

# Lecture 6

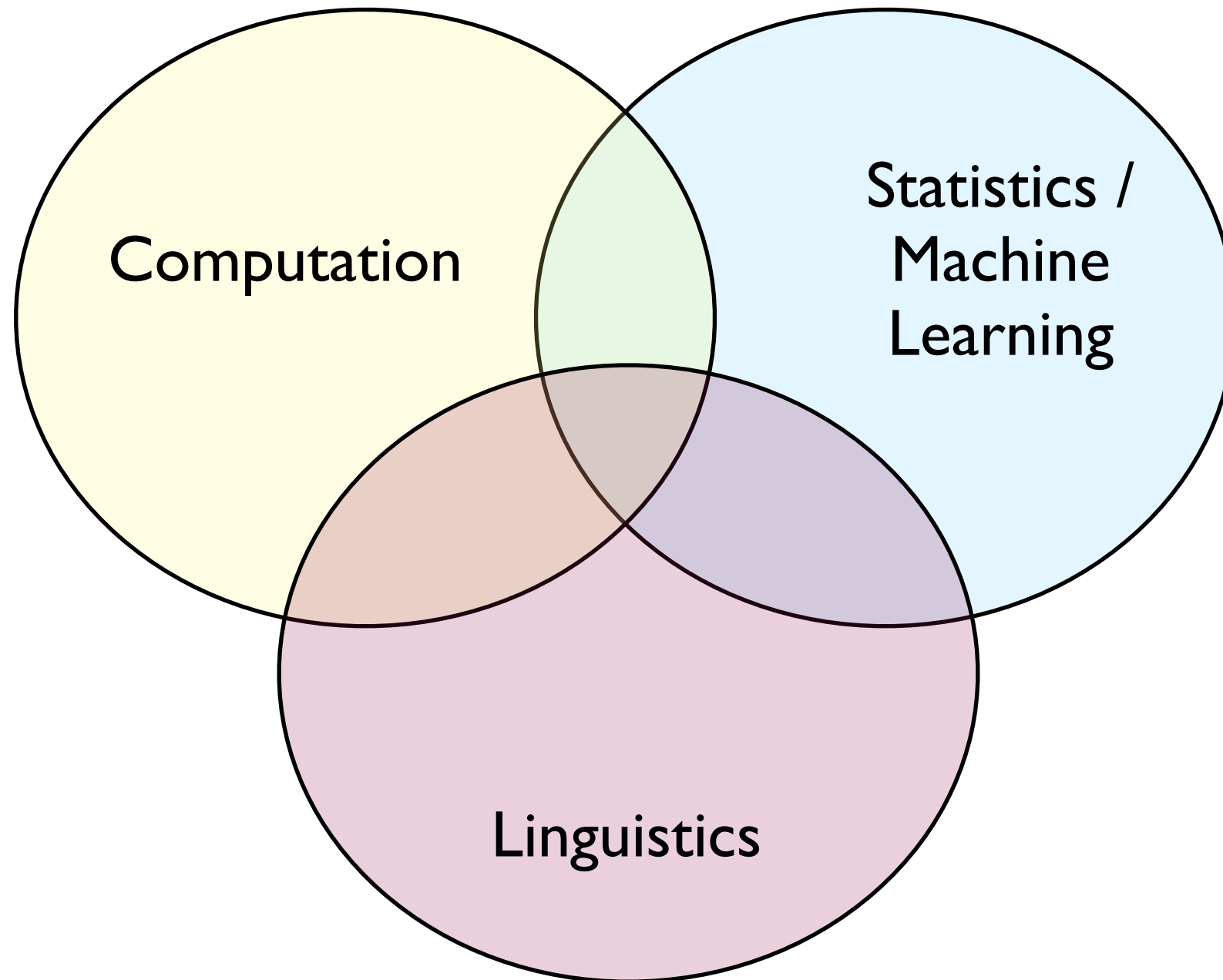
## Classification: Naive Bayes

Intro to NLP, CS585, Fall 2014

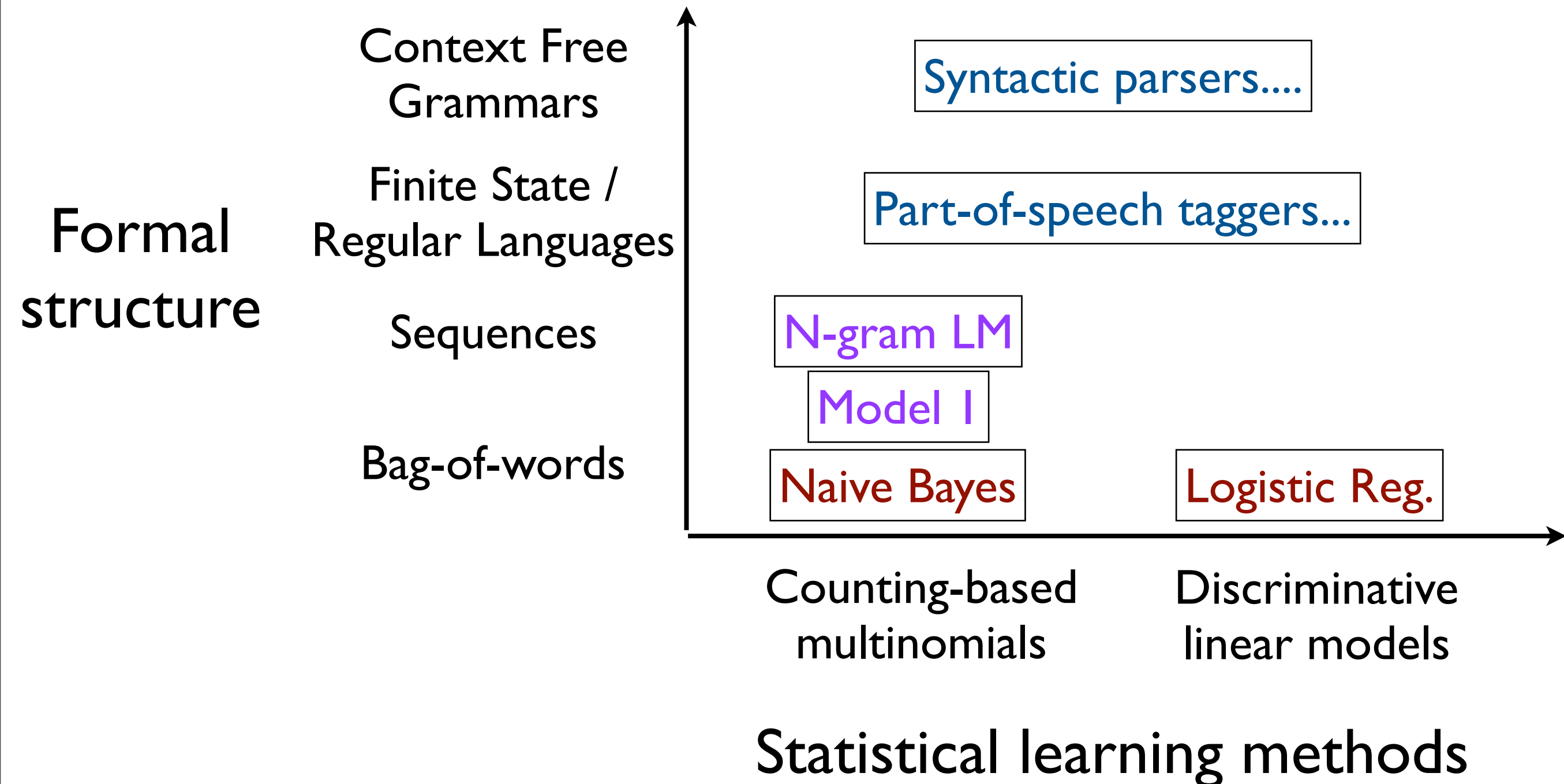
<http://people.cs.umass.edu/~brenocon/inlp2014/>

Brendan O'Connor (<http://brenocon.com>)

# This course includes



# Computation/Statistics in NLP (in this course)

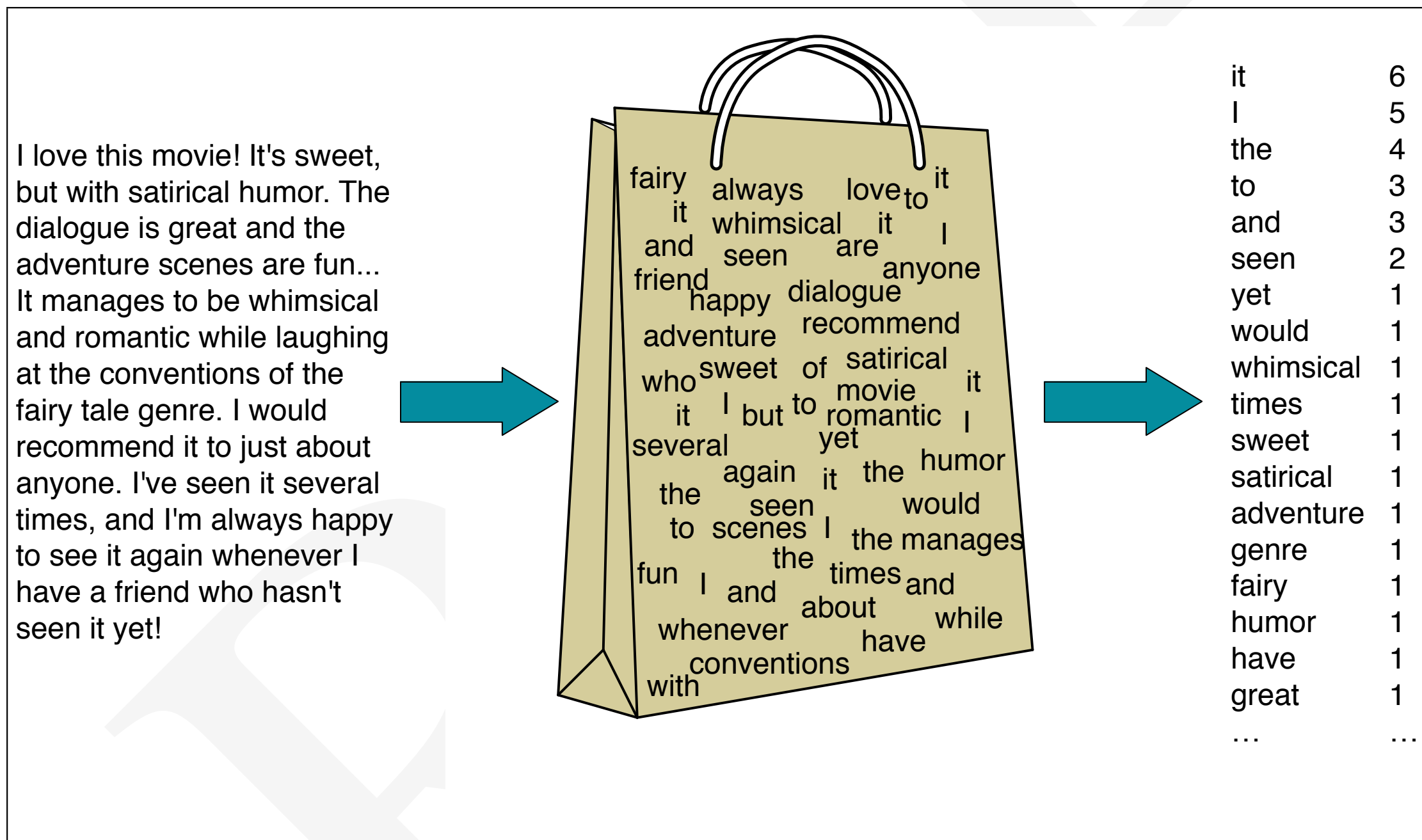


# Classification problems

- Given text ***d***, want to predict label ***y***
  - Is this restaurant review positive or negative?
  - Is this email spam or not?
  - Which author wrote this text?
  - (Is this word a noun or verb?)
- ***d***: documents, sentences, etc.
- ***y***: discrete/categorical variable

Goal: from training set of  $(d, y)$  pairs, *learn*  
a probabilistic classifier  $f(d) = P(y|d)$   
("supervised learning")

# Features for model: Bag-of-words



**Figure 6.1** Intuition of the multinomial naive Bayes classifier applied to a movie review. The position of the words is ignored (the *bag of words* assumption) and we make use of the frequency of each word.

# Generative vs. Discriminative approaches

Goal: from training set of  $(d,y)$  pairs, *learn* a probabilistic “classifier”  $f(d) = P(y|d)$

---

Generative model: use the “noisy channel” idea.

$$P(y \mid d) \propto P(y) P(d \mid y; \theta)$$

*Learning:*  $\max_{\theta} \prod_{i \in \text{train}} P(d_i \mid y_i; \theta)$  (where it's just counting)

**Naive Bayes**

---

Discriminative model: directly learn this function

$$P(y \mid d) = f(d; \theta)$$

*Learning:*  $\max_{\theta} \prod_{i \in \text{train}} P(y_i \mid d_i)$  (where it's harder than counting)

**Logistic Regression**

# Multinomial Naive Bayes: Unigram LM

Tokens in doc



$$P(y \mid w_1..w_T) \propto P(y) P(w_1..w_T \mid y)$$

conditional  
independence  
assumption



$$\prod_t P(w_t \mid y)$$

- Generative story:
- Choose doc category **y**
- For each token position in doc:
  - Draw **w<sub>t</sub>**

Parameters:  $P(w \mid y)$  for each document category **y** and wordtype **w**  
 $P(y)$  prior distribution over document categories **y**

Learning: with pseudocount smoothing,

$$P(w \mid y, \alpha) = \frac{\#(w \text{ occurrences in docs with label } y) + \alpha}{\#(\text{tokens total across docs with label } y) + V\alpha}$$

# Multinomial Naive Bayes: Unigram LM

## Prediction

Infer most likely class for new document

$$\arg \max_k P(y = k) \prod_t P(w_t \mid y = k)$$

Infer posterior probabilities for new document

$$P(y = k \mid w_1..w_T) = \frac{P(y = k) \prod_t P(w_t \mid y = k)}{\sum_{k'} P(y = k') \prod_t P(w_t \mid y = k')}$$



# Example

## Learning

### Estimate prior

$$P(-) = \frac{3}{5} \quad P(+) = \frac{2}{5}$$

### Estimate word likelihoods with pseudocount=1

## Prediction/Inference

$$P(S|-)P(-) = \frac{3}{5} \times \frac{2 \times 1 \times 2 \times 1}{34^4} = 1.8 \times 10^{-6}$$

$$P(S|+)P(+) = \frac{2}{5} \times \frac{1 \times 1 \times 1 \times 1}{29^4} = 5.7 \times 10^{-7}$$

	Cat	Documents
Training	-	just plain boring
	-	entirely predictable and lacks energy
	-	no surprises and very few laughs
	+	very powerful
	+	the most fun film of the summer
Test	?	predictable with no originality

$$P(\text{"predictable"}|-) = \frac{1+1}{14+20}$$

$$P(\text{"predictable"}|+) = \frac{0+1}{9+20}$$

$$P(\text{"with"}|-) = \frac{0+1}{14+20}$$

$$P(\text{"with"}|+) = \frac{0+1}{9+20}$$

$$P(\text{"no"}|-) = \frac{1+1}{14+20}$$

$$P(\text{"no"}|+) = \frac{0+1}{9+20}$$

$$P(\text{"originality"}|-) = \frac{0+1}{14+20}$$

$$P(\text{"originality"}|+) = \frac{0+1}{9+20}$$

# NB as a Linear Model

Consider: ratio of posterior probs

$$\frac{P(+ \mid w_1..w_T)}{P(- \mid w_1..w_T)} \quad \begin{array}{l} > 1 \text{ then } + \text{ more likely} \\ < 1 \text{ then } - \text{ more likely} \end{array}$$

Odds form of Bayes Rule:

$$\begin{aligned} & \begin{array}{cc} \text{prior ratio} & \text{likelihood ratio} \end{array} \\ = & \frac{P(+)}{P(-)} \frac{P(w_1..w_T \mid +)}{P(w_1..w_T \mid -)} \frac{\cancel{1/P(w_1..w_T)}}{\cancel{1/P(w_1..w_T)}} \\ = & \frac{P(+)}{P(-)} \frac{\prod_t P(w_t \mid +)}{\prod_t P(w_t \mid -)} \\ = & \frac{3}{5} \frac{2/34 \times 1/34 \times 2/34 \times 1/34}{1/29 \times 1/29 \times 1/29 \times 1/29} \end{aligned}$$

# NB as a Linear Model

$$\frac{P(+ | w_1..w_T)}{P(- | w_1..w_T)} = \frac{P(+)}{P(-)} \frac{\prod_t P(w_t|+)}{\prod_t P(w_t|-)}$$

>1 then + more likely

<1 then - more likely

$$= \frac{P(+)}{P(-)} \prod_t \frac{P(w_t|+)}{P(w_t|-)}$$

$$\log \frac{P(+ | w_1..w_T)}{P(- | w_1..w_T)} = \log \frac{P(+)}{P(-)} + \sum_t \log \frac{P(w_t|+)}{P(w_t|-)}$$

>0 then + more likely

<0 then - more likely

$$= \log \frac{P(+)}{P(-)} + \sum_w^V n_w \log \frac{P(w|+)}{P(w|-)}$$

$$= \log \frac{3}{5} + \log \frac{2/34}{1/29} + \log \frac{1/34}{1/29} + \log \frac{2/34}{1/29} + \log \frac{1/34}{1/29}$$

# NB as a Linear Model

$$\log \frac{P(+ | w_1..w_T)}{P(- | w_1..w_T)} = \log \frac{P(+)}{P(-)} + \sum_w^V n_w \log \frac{P(w|+)}{P(w|-)}$$

>0 then + more likely

<0 then - more likely

$$= \frac{\beta_0}{\beta_0} + \frac{(\beta_{1:V})^T \mathbf{n}}{(\beta_{1:V})^T \mathbf{n}}$$

$$= \beta^T \mathbf{x}$$

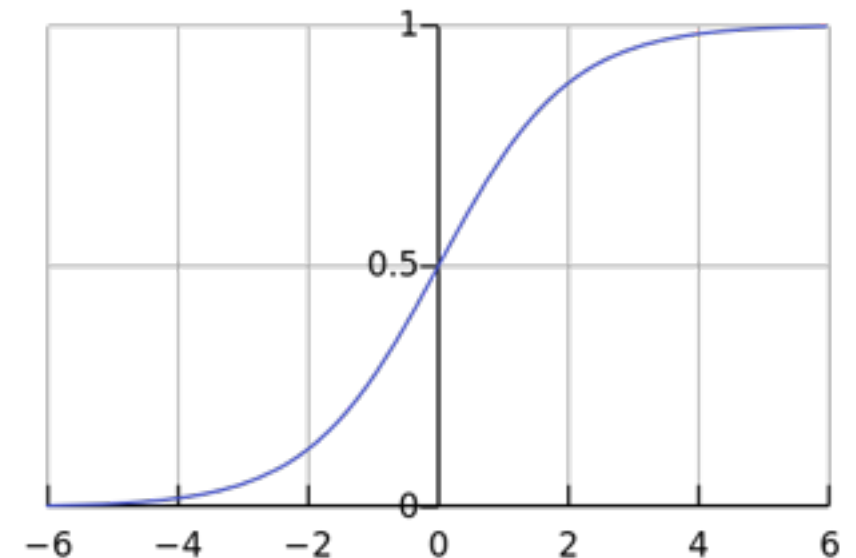
Where

$\mathbf{x} = (1, \text{count "happy", count "sad", ....})$  *Feature vector*

$$P(+ | w_1..w_T) = \frac{\exp(\beta^T \mathbf{x})}{1 + \exp(\beta^T \mathbf{x})}$$

*Logistic sigmoid function*

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



# Logistic regression

$$P(+ \mid w_1..w_T) = \frac{\exp(\beta^T \mathbf{x})}{1 + \exp(\beta^T \mathbf{x})}$$

- NB (decision between unigram LMs) prescribes one particular formula for the beta weights.
- Can we just fit the beta weights to maximize likelihood of the training data?