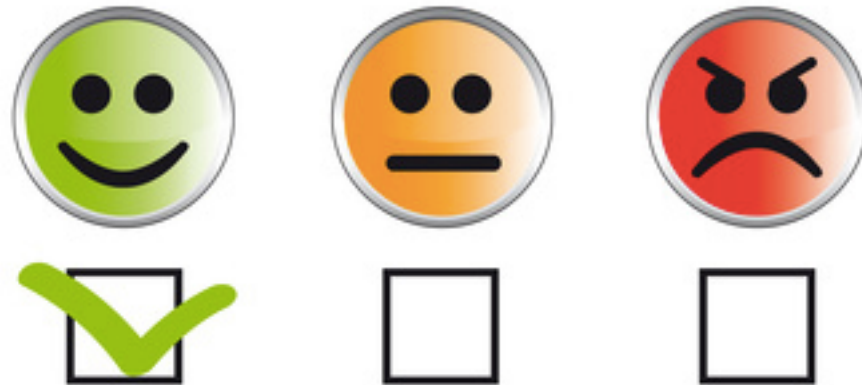


# ML4NLP

## Introduction to Classification



Dan Goldwasser

*Purdue University*

[dgoldwas@purdue.edu](mailto:dgoldwas@purdue.edu)

# Statistical Language Modeling

- **Intuition:** by looking at large quantities of text we can find statistical regularities
  - *Distinguish between correct and incorrect sentences*
- Language models define a probability distribution over strings (e.g., sentences) in a language.
- ***We can use language model to score and rank sentences***

*“I don’t know {whether,weather} to laugh or cry”*

$P(\text{“I don’t.. weather to laugh..”}) >< P(\text{“I don’t.. whether to laugh..”})$

# Language Modeling with N-grams

A language model over a vocabulary  $V$  assigns probabilities to strings drawn from  $V^*$ .

Recall the chain rule:

$$P(w_1 \dots w_i) = P(w_1)P(w_2|w_1)P(w_3|w_1w_2)\dots P(w_i|w_1 \dots w_{i-1})$$

An n-gram language model assumes each word depends only on the last n-1 words:

$$P_{n\text{gram}}(w_1 \dots w_i) := P(w_1)P(w_2|w_1)\dots P(\underbrace{w_i}_{\text{nth word}} \mid \underbrace{w_{i-n-1} \dots w_{i-1}}_{\text{prev. } n-1 \text{ words}})$$

Unigram model  $P(w_1)P(w_2)\dots P(w_i)$

Bigram model  $P(w_1)P(w_2|w_1)\dots P(w_i|w_{i-1})$

Trigram model  $P(w_1)P(w_2|w_1)\dots P(w_i|w_{i-2} w_{i-1})$

# Evaluating Language Models

- Assuming that we have a language model, how can we tell if it's good?
- **Option 1: *try to generate Shakespeare..***
  - This is know as *Qualitative evaluation*
- **Option 2: *Quantitative evaluation***
  - **Option 2.1: *See how well you do on Spelling correction***
    - This is known as ***Extrinsic Evaluation***
  - **Option 2.2: *Find an independent measure for LM quality***
    - This is known as ***Intrinsic Evaluation***

# When are LM applicable?

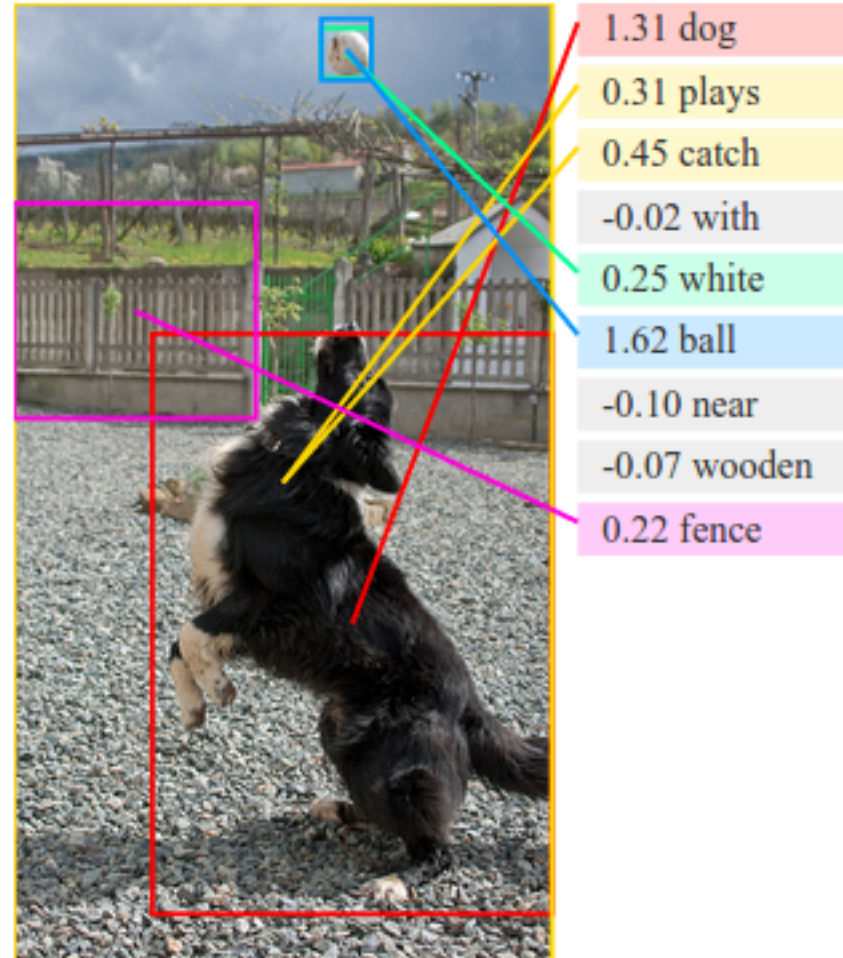
Finding regularity in language is surprisingly useful!

**Easy example:** weather/whether

But also-

**Translation** (can you produce “legal” French from source English?)

**Caption Generation** (combine output of visual sensors into a grammatical sentence)



# Classification

- *A fundamental machine learning tool*
  - Widely applicable in NLP
- **Supervised learning:** Learner is given a collection of labeled documents
  - Emails: Spam/not spam; Reviews: Pos/Neg
- Build a **function** mapping documents to labels
  - Key property: **Generalization**
    - function should work well on new data

# Sentiment Analysis

Dude, I just watched this horror flick! Selling points: nightmares scenes, torture scenes, terrible monsters that was so bad a##!

Don't buy the popcorn it was terrible, the monsters selling it must have wanted to torture me, it was so bad it gave me nightmares!

*What should your learning algorithm look at?*

# Deceptive Reviews

Which of these two hotel reviews is deceptive opinion spam?

My husband and I stayed at the James Chicago Hotel for our anniversary. This place is fantastic! We knew as soon as we arrived we made the right choice! The rooms are BEAUTIFUL and the staff very attentive and wonderful!! The area of the hotel is great, since I love to shop I couldn't ask for more!! We will definatly be back to Chicago and we will for sure be back to the James Chicago.

I have stayed at many hotels traveling for both business and pleasure and I can honestly say that The James is tops. The service at the hotel is first class. The rooms are modern and very comfortable. The location is perfect within walking distance to all of the great sights and restaurants. Highly recommend to both business travellers and couples.

*What should your learning algorithm look at?*



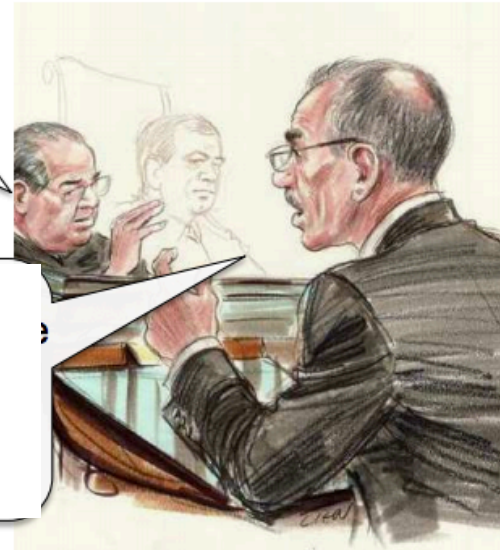
# Power Relations

Can subtle, *domain-independent* linguistic cues reveal (situational) power?

*Who's in charge?*

Blah  
Unaccep  
table  
blah

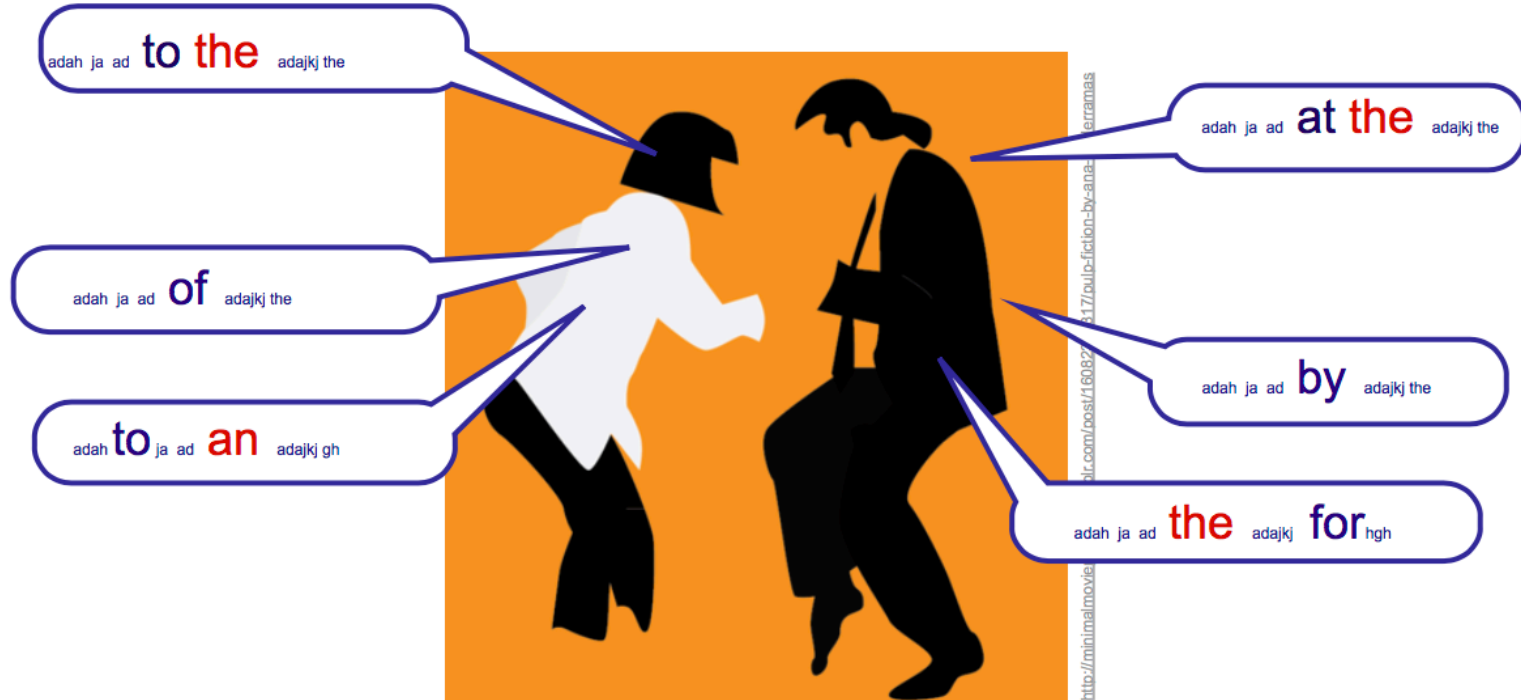
Your honor, I  
agree blah blah  
blah



*What should your learning algorithm look at?*

# Power Relations

Communicative behaviors are “patterned and coordinated, like a dance” [Niederhoffer and Pennebaker 2002]



Echoes of Power: Language Effects and Power Differences in Social Interaction. Danescu-Niculescu-Mizil et-al. WWW 2012.

# Classification

- We assume we have a labeled dataset. **How can we build a classifier?**
- Decide on a **representation** and a **learning algorithm**
  - Essentially: **Function approximation**
    - **Representation:** What is the domain of the function
    - **Learning:** How to find a **good** approximation
  - We will look into several simple examples
    - Naïve Bayes, Perceptron
- Let's start with some definitions..

# Basic Definitions

- Given:  $D$  a set of labeled examples  $\{ \langle x, y \rangle \}$
- **Goal:** *Learn a function  $f(x)$  s.t.  $f(x) = y$* 
  - **Note:**  *$y$  can be binary, or categorical*
  - Typically *the input  $x$  is represented as a vector of **features***
- Break  $D$  into three parts:
  - **Training** set (used by the learning algorithm)
  - **Test** set (evaluate the learned model)
  - **Development** set (tuning the learning algorithm)
- **Evaluation:**
  - performance measure over the test set
  - **Accuracy:** proportion of correct predictions (test data)

# Precision and Recall

- Given a dataset, we train a classifier that gets 99% accuracy
- **Did we do a good job?**
- Build a classifier for brain tumor:
  - 99.9% of brain scans do not show signs of tumor
  - *Did we do a good job?*
- By simply saying “NO” to all examples we reduce the error by a factor of 10!
  - *Clearly Accuracy is not the best way to evaluate the learning system when the data is heavily skewed!*
- **Intuition:** we need a measure that captures the class we care about! (rare)

# Precision and Recall

- The learner can make two kinds of mistakes:
  - False Positive
  - False Negative

|                     |                      |                      |
|---------------------|----------------------|----------------------|
|                     | True Label: <b>1</b> | True Label: <b>0</b> |
| Predicted: <b>1</b> | True Positive        | False Positive       |
| Predicted: <b>0</b> | False Negative       | True Negative        |

- Precision:**

- “when we predicted the rare class, how often are we right?”

$$\frac{\text{True Pos}}{\text{Predicted Pos}} = \frac{\text{True Pos}}{\text{True Pos} + \text{False Pos}}$$

- Recall**

- “Out of all the instances of the rare class, how many did we catch?”

$$\frac{\text{True Pos}}{\text{Actual Pos}} = \frac{\text{True Pos}}{\text{True Pos} + \text{False Neg}}$$

# F-Score

- Precision and Recall give us two reference points to compare learning performance

|             | Precision | Recall |
|-------------|-----------|--------|
| Algorithm 1 | 0.5       | 0.4    |
| Algorithm 2 | 0.7       | 0.1    |
| Algorithm 3 | 0.02      | 1      |

- Which algorithm is better?* **We need a single score**

- Option 1: **Average**  $\frac{P + R}{2}$
- Option 2: **F-Score**  $2 \frac{PR}{P + R}$

## Properties of f-score:

- Ranges between 0-1
- Prefers precision and recall with similar values

# Simple Example: Naïve Bayes

- **Naïve Bayes**: simple probabilistic classifier
  - Given a set of labeled data:
    - Documents **D**, each associated with a label **v**
    - Simple feature representation: BoW
  - **Learning**:
    - Construct a probability distribution  $P(v|d)$
  - **Prediction**:
    - Assign the label with the highest probability
- Relies on **strong** simplifying assumptions



# Simple Representation: BoW

- Basic idea: (sentiment analysis)
  - “I loved this movie, it’s awesome! I couldn’t stop laughing for two hours!”
  - Mapping input to label can be done by *representing the frequencies of individual words*
  - **Document = word counts**
- ***Simple, yet surprisingly powerful representation!***

# Bayes Rule

- Naïve Bayes is a simple probabilistic classification method, based on Bayes rule.

$$P(\mathbf{v} | \mathbf{d}) = P(\mathbf{d} | \mathbf{v}) \frac{P(\mathbf{v})}{P(\mathbf{d})}$$

# Basics of Naïve Bayes

- $P(v)$  - the **prior probability** of a label  $v$   
Reflects background knowledge; before data is observed. If no information - uniform distribution.
- $P(D)$  - The probability that this sample of the Data is observed. (*No knowledge of the label*)
- $P(D | v)$ : The probability of observing the sample  $D$ , given that the label  $v$  is the target (*Likelihood*)  
—
- $P(v | D)$ : The **posterior probability** of  $v$ . The probability that  $v$  is the target, given that  $D$  has been observed.

# Bayes Rule

- Naïve Bayes is a simple classification method, based on Bayes rule.

$$P(\mathbf{v} \mid \mathbf{d}) = P(\mathbf{d} \mid \mathbf{v}) \frac{P(\mathbf{v})}{P(\mathbf{d})}$$

*Check your intuition:*

$P(\mathbf{v} \mid \mathbf{d})$  increases with  $P(\mathbf{v})$  and with  $P(\mathbf{d} \mid \mathbf{v})$

$P(\mathbf{v} \mid \mathbf{d})$  decreases with  $P(\mathbf{d})$

# Naïve Bayes

$$P(v | D) = P(D | v) P(v) / P(D)$$

- The learner considers a set of candidate labels, and attempts to find the most probable one  $v \in V$ , given the observed data.
- Such maximally probable assignment is called maximum a posteriori assignment (MAP); Bayes theorem is used to compute it:

$$v_{MAP} = \operatorname{argmax}_{v \in V} P(v | D) = \operatorname{argmax}_{v \in V} P(D | v) P(v) / P(D)$$

$$= \operatorname{argmax}_{v \in V} P(D | v) P(v)$$

Since  $P(D)$  is the same for all  $v \in V$

# Naïve Bayes

- How can we compute  $P(v | D)$ ?
  - Basic idea: represent document as a set of features, such as BoW features

$$\mathbf{v}_{\text{MAP}} = \operatorname{argmax}_{\mathbf{v}_j \in V} \mathbf{P}(\mathbf{v}_j | \mathbf{x}) = \operatorname{argmax}_{\mathbf{v}_j \in V} \mathbf{P}(\mathbf{v}_j | \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n)$$

$$\begin{aligned} \mathbf{v}_{\text{MAP}} &= \operatorname{argmax}_{\mathbf{v}_j \in V} \frac{\mathbf{P}(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n | \mathbf{v}_j) \mathbf{P}(\mathbf{v}_j)}{\mathbf{P}(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n)} \\ &= \operatorname{argmax}_{\mathbf{v}_j \in V} \mathbf{P}(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n | \mathbf{v}_j) \mathbf{P}(\mathbf{v}_j) \end{aligned}$$

# NB: Parameter Estimation

$$V_{MAP} = \operatorname{argmax}_v P(x_1, x_2, \dots, x_n | v) P(v)$$

- Given training data we can estimate the two terms
- Estimating  $P(v)$  is **easy**. For each value  $v$  count how many times it appears in the training data.

**Question:** Assume binary  $x_i$ 's. How many parameters does the model require?

- **However, it is not feasible to estimate  $P(x_1, \dots, x_n | v)$** 
  - In this case we have to estimate, for each target value, the probability of each instance (most of which will not occur)
- In order to use a Bayesian classifiers in practice, we *need to make assumptions* that will allow us to estimate these quantities.

# NB: Independence Assumption

- **Bag of words representation:**
  - Word position can be ignored
- **Conditional Independence:** Assume feature probabilities are independent given the label
  - $P(x_i | v_j) = P(x_i | x_{i-1}; v_j)$
- Both assumptions are **not true**
  - Help simplify the model
  - *Simple models work well*

*how many features are needed now?*



# Naive Bayes

$$V_{MAP} = \operatorname{argmax}_v P(x_1, x_2, \dots, x_n | v)P(v)$$

$$P(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n | \mathbf{v}_j) =$$

$$= P(\mathbf{x}_1 | \mathbf{x}_2, \dots, \mathbf{x}_n, \mathbf{v}_j)P(\mathbf{x}_2, \dots, \mathbf{x}_n | \mathbf{v}_j)$$

$$= P(\mathbf{x}_1 | \mathbf{x}_2, \dots, \mathbf{x}_n, \mathbf{v}_j)P(\mathbf{x}_2 | \mathbf{x}_3, \dots, \mathbf{x}_n, \mathbf{v}_j)P(\mathbf{x}_3, \dots, \mathbf{x}_n | \mathbf{v}_j)$$

$$= \dots\dots$$

$$= P(\mathbf{x}_1 | \mathbf{x}_2, \dots, \mathbf{x}_n, \mathbf{v}_j)P(\mathbf{x}_2 | \mathbf{x}_3, \dots, \mathbf{x}_n, \mathbf{v}_j)P(\mathbf{x}_3 | \mathbf{x}_4, \dots, \mathbf{x}_n, \mathbf{v}_j) \dots P(\mathbf{x}_n | \mathbf{v}_j)$$

**Assumption:** feature values are independent given the target value

$$= \prod_{i=1}^n P(\mathbf{x}_i | \mathbf{v}_j)$$

# Estimating Probabilities (MLE)

Assume a document classification problem, using word features

$$\mathbf{v}_{\text{NB}} = \operatorname{argmax}_{\mathbf{v} \in \{\text{like}, \text{dislike}\}} \mathbf{P}(\mathbf{v}) \prod_i \mathbf{P}(\mathbf{x}_i = \text{word}_i \mid \mathbf{v})$$

*How do we estimate  $\mathbf{P}(\text{word}_k \mid \mathbf{v})$ ?*

$$\mathbf{P}(\text{word}_k \mid \mathbf{v}) = \frac{\#(\text{word}_k \text{ appears in training in } \mathbf{v} \text{ documents})}{\#(\mathbf{v} \text{ documents})} = \frac{\mathbf{n}_k}{\mathbf{n}}$$

## Sparsity of data is a problem

- if  $\mathbf{n}$  is small, the estimate is not accurate
- if  $\mathbf{n}_k$  is 0, it will dominate the estimate: we will never predict  $\mathbf{v}$  if a word that never appeared in training (with  $\mathbf{v}$ ) appears in the test data

# Robust Estimation of Probabilities

$$\mathbf{v}_{\text{NB}} = \operatorname{argmax}_{\mathbf{v} \in \{\text{like}, \text{dislike}\}} \mathbf{P}(\mathbf{v}) \prod_i \mathbf{P}(\mathbf{x}_i = \text{word}_i \mid \mathbf{v})$$

- This process is called **smoothing**.
- There are many ways to do it, some better justified than others;
- An empirical issue.

$$\mathbf{P}(\mathbf{x}_k \mid \mathbf{v}) = \frac{\mathbf{n}_k + \mathbf{m}p}{\mathbf{n} + \mathbf{m}}$$

- **Here:**
  - $n_k$  is #(of occurrences of the word in the presence of  $\mathbf{v}$ )
  - $n$  is #(of occurrences of the label  $\mathbf{v}$ )
  - $p$  is a *prior estimate of  $\mathbf{v}$*  (e.g., uniform)
  - $m$  is *equivalent sample size* (# of labels)
- **Laplace Rule:** for the Boolean case,  $p=1/2$ ,  $m=2$

$$\mathbf{P}(\mathbf{x}_k \mid \mathbf{v}) = \frac{\mathbf{n}_k + 1}{\mathbf{n} + 2}$$

# Naïve Bayes

- **Very easy to implement**
- **Converges very quickly**
  - Learning is just counting
- **Performs well in practice**
  - Applied to many document classification tasks
  - If data set is small, NB can perform better than sophisticated algorithms
- **Strong independence assumptions**
  - If assumptions hold: **NB is the optimal classifier**
  - Even if not, can perform well
- **Next: *from NB to learning linear threshold functions***

# Naïve Bayes: Two Classes

- *Notice that the naïve Bayes method gives a method for predicting rather than an explicit classifier*
- In the case of two classes,  $v \in \{0,1\}$  we predict that  $v=1$  iff:

$$\frac{P(v_j = 1) \cdot \prod_{i=1}^n P(x_i | v_j = 1)}{P(v_j = 0) \cdot \prod_{i=1}^n P(x_i | v_j = 0)} > 1$$

# Naïve Bayes: Two Classes

- Notice that the naïve Bayes method gives a method for predicting rather than an explicit classifier.
- In the case of two classes,  $v \in \{0,1\}$  we predict that  $v=1$  iff:

$$\frac{P(v_j = 1) \cdot \prod_{i=1}^n P(x_i | v_j = 1)}{P(v_j = 0) \cdot \prod_{i=1}^n P(x_i | v_j = 0)} > 1$$

Denote:  $p_i = P(x_i = 1 | v = 1)$ ,  $q_i = P(x_i = 1 | v = 0)$

$$\frac{P(v_j = 1) \cdot \prod_{i=1}^n p_i^{x_i} (1 - p_i)^{1-x_i}}{P(v_j = 0) \cdot \prod_{i=1}^n q_i^{x_i} (1 - q_i)^{1-x_i}} > 1$$

# Naïve Bayes: Two Classes

In the case of two classes,  $v \in \{0,1\}$  we predict that  $v=1$  iff:

$$\frac{P(v_j = 1) \cdot \prod_{i=1}^n p_i^{x_i} (1-p_i)^{1-x_i}}{P(v_j = 0) \cdot \prod_{i=1}^n q_i^{x_i} (1-q_i)^{1-x_i}} = \frac{P(v_j = 1) \cdot \prod_{i=1}^n (1-p_i) \left(\frac{p_i}{1-p_i}\right)^{x_i}}{P(v_j = 0) \cdot \prod_{i=1}^n (1-q_i) \left(\frac{q_i}{1-q_i}\right)^{x_i}} > 1$$

# Naïve Bayes: Two Classes

In the case of two classes,  $v \in \{0,1\}$  we predict that  $v=1$  iff:

$$\frac{P(v_j = 1) \cdot \prod_{i=1}^n p_i^{x_i} (1-p_i)^{1-x_i}}{P(v_j = 0) \cdot \prod_{i=1}^n q_i^{x_i} (1-q_i)^{1-x_i}} = \frac{P(v_j = 1) \cdot \prod_{i=1}^n (1-p_i) \left(\frac{p_i}{1-p_i}\right)^{x_i}}{P(v_j = 0) \cdot \prod_{i=1}^n (1-q_i) \left(\frac{q_i}{1-q_i}\right)^{x_i}} > 1$$

**Take logarithm; we predict  $v = 1$  iff :**

$$\log \frac{P(v_j = 1)}{P(v_j = 0)} + \sum_i \log \frac{1-p_i}{1-q_i} + \sum_i \left( \log \frac{p_i}{1-p_i} - \log \frac{q_i}{1-q_i} \right) x_i > 0$$



# Naïve Bayes: Two Classes

In the case of two classes,  $v \in \{0,1\}$  we predict that  $v=1$  iff:

$$\frac{P(v_j = 1) \cdot \prod_{i=1}^n p_i^{x_i} (1-p_i)^{1-x_i}}{P(v_j = 0) \cdot \prod_{i=1}^n q_i^{x_i} (1-q_i)^{1-x_i}} = \frac{P(v_j = 1) \cdot \prod_{i=1}^n (1-p_i) \left(\frac{p_i}{1-p_i}\right)^{x_i}}{P(v_j = 0) \cdot \prod_{i=1}^n (1-q_i) \left(\frac{q_i}{1-q_i}\right)^{x_i}} > 1$$

**Take logarithm; we predict  $v = 1$  iff :**

$$\log \frac{P(v_j = 1)}{P(v_j = 0)} + \sum_i \log \frac{1-p_i}{1-q_i} + \sum_i \left( \log \frac{p_i}{1-p_i} - \log \frac{q_i}{1-q_i} \right) x_i > 0$$

• **We get that naive Bayes is a linear separator with :**

$$w_i = \log \frac{p_i}{1-p_i} - \log \frac{q_i}{1-q_i} = \log \frac{p_i}{q_i} \frac{1-q_i}{1-p_i}$$

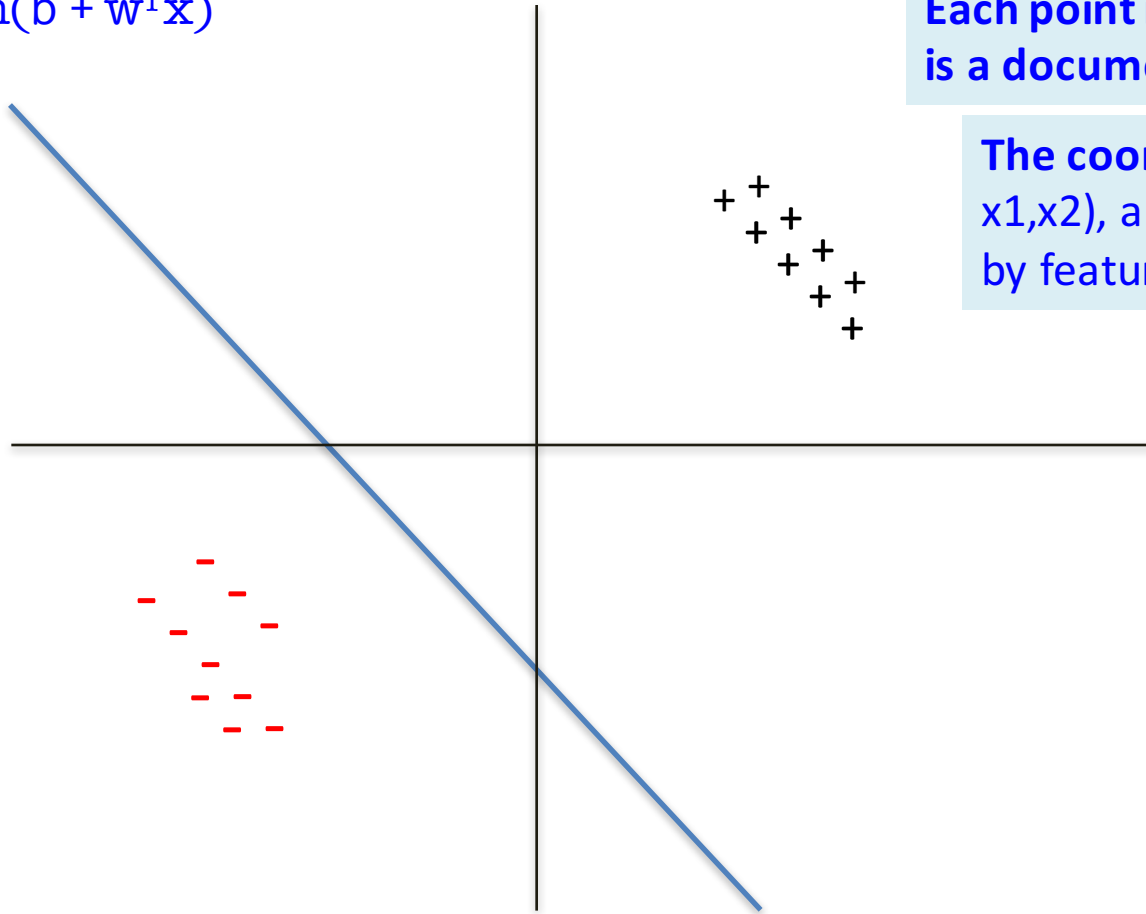
if  $p_i = q_i$  then  $w_i = 0$  and the feature is irrelevant

# Linear Classifiers

- Linear threshold functions
  - Associate a weight ( $w_i$ ) with each feature ( $x_i$ )
  - **Prediction:**  $\text{sign}(b + w^T \mathbf{x}) = \text{sign}(b + \sum w_i x_i)$ 
    - $b + w^T \mathbf{x} \geq 0$  predict  $y=1$
    - Otherwise, predict  $y=-1$
- *NB is a linear threshold function*
  - Weight vector ( $w$ ) is assigned by computing conditional probabilities
- *In fact, Linear threshold functions are a very popular representation!*

# Linear Classifiers

$$\text{sign}(b + w^T x)$$

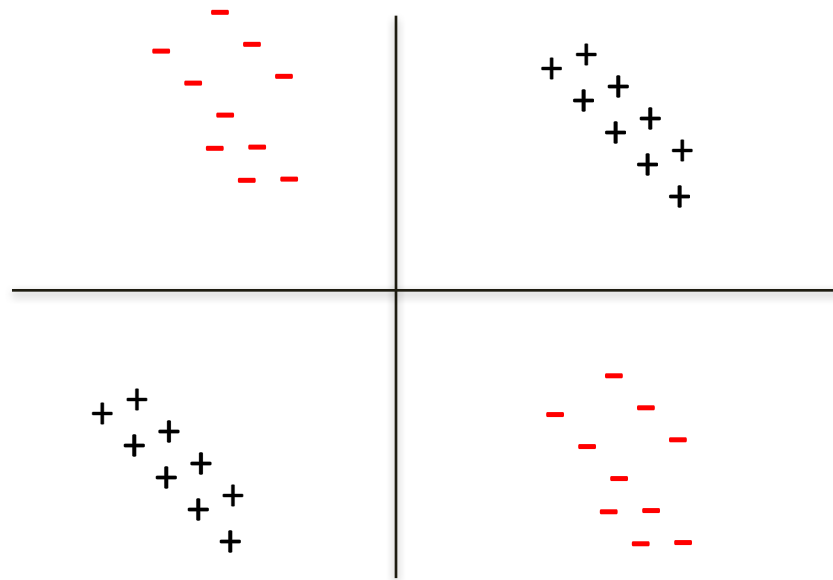


Each point in this space is a document.

The coordinates (e.g.,  $x_1, x_2$ ), are determined by feature activations

# Expressivity

- Linear functions are quite expressive
  - *Exists a linear function that is consistent with the data*
- A famous negative examples (XOR):

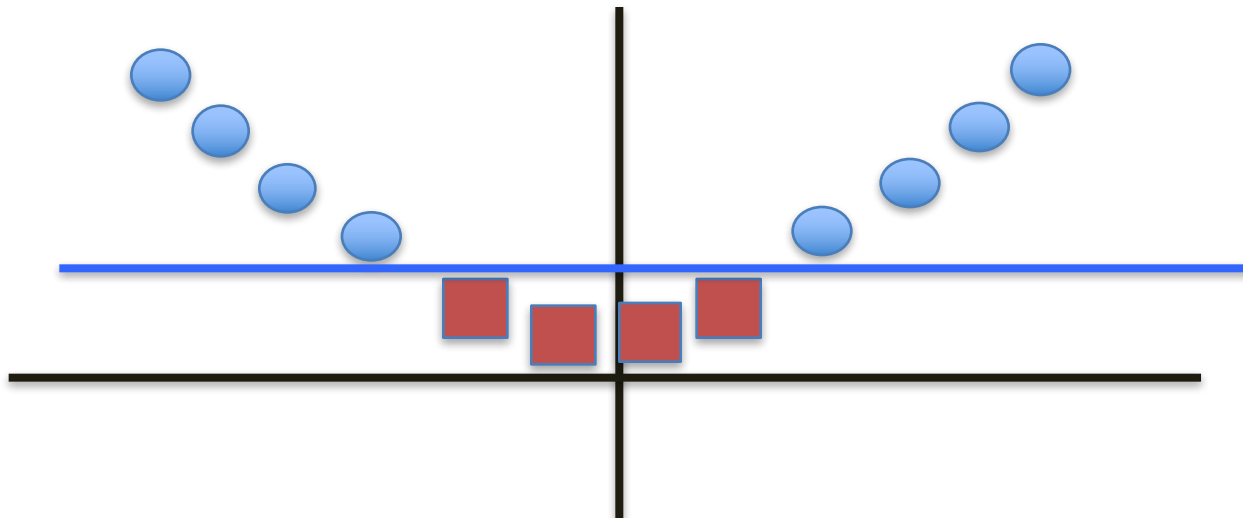


# Expressivity

By transforming the feature space these functions can be made linear

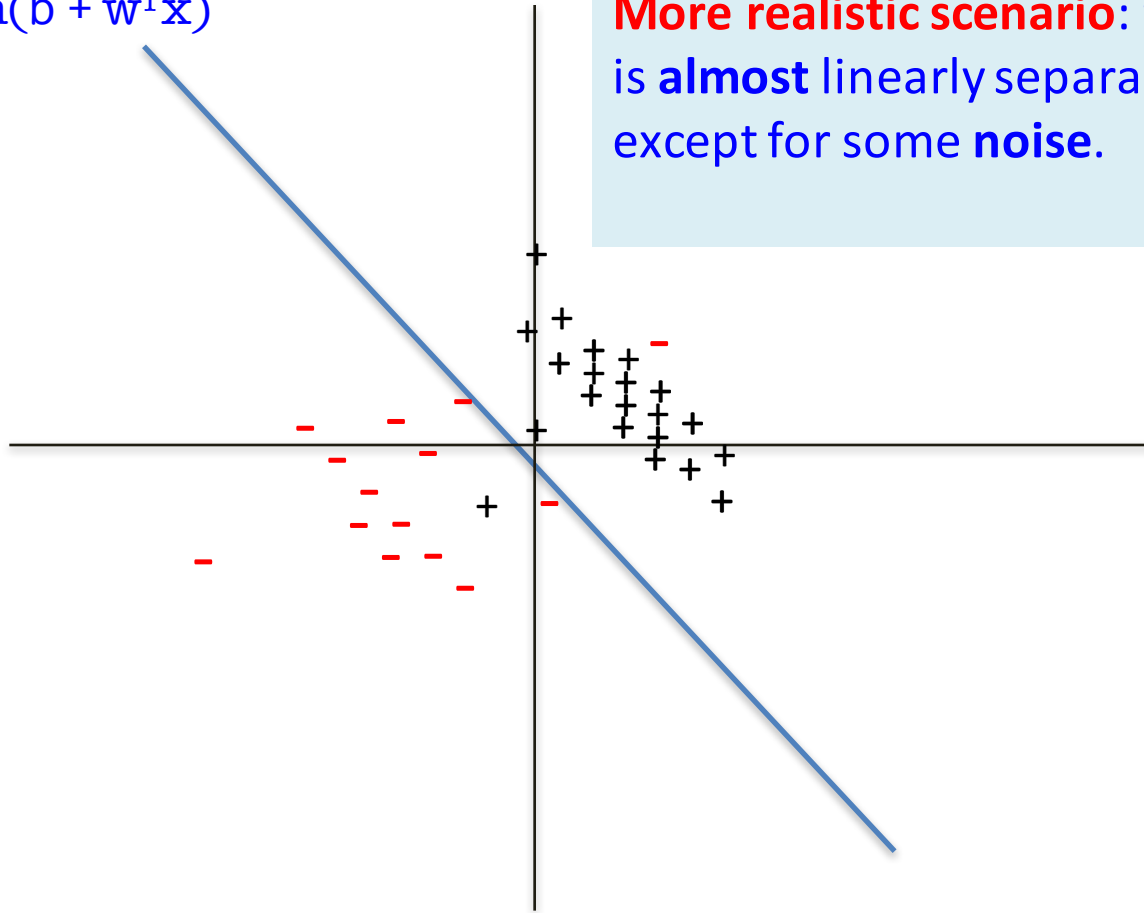


Represent each point in 2D as  $(x, x^2)$



# Expressivity

$\text{sign}(b + w^T x)$



**More realistic scenario:** the data is **almost** linearly separable, except for some **noise**.

# Features

- So far we have discussed BoW representation
  - *In fact, you can use a very rich representation*
- **Broader definition**
  - Functions mapping attributes of the input to a Boolean/categorical/numeric value

$$\phi_1(x) = \begin{cases} 1 & x_1 \text{ is capitalized} \\ 0 & \text{otherwise} \end{cases}$$

$$\phi_k(x) = \begin{cases} 1 & x \text{ contains "good" more than twice} \\ 0 & \text{otherwise} \end{cases}$$

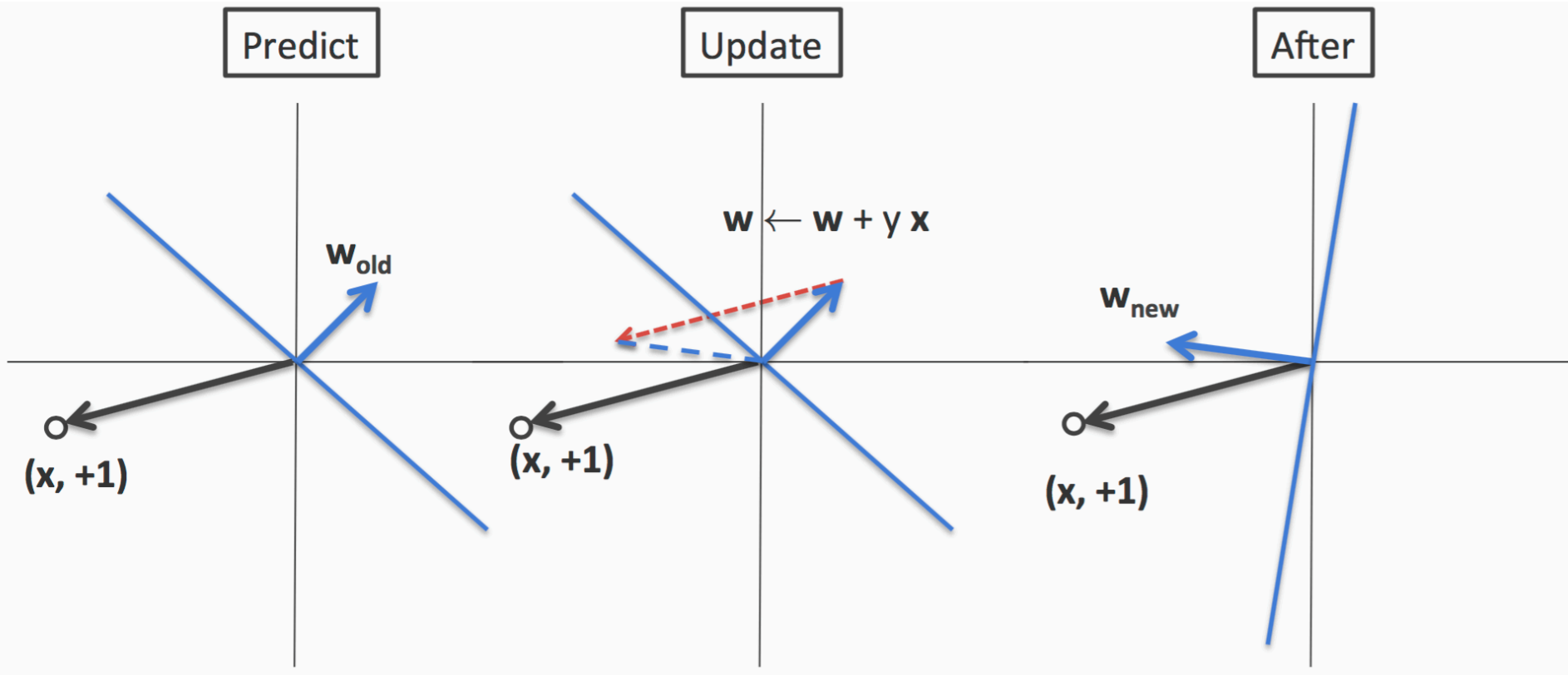
- **Question:** *assume that you have a lexicon, containing positive and negative sentiment words. How can you use it to improve over BoW?*

# Perceptron

- One of the earliest learning algorithms
  - Introduced by Rosenblatt 1958 to model neural learning
- **Goal:** directly search for a separating hyperplane
  - If one exists, perceptron will find it
  - If not, ...
- **Online** algorithm
  - Considers one example at a time (NB – looks at entire data)
- **Error driven** algorithm
  - Updates the weights only when a mistake is made



# Perceptron Intuition



# Perceptron

- We learn  $f: X \rightarrow \{-1, +1\}$  represented as  $f = \text{sgn}\{w \bullet x\}$
- Where  $X = \{0, 1\}^n$  or  $X = \mathbb{R}^n$  and  $w \in \mathbb{R}^n$
- Given Labeled examples:  $\{(x_1, y_1), (x_2, y_2), \dots, (x_m, y_m)\}$

1. Initialize  $w = 0 \in \mathbb{R}^n$

2. Cycle through all examples

- a. **Predict** the label of instance  $x$  to be  $y' = \text{sgn}\{w \bullet x\}$

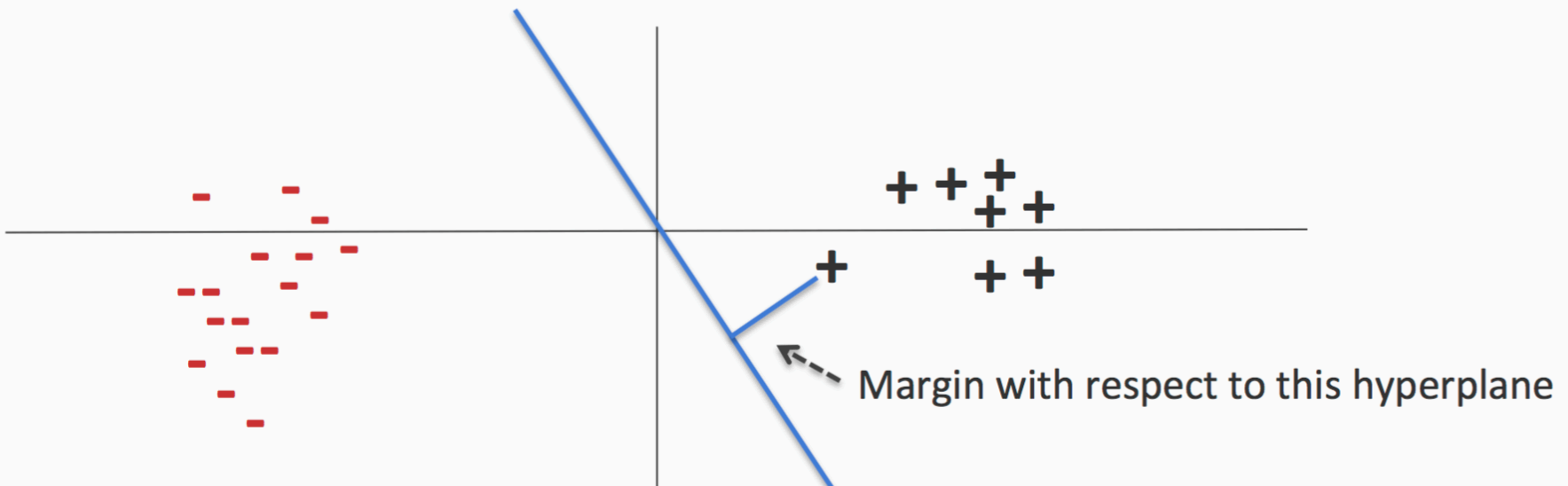
- b. If  $y' \neq y$ , **update** the weight vector:

$$w = w + r y x \quad (r - \text{a constant, learning rate})$$

Otherwise, if  $y' = y$ , leave weights unchanged.

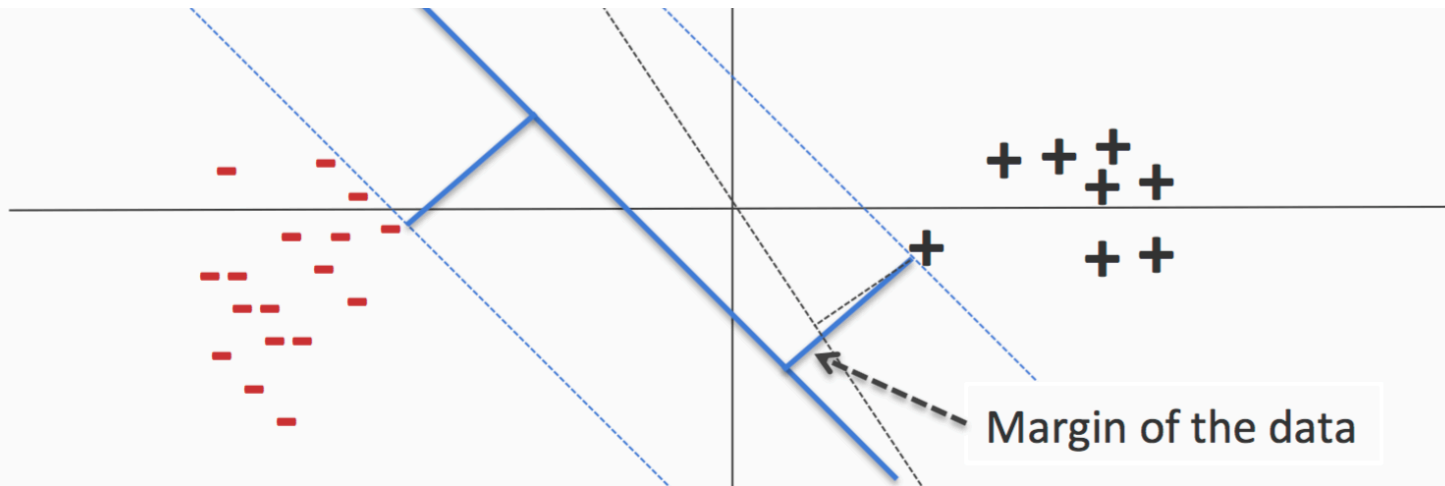
# Margin

- The **margin of a hyperplane** for a dataset is the distance between the hyperplane and the data point nearest to it.



# Margin

- The **margin of a hyperplane** for a dataset is the distance between the hyperplane and the data point nearest to it.
- The **margin of a data set** ( $\gamma$ ) is the maximum margin possible for that dataset using any weight vector.

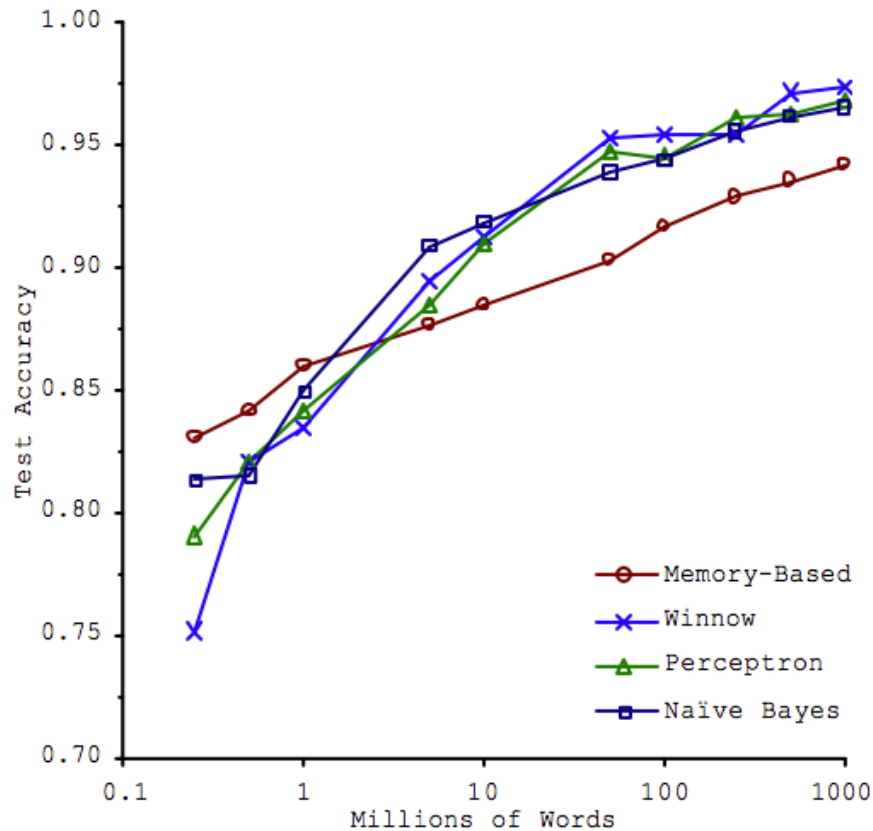


# Mistake Bound for Perceptron

- Let  $D = \{(x_i, y_i)\}$  be a labeled dataset that is separable
- Let  $\|x_i\| < R$  for all examples.
- Let  $\gamma$  be the margin of the dataset  $D$ .
- Then, the perceptron algorithm will make at most  $R^2 / \gamma^2$  mistakes on the data.

# Practical Example

Task: **context sensitive spelling**  
{principle,  
principal},{weather,whet  
her}.



Source: *Scaling to very very large corpora for natural language disambiguation*  
Michele Banko, Eric Brill. Microsoft Research, Redmond, WA. 2001.

# Deceptive Reviews

Which of these two hotel reviews is deceptive opinion spam?

My husband and I stayed at the James Chicago Hotel for our anniversary. This place is fantastic! We knew as soon as we arrived we made the right choice! The rooms are BEAUTIFUL and the staff very attentive and wonderful!! The area of the hotel is great, since I love to shop I couldn't ask for more!! We will definatly be back to Chicago and we will for sure be back to the James Chicago.

I have stayed at many hotels traveling for both business and pleasure and I can honestly say that The James is tops. The service at the hotel is first class. The rooms are modern and very comfortable. The location is perfect within walking distance to all of the great sights and restaurants. Highly recommend to both business travellers and couples.

*What should your learning algorithm look at?*

# Deception Classification

| Approach                             | Features                                 | Accuracy     | TRUTHFUL    |             |             | DECEPTIVE   |             |             |
|--------------------------------------|--|--------------|-------------|-------------|-------------|-------------|-------------|-------------|
|                                      |  |              | P           | R           | F           | P           | R           | F           |
| GENRE IDENTIFICATION                 | POS <sub>SVM</sub>                       | 73.0%        | 75.3        | 68.5        | 71.7        | 71.1        | 77.5        | 74.2        |
| PSYCHOLINGUISTIC DECEPTION DETECTION | LIWC <sub>SVM</sub>                      | 76.8%        | 77.2        | 76.0        | 76.6        | 76.4        | 77.5        | 76.9        |
| TEXT CATEGORIZATION                  | UNIGRAMS <sub>SVM</sub>                  | 88.4%        | 89.9        | 86.5        | 88.2        | 87.0        | 90.3        | 88.6        |
|                                      | BIGRAMS <sub>SVM</sub> <sup>+</sup>      | 89.6%        | 90.1        | 89.0        | 89.6        | 89.1        | 90.3        | 89.7        |
|                                      | LIWC+BIGRAMS <sub>SVM</sub> <sup>+</sup> | <b>89.8%</b> | 89.8        | <b>89.8</b> | <b>89.8</b> | <b>89.8</b> | 89.8        | <b>89.8</b> |
|                                      | TRIGRAMS <sub>SVM</sub> <sup>+</sup>     | 89.0%        | 89.0        | 89.0        | 89.0        | 89.0        | 89.0        | 89.0        |
|                                      | UNIGRAMS <sub>NB</sub>                   | 88.4%        | <b>92.5</b> | 83.5        | 87.8        | 85.0        | <b>93.3</b> | 88.9        |
|                                      | BIGRAMS <sub>NB</sub> <sup>+</sup>       | 88.9%        | 89.8        | 87.8        | 88.7        | 88.0        | 90.0        | 89.0        |
|                                      | TRIGRAMS <sub>NB</sub> <sup>+</sup>      | 87.6%        | 87.7        | 87.5        | 87.6        | 87.5        | 87.8        | 87.6        |
| HUMAN / META                         | JUDGE 1                                  | <b>61.9%</b> | 57.9        | 87.5        | <b>69.7</b> | 74.4        | 36.3        | 48.7        |
|                                      | JUDGE 2                                  | 56.9%        | 53.9        | <b>95.0</b> | 68.8        | <b>78.9</b> | 18.8        | 30.3        |
|                                      | SKEPTIC                                  | 60.6%        | <b>60.8</b> | 60.0        | 60.4        | 60.5        | <b>61.3</b> | <b>60.9</b> |



# Summary

- Classification is a basic tool for NLP
  - *E.g., What is the topic of a document?*
- **Classifier: *mapping from input to label***
  - *Label: Binary or Categorical*
- We saw two simple learning algorithms for finding the parameters of linear classification functions
  - Naïve Bayes and Perceptron
- **Next:**
  - More sophisticated algorithms
  - Applications (or – *how to get it to work!*)

*Questions?*