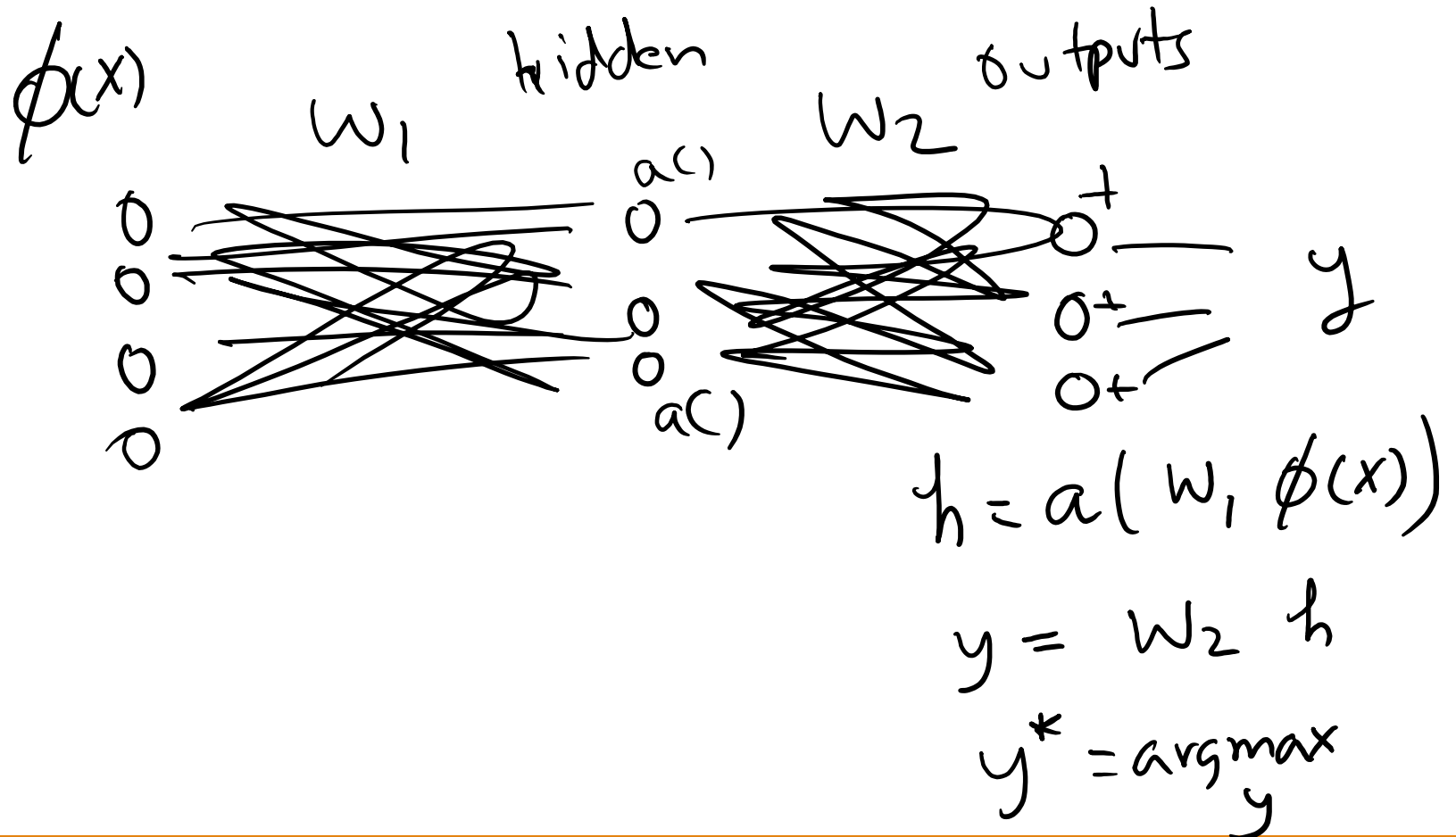
---



# Linear vs Non-linear Model



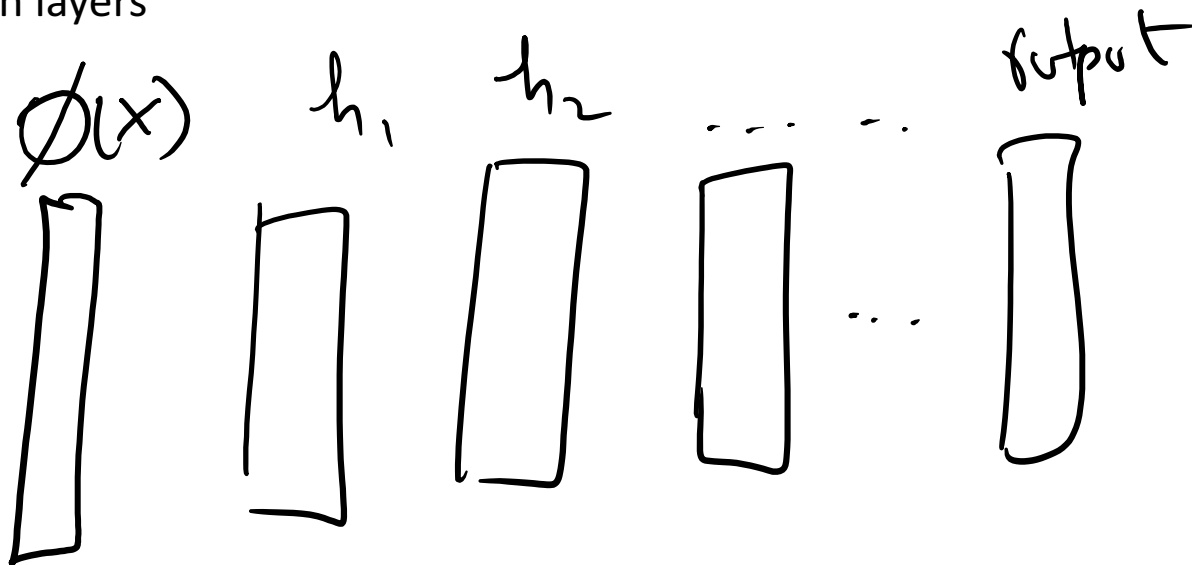
# Introducing a Hidden Layer



# What is Deep Learning?

---

Many hidden layers



In NLP, utilize unlabeled data to learn representations... (next lecture)

# Outline

---

Logistic Regression

Brief Intro to Neural Networks

Document Representations

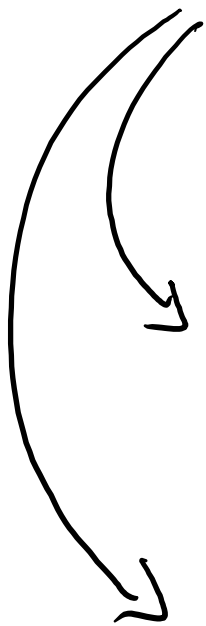
# Document Similarity

---

A survey of user opinion of computer system response time

Relation of user perceived response time to error measurement

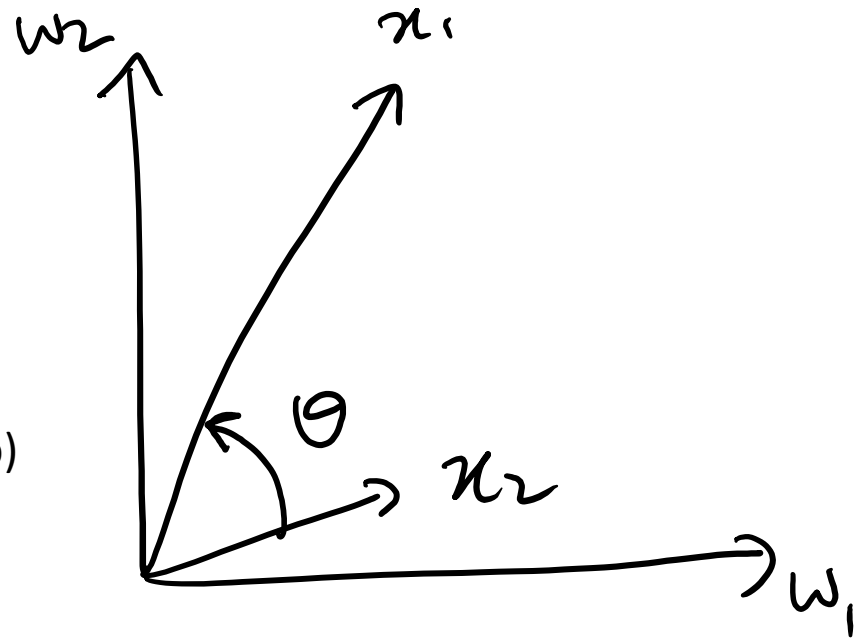
The generation of random, binary, ordered trees





# Cosine Distance

$d_1$   $d_2$   
 $\vec{x}_1$   $\vec{x}_2$



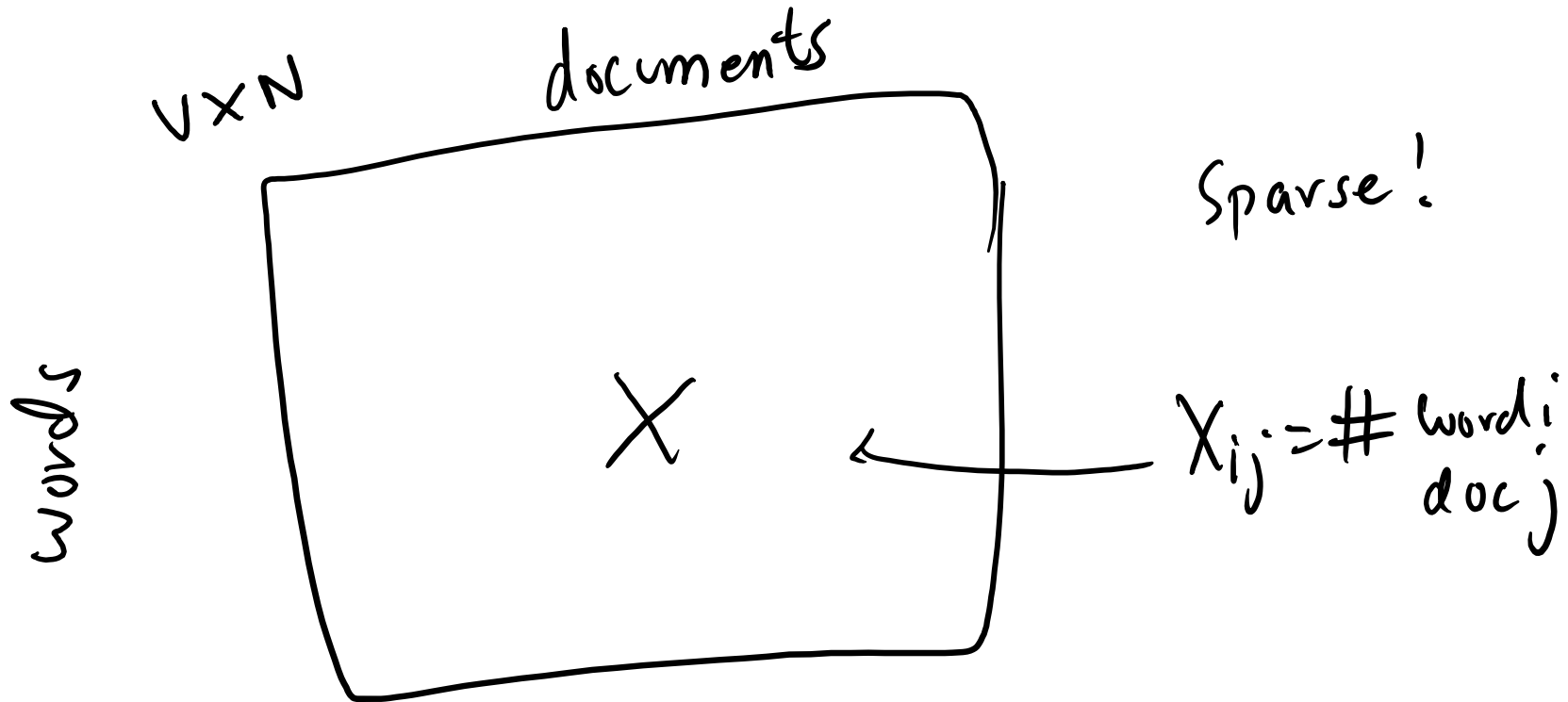
## Advantages

- Between -1 and 1 (0 means no overlap)
- If all >0, it is between 0 and 1
- Size of vectors don't matter

$$\text{dist}(x,y) = \cos \theta = \frac{x_1 \cdot x_2}{\|x_1\| \|x_2\|} = \frac{\sum_k x_{1k} \cdot x_{2k}}{\sqrt{\sum_k x_{1k}^2} \sqrt{\sum_k x_{2k}^2}}$$

# Term Document Matrix

---



# Local and Global Weighting

---

$$X_{ij} = l_{ij} g_i$$

## Local Weighting

- Binary: 1
- Term Freq:
- Log:  $\log(1 + \#_{ij})$

## Global Weighting

- Binary: 1
- Normal:  $\frac{1}{D_i}$
- IDF:  $\text{idf}(i)$

# Example: Documents

---

- c1: Human machine interface for ABC computer applications
- c2: A survey of user opinion of computer system response time
- c3: The EPS user interface management system
- c4: System and human system engineering testing of EPS
- c5: Relation of user perceived response time to error measurement

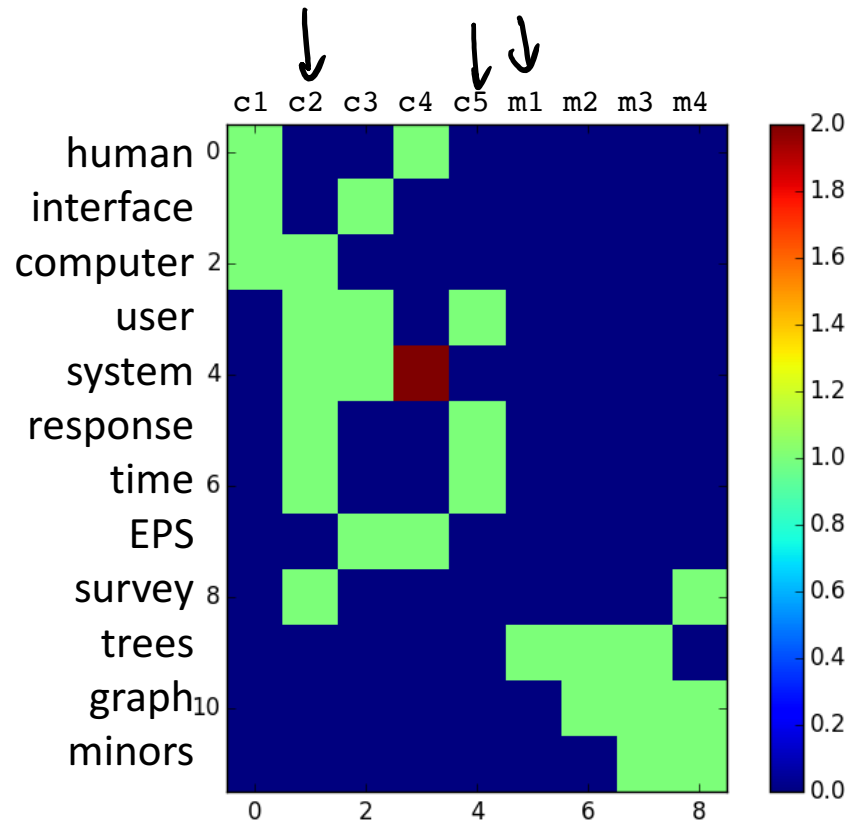
HCI

- m1: The generation of random, binary, ordered trees
- m2: The intersection graph of paths in trees
- m3: Graph minors IV: Widths of trees and well-quasi-ordering
- m4: Graph minors: A survey

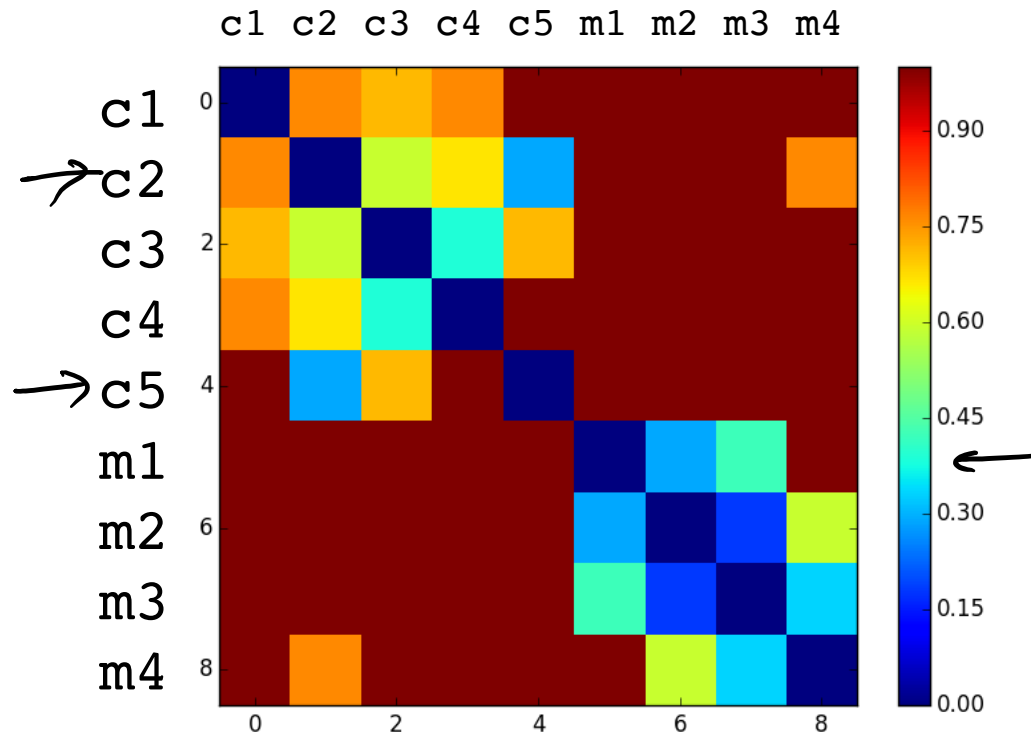
Theory

From <http://lsa.colorado.edu/papers/dp1.LSAintro.pdf>

# Example: Term-Doc Matrix



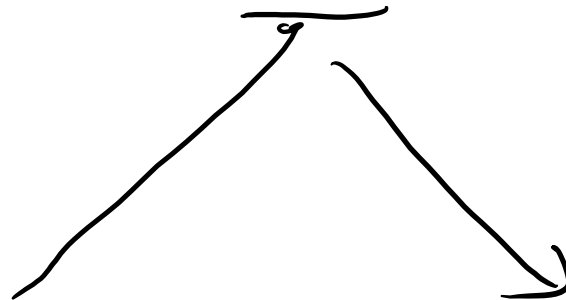
# Example: Distance Matrix



# Problems with Sparse Vectors

---

c2: A survey of user opinion of computer system response time

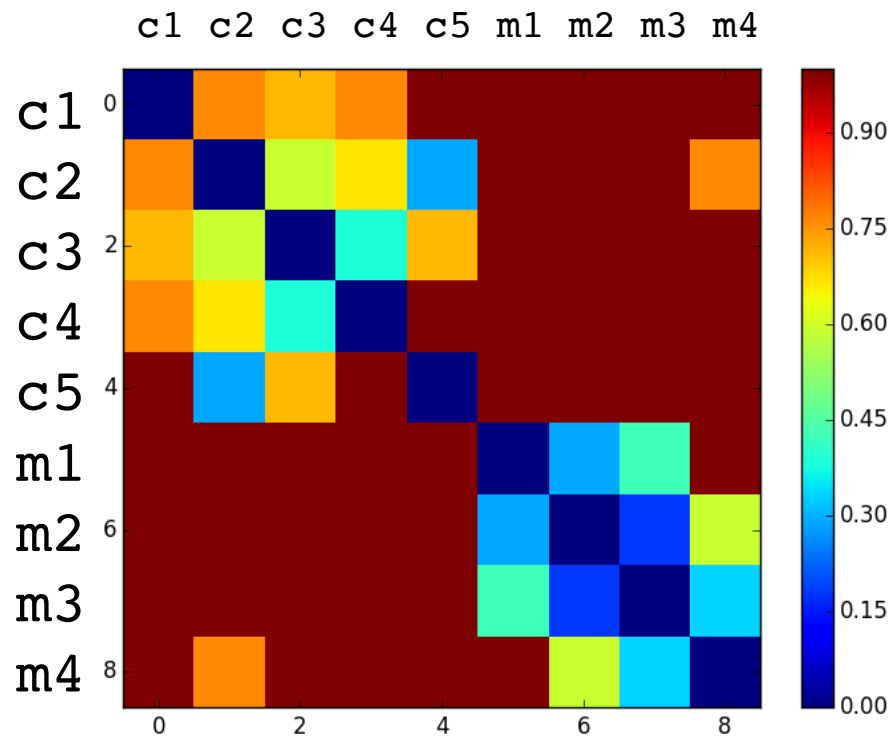


m4: Graph minors: A survey

c1: Human machine interface  
for ABC computer applications

# Example: Distance Matrix

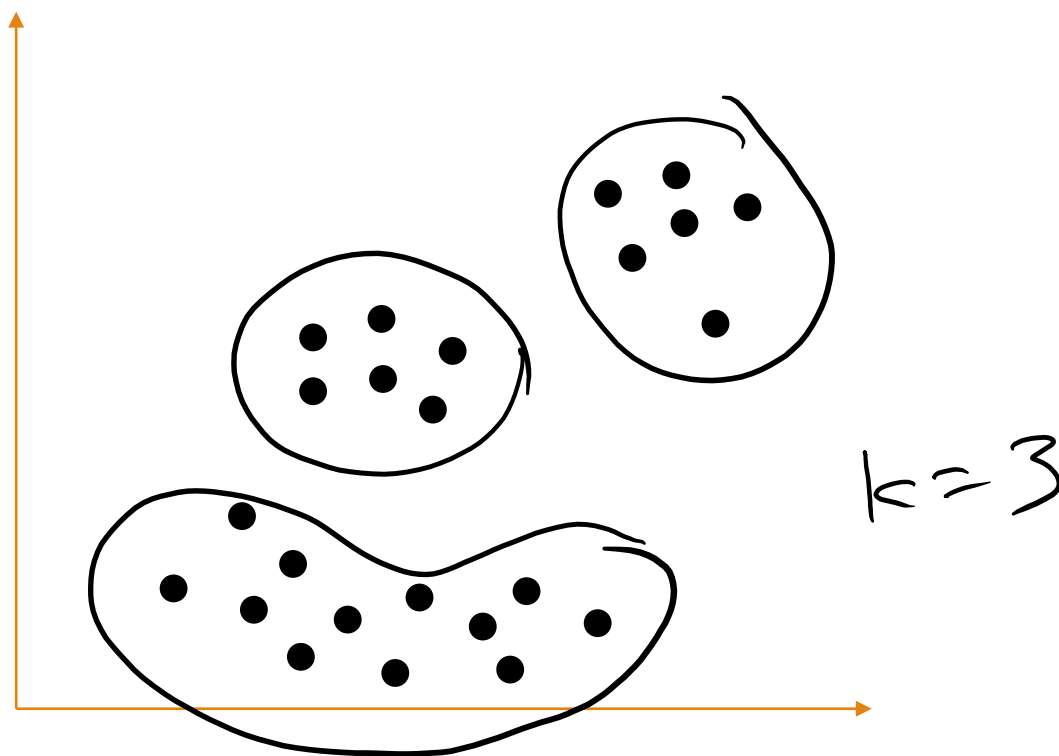
---



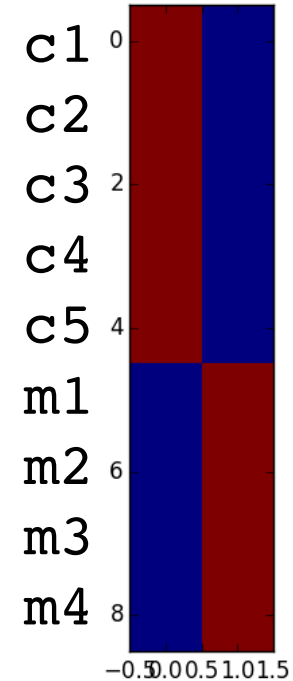
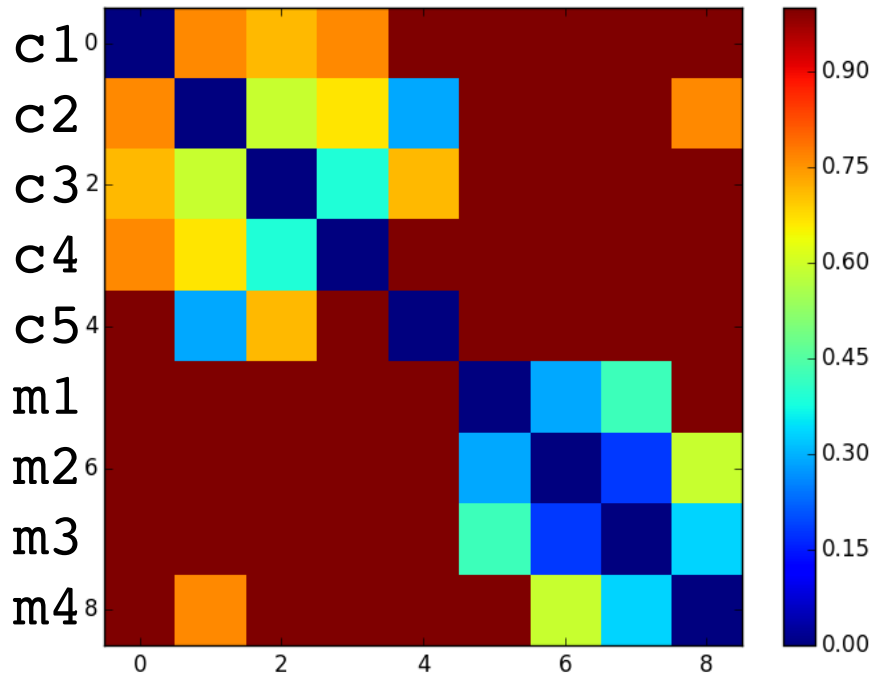
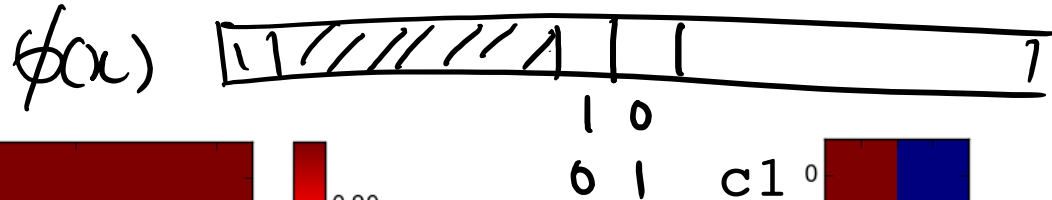


# Option 1: Clustering

---



# Example: Clustering



$k=2$

# Upcoming...

---

## Homework

- Homework 1 is up!
- No more material will be covered
- Due: **January 26, 2017**

## Project

- Project pitch is due **January 23, 2017!**
- Start assembling teams now
- Tons of datasets on the “projects” page on website