## Master in Artificial Intelligence

Statistical Models for NI P

Maximum Likelihood Estimation (MLE)

Maximum Entropy Modeling

Log-Linear Models

# IIP

Advanced Human Language Technologies Statistical Models of Language



UNIVERSITAT POLITÈCNICA DE CATALUNYA BARCELONATECH

Facultat d'Informàtica de Barcelona

FIR

## Outline

#### Statistical Models for NLP

Maximum Likelihood Estimation (MLE)

Maximum Entropy Modeling

Log-Linear Models

## 1 Statistical Models for NLP

- Why modeling
- Prediction & Similarity Models
- Maximum Likelihood Estimation (MLE)
  Working example
  - Smoothing & Estimator Combination
- 3 Maximum Entropy Modeling
  - Overview
  - Building ME Models



## Outline

Statistical Models for NLP

Why modeling

Maximum Likelihood Estimation (MLE)

Maximum Entropy Modeling

Log-Linear Models

## 1 Statistical Models for NLP

- Why modeling
- Prediction & Similarity Models
- g
- Working example
  - Smoothing & Estimator Combination
- Maximum Entropy Modeling
  - Overview
  - Building ME Models





Statistical Models for NLP Why modeling Maximum Likelihood

Estimation (MLE)

Maximum Entropy Modeling







Statistical

Maximum Entropy Modeling





Maximum Entropy Modeling



## Outline



Prediction & Similarity Models

Maximum Likelihood Estimation (MLE)

Maximum Entropy Modeling

Log-Linear Models

### Statistical Models for NLP

Why modeling

Prediction & Similarity Models

Working example

- Smoothing & Estimator Combination

- Overview
- Building ME Models



## Prediction Models & Similarity Models

Statistical Models for NLP

Prediction & Similarity Models

Maximum Likelihood Estimation (MLE)

Maximum Entropy Modeling

- Prediction Models: Oriented to *predict* probabilities of future events, knowing past and present.
- **Similarity Models**: Oriented to compute *similarities* between objects (may be used to predict, EBL).

## Similarity Models

 Objects represented as feature-vectors, feature-sets, distribution-vectors, ...

Statistical Models for NLP

Prediction & Similarity Models

Maximum Likelihood Estimation (MLE)

Maximum Entropy Modeling

- Used to group objects (clustering, data analysis, pattern discovery, ...)
- If classified objects are available, similarity may be used as a prediction (example-based ML techniques).
- Example: Document representation
  - Documents are represented as vectors in a high dimensional R<sup>n</sup> space.
  - Dimensions are word forms, lemmas, NEs, n-grams, ...
  - Values may be either binary or real-valued (count, frequency, ...)
  - Vector-space algebra and metrics can be used

Statistical Models for NLP

Prediction & Similarity Models

Maximum Likelihood Estimation (MLE)

Maximum Entropy Modeling

Log-Linear Models

## **Prediction Models**

Estimation: Using data to infer information about distributions

- Parametric / non-parametric estimation
- Finding good estimators: MLE, MEE, ...
- Explicit / implicit models
- Classification: Predictions based on past behaviour
  - Predict most likely target given classification features (implies independence assumptions!)
  - Granularity of equivalence classes (bins): discrimination power *vs.* statistical reliability
- In general, ML models estimate (i.e. *learn*) conditional probability distributions P(target|features)
- Many NLP tasks require a posterior search step to find the best combination of predictions.

## Prediction Models

**NLP** Applications





#### Statistical Models for NLP

Prediction & Similarity Models

Maximum Likelihood Estimation (MLE)

Maximum Entropy Modeling

Log-Linear Models

Appl.	Input	Output	p(i)	p(o   i)
MT	L word	M word	p(L)	Translation
	sequence	sequence		model
OCR	Actual text	Text with	prob. of	model of
		mistakes	language text	OCR errors
PoS	PoS tags	word	prob. of PoS	$p(w \mid t)$
tagging	sequence	sequence	sequence	
Speech	word	speech	prob. of word	acoustic
recog.	sequence	signal	sequence	model

Given o, we want to find the most likely i

$$\underset{i}{\operatorname{argmax}} P(\mathbf{i} \mid \mathbf{o}) = \underset{i}{\operatorname{argmax}} P(\mathbf{o}, \mathbf{i}) = \underset{i}{\operatorname{argmax}} P(\mathbf{i})P(\mathbf{o} \mid \mathbf{i})$$

## Finding good estimators: MLE

#### Maximum Likelihood Estimation (MLE)

- Choose the alternative that maximizes the probability of the observed outcome.
- $\bar{\mu}_n$  is a MLE for E(X)
- $s_n^2$  is a MLE for  $\sigma^2$
- Zipf's Laws. Data sparseness. Smoothing tecnhiques.

P(a, b)	dans	en	à	sur	au-cours-de	pendant	selon	
in	0.04	0.10	0.15	0	0.08	0.03	0	0.40
on	0.06	0.25	0.10	0.15	0	0	0.04	0.60
total	0.10	0.35	0.25	0.15	0.08	0.03	0.04	1.0

Statistical Models for NLP

Prediction & Similarity Models

Maximum Likelihood Estimation (MLE)

Maximum Entropy Modeling

## Finding good estimators: MEE

#### Maximum Entropy Estimation (MEE)

 Choose the alternative that maximizes the entropy of the obtained distribution, maintaining the observed probabilities.

```
Observations:
```

$$p(en \lor a) = 0.6$$

P(a, b)	dans	en	à	sur	au-cours-de	pendant	selon	
in	0.04	0.15	0.15	0.04	0.04	0.04	0.04	
on	0.04	0.15	0.15	0.04	0.04	0.04	0.04	
total			_					1.0
		õ	.6					

Statistical Models for NLP

Prediction & Similarity Models

Maximum Likelihood Estimation (MLE)

Maximum Entropy Modeling

## Finding good estimators: MEE

#### Maximum Entropy Estimation (MEE)

 Choose the alternative that maximizes the entropy of the obtained distribution, maintaining the observed probabilities.

```
Observations:
```

$$p(en \lor a) = 0.6;$$
  $p((en \lor a) \land in) = 0.4$ 

P(a, b)	dans	en	à	sur	au-cours-de	pendant	selon	
in	0.04	0.20	0.20	0.04	0.04	0.04	0.04	
on	0.04	0.10	0.10	0.04	0.04	0.04	0.04	
total			_					1.0
		0	.6					

Statistical Models for NLP

Prediction & Similarity Models

Maximum Likelihood Estimation (MLE)

Maximum Entropy Modeling

## Finding good estimators: MEE

#### Maximum Entropy Estimation (MEE)

 Choose the alternative that maximizes the entropy of the obtained distribution, maintaining the observed probabilities.

```
Observations:
```

$$p(\mathsf{en} \lor \grave{\mathsf{a}}) = \mathsf{0.6}; \qquad p((\mathsf{en} \lor \grave{\mathsf{a}}) \land \mathsf{in}) = \mathsf{0.4}; \qquad p(\mathsf{in}) = \mathsf{0.5}$$

P(a, b)	dans	en	à	sur	au-cours-de	pendant	selon	
in	0.02	0.20	0.20	0.02	0.02	0.02	0.02	0.5
on	0.06	0.10	0.10	0.06	0.06	0.06	0.06	
total								1.0
		0	.6					

Statistical Models for NLP

Prediction & Similarity Models

Maximum Likelihood Estimation (MLE)

Maximum Entropy Modeling

## Outline

Statistical Models for NLP

Maximum Likelihood Estimation (MLE)

Maximum Entropy Modeling

Log-Linear Models

#### Statistical Models for NLP

- Why modeling
- Prediction & Similarity Models

## 2 Maximum Likelihood Estimation (MLE)

- Working example
- Smoothing & Estimator Combination
- 8 Maximum Entropy Modeling
  - Overview
  - Building ME Models



## Outline

Statistical Models for NLP

Maximum Likelihood Estimation (MLE)

Working example

Maximum Entropy Modeling

Log-Linear Models

# 2 Maximum Likelihood Estimation (MLE)

Prediction & Similarity Models

#### Working example

Why modeling

Smoothing & Estimator Combination

Maximum Entropy Modeling

- Overview
- Building ME Models



## Working Example: N-gram models

- Predict the next element in a sequence (e.g. next character, next word, next PoS, next stock value, ... ), given the *history* of previous elements:  $P(w_n | w_1 \dots w_{n-1})$
- Markov assumption: Only *local* context (of size n − 1) is taken into account. P(w<sub>i</sub> | w<sub>i-n+1</sub>...w<sub>i-1</sub>)
- bigrams, trigrams, four-grams (n = 2, 3, 4).
   Sue swallowed the large green <?>
- Parameter estimation (number of equivalence classes)
- Parameter reduction: stemming, semantic classes, PoS, ...

Model	Parameters
bigram	$20,000^2 = 4  imes 10^8$
trigram	$20,000^3 = 8  imes 10^{12}$
four-gram	$20,000^4 = 1.6  imes 10^{17}$

Language model sizes for a 20,000 words vocabulary

Statistical Models for NLP

Maximum Likelihood Estimation (MLE)

Working example

Maximum Entropy Modeling

## N-gram model estimation

Estimate the probability of the target feature based on observed data. The prediction task can be reduced to having good estimations of the n-gram distribution:

$$\mathsf{P}(w_n \mid w_1 \dots w_{n-1}) = \frac{\mathsf{P}(w_1 \dots w_n)}{\mathsf{P}(w_1 \dots w_{n-1})}$$

## • MLE (Maximum Likelihood Estimation) $P_{MLE}(w_1 \dots w_n) = \frac{C(w_1 \dots w_n)}{N}$ $P_{MLE}(w_n \mid w_1 \dots w_{n-1}) = \frac{C(w_1 \dots w_n)}{C(w_1 \dots w_{n-1})}$

- No probability mass for unseen events
- Data sparseness, Zipf's Law
- Unsuitable for NLP (widely used, though)

Statistical Models for NLP

Maximum Likelihood Estimation (MLE)

Working example

Maximum Entropy Modeling

## Brief Parenthesis: Zipf's Laws

#### Zipf's Laws (1929)

- $\blacksquare$  Word frequency is inversely proportional to its rank (speaker/hearer minimum effort)  $f\sim 1/r$
- $\blacksquare$  Number of senses is proportional to frequency root  $m\sim \sqrt{f}$
- $\blacksquare$  Frequency of intervals between repetitions is inversely proportional to the length of the interval F  $\sim 1/I$
- Frequency based approaches are hard, since most words are rare
  - Most common 5% words account for about 50% of a text
  - 90% least common words account for less than 10% of the text
  - Almost half of the words in a text occurr only once

Statistical Models for NLP

Maximum Likelihood Estimation (MLE)

Working example

Maximum Entropy Modeling

## Outline

Statistical Models for NLP

Maximum Likelihood Estimation (MLE)

Smoothing & Estimator Combination

Maximum Entropy Modeling

Log-Linear Models

## 1 Statistical Models for NLP

- Why modeling
- Prediction & Similarity Models

# Maximum Likelihood Estimation (MLE) Working example

Smoothing & Estimator Combination

Maximum Entropy Modeling

- Overview
- Building ME Models



## Notation

Statistical Models for NLP

Maximum Likelihood Estimation (MLE)

Smoothing & Estimator Combination

Maximum Entropy Modeling

- $C(w_1 \dots w_n)$ : Observed occurrence count for n-gram  $w_1 \dots w_n$ .
- N: Number of observed n-gram occurrences

$$\mathsf{N} = \sum_{w_1 \dots w_n} \mathsf{C}(w_1 \dots w_n)$$

- N<sub>k</sub>: Number of classes (n-grams) observed k times.
- B: Number of equivalence classes or bins (number of potentially observable n-grams).

## Smoothing 1 - Adding Counts

Statistical Models for NLP

Maximum Likelihood Estimation (MLE)

Smoothing & Estimator Combination

Maximum Entropy Modeling

Log-Linear Models

## • Laplace's Law (adding one) $P_{LAP}(w_1 \dots w_n) = \frac{C(w_1 \dots w_n) + 1}{N + B}$

 For large values of B too much probability mass is assigned to unseen events

#### Lidstone's Law

$$\mathsf{P}_{\mathrm{LID}}(w_1 \dots w_n) = \frac{\mathsf{C}(w_1 \dots w_n) + \lambda}{\mathsf{N} + \mathsf{B}\lambda}$$

• Usually  $\lambda = 0.5$ , *Expected Likelihood Estimation*.

Equivalent to linear interpolation between MLE and uniform prior, with  $\mu = N/(N + B\lambda)$ ,

$$P_{LID}(w_1 \dots w_n) = \mu \frac{C(w_1 \dots w_n)}{N} + (1-\mu) \frac{1}{B}$$

## Smoothing 2 - Discounting Counts

#### Absolute Discounting

Statistical Models for NLP

Maximum Likelihood Estimation (MLE)

Smoothing & Estimator Combination

Maximum Entropy Modeling

Log-Linear Models

$$P_{ABS}(w_1 \dots w_n) = \begin{cases} \frac{C(w_1 \dots w_n) - \delta}{N} & \text{if } C(w_1 \dots w_n) > 0\\ \frac{(B - N_0)\delta/N_0}{N} & \text{otherwise} \end{cases}$$

Linear Discounting

$$P_{LIN}(w_1 \dots w_n) = \left\{ \begin{array}{ll} (1-\alpha) \frac{C(w_1 \dots w_n)}{N} & \text{if } C(w_1 \dots w_n) > 0 \\ \\ \alpha/N_0 & \text{otherwise} \end{array} \right.$$

Statistical Models for NLP

Maximum Likelihood Estimation (MLE)

Smoothing & Estimator Combination

Maximum Entropy Modeling

Log-Linear Models

## Combining Estimators

#### Simple Linear Interpolation

$$P_{LI}(w_{n} | w_{n-2}, w_{n-1}) = \lambda_{1} P_{1}(w_{n}) + \lambda_{2} P_{2}(w_{n} | w_{n-1}) + \lambda_{3} P_{3}(w_{n} | w_{n-2}, w_{n-1})$$

Backing-off

$$P_{BO}(w_i \mid h) = \begin{cases} (1 - \alpha_h) \frac{C(h, w_i)}{C(h)} & \text{if } C(h, w_i) > k \\ \delta_{h'} P_{BO}(w_i \mid h') & \text{otherwise} \end{cases}$$

 $\begin{array}{ll} \left( \text{where } h = w_{i-n+1} \dots w_{i-1}, & h' = w_{i-n+2} \dots w_{i-1} \right) \\ \text{Different options to determine } \alpha_h \text{ and } \delta_{h'} \ \left( \text{e.g. } \alpha_h = \delta_{h'} & \forall h \right) \end{array}$ 

## Outline

Statistical Models for NLP

Maximum Likelihood Estimation (MLE)

Maximum Entropy Modeling

Log-Linear Models

#### Statistical Models for NLP

- Why modeling
- Prediction & Similarity Models
- Maximum Likelihood Estimation (MLE Working example
- Smoothing & Estimator Combination

### 3 Maximum Entropy Modeling

- Overview
- Building ME Models



## Outline

Statistical Models for NLP

Maximum Likelihood Estimation (MLE)

Maximum Entropy Modeling Overview

Log-Linear Models

#### Statistical Models for NLP

- Why modeling
- Prediction & Similarity Models
- Maximum Likelihood Estimation (MLE Working example
- Smoothing & Estimator Combination

Maximum Entropy Modeling
 Overview

Building ME Models



## **MEM** Overview

- Maximum Entropy: alternative estimation technique.
- Able to deal with different kinds of evidence
- ME principle:
  - Do not assume anything about non-observed events.
  - Find the most uniform (maximum entropy, less informed) probability distribution that matches the observations.
- Example:

p(x, y)	0	1		p(x, y)	0	1		p(x, y)	0	1	
а	?	?		а	0.5	0.1		а	0.3	0.2	
b	?	?		b	0.1	0.3		b	0.3	0.2	
total	0.6		1.0	total	0.6		1.0	total	0.6		1.0
Obse	 ervat	ion	5	One po	ssible	e p(x	.,y)	Max.Er	htrop	y p()	(, y)

Statistical Models for NLP

Maximum Likelihood Estimation (MLE)

Maximum Entropy Modeling Overview

## ME Modeling

Observed facts are constraints for the desired model p.
 Constraints take the form of feature functions:

$$f_i: \varepsilon \to \{0, 1\}$$

 The desired model p must satisfy the constraints: The expectation predicted by model p for any feature f<sub>i</sub> must match the observed expectation for f<sub>i</sub> i.e.:

$$\begin{array}{rcl} \mathsf{E}_{p}(\mathsf{f}_{\mathfrak{i}}) &=& \mathsf{E}_{\widetilde{p}}(\mathsf{f}_{\mathfrak{i}}) & \forall \mathfrak{i} \\ \displaystyle \sum_{x \in \varepsilon} p(x) \mathsf{f}_{\mathfrak{i}}(x) &=& \displaystyle \sum_{x \in \varepsilon} \widetilde{p}(x) \mathsf{f}_{\mathfrak{i}}(x) & \forall \mathfrak{i} \end{array}$$

Statistical Models for NLP

Maximum Likelihood Estimation (MLE)

Maximum Entropy Modeling Overview

## Example

Example:

Statistical Models for NLP

Maximum Likelihood Estimation (MLE)

Maximum Entropy Modeling Overview

Log-Linear Models

	p(x, y)	0	1	
$c = [a, b] \times [0, 1]$	a	?	?	
$\mathcal{E} = \{\mathbf{u}, \mathbf{v}\} \times \{\mathbf{v}, \mathbf{I}\}$	b	?	?	
	total	0.6		1.0

• Observed fact: p(a, 0) + p(b, 0) = 0.6

• Encoded as a constraint:  $E_p(f_1) = 0.6$  where:

• 
$$f_1(x, y) = \begin{cases} 1 & \text{if } y = 0 \\ 0 & \text{otherwise} \end{cases}$$
  
•  $E_p(f_1) = \sum_{(x,y) \in \{\alpha,b\} \times \{0,1\}} p(x,y) f_1(x,y)$ 

## Outline

Statistical Models for NLP

Maximum Likelihood Estimation (MLE)

Maximum Entropy Modeling Building ME Models

Log-Linear Models

#### Statistical Models for NLP

- Why modeling
- Prediction & Similarity Models

Maximum Likelihood Estimation (MLE) Working example

Smoothing & Estimator Combination

#### 3 Maximum Entropy Modeling

- Overview
- Building ME Models



## **Probability Model**

Statistical Models for NLP

Maximum Likelihood Estimation (MLE)

Maximum Entropy Modeling Building ME Models

Log-Linear Models There is an infinite set P of probability models consistent with observations:

$$\mathbf{P} = \{ \mathbf{p} \mid \mathsf{E}_{\mathbf{p}}(\mathsf{f}_{\mathfrak{i}}) = \mathsf{E}_{\widetilde{\mathbf{p}}}(\mathsf{f}_{\mathfrak{i}}), \ \forall \mathfrak{i} \}$$

Maximum entropy model

$$\begin{aligned} p^* &= & \underset{p \in P}{\operatorname{argmax}} H(p) \\ &= & \underset{p \in P}{\operatorname{argmax}} \left( -\sum_{x \in \varepsilon} p(x) \log p(x) \right) \end{aligned}$$

## Conditional Probability Model

 For NLP applications, we are usually interested in conditional distributions P(Y|X), thus, the ME model is

$$p^* = \mathop{\text{argmax}}_{p \in P} H(p) = \mathop{\text{argmax}}_{p \in P} H(Y \mid X)$$

where:

Н

$$\begin{aligned} (Y \mid X) &= \sum_{x \in X} p(x) H(Y \mid X = x) \\ &= -\sum_{x \in X} p(x) \sum_{y \in Y} p(y \mid x) \log p(y \mid x) \\ &= -\sum_{x \in X, y \in Y} p(x, y) \log p(y \mid x) \\ &= -\sum_{x \in X, y \in Y} p(x, y) \log \frac{p(x, y)}{p(x)} \end{aligned}$$

Statistical Models for NLP

Maximum Likelihood Estimation (MLE)

Maximum Entropy Modeling Building ME Models

## Parameter Estimation

	Example: Maximum entropy model for translating in to French				
	No constraints				
Statistical Models for	P(x) dans en à au-cours-de pendant				
NLP	0.2 0.2 0.2 0.2 0.2				
Maximum	total 1.0				
Estimation (MLE)	• With constraint $p(ans) + p(en) = 0.3$				
Maximum	P(x) dans en à au-cours-de pendant				
Entropy Modeling	0.15 0.15 0.233 0.233 0.233				
Building ME Models	total <b>0.3</b> 1.0				
Log-Linear Models	<ul> <li>With constraints</li> </ul>				
	p(dans) + p(en) = 0.3; p(en) + p(a) = 0.5				
	Not so easy !				

## Parameter estimation

Exponential models

$$p(y \mid x) = \frac{1}{\mathsf{Z}(x)} \prod_{j=1}^{k} \alpha_{j}^{f_{j}(x,y)} \quad \alpha_{j} > 0, \quad \mathsf{Z}(x) = \sum_{y} \prod_{i=1}^{k} \alpha_{i}^{f_{i}(x,y)}$$

- Can also be formuled as  $p(y \mid x) = \frac{1}{\mathsf{Z}(x)} \exp\left(\sum_{j=1}^k \lambda_j f_j(x, y)\right) \qquad (i.e. \ \lambda_i = \ln \alpha_i)$
- Each model parameter weights the influence of a feature.
- Optimal parameters can be computed with:
  - Generalized Iterative Scaling (GIS) [Darroch & Ratcliff 72]
  - Improved Iterative Scaling (IIS) [Della Pietra et al. 96]
  - Limited Memory BFGS (LM-BFGS) [Malouf 03]

Statistical Models for NLP

Maximum Likelihood Estimation (MLE)

Maximum Entropy Modeling Building ME Models

## Example: Text Categorization

Statistical Models for NLP

Maximum Likelihood Estimation (MLE)

Maximum Entropy Modeling Building ME Models

Log-Linear Models ■ Probabilistic model over W × C (Words × Categories).
 A document is a set of words: d = (w<sub>1</sub>, w<sub>2</sub>...w<sub>N</sub>).
 Each combination w, c ∈ W × C is a feature:

$$f_{w,c}(d,c') = \begin{cases} \frac{N(w,d)}{N(d)} & \text{if } c = c' \\ 0 & \text{otherwise} \end{cases}$$

Disambiguation: Select class with highest  $P(c \mid d)$ 

$$P(c \mid d) = \frac{1}{Z(d)} \exp(\sum_{i} \lambda_i f_i(d, c))$$

## **MEM Summary**

Advantages

- Teoretically well founded
- Enables combination of random context features
- Better probabilistic models than MLE (no smoothing needed)
- General approach (features, events and classes)
- Disadvantages
  - Implicit probabilistic model (joint or conditional probability distribution obtained from model parameters).

ME Models are a particular case of Log-Linear models

Statistical Models for NLP

Maximum Likelihood Estimation (MLE)

Maximum Entropy Modeling Building ME Models

## Outline

Statistical Models for NI P

Maximum Likelihood Estimation (MLE)

Maximum Entropy Modeling

Log-Linear Models

- Why modeling
- Prediction & Similarity Models

Working example

- Smoothing & Estimator Combination
- - Overview
  - Building ME Models



## Log-Linear Models

$$\mathsf{P}(\mathsf{y} \mid \mathsf{x}; \mathbf{w}) = \frac{\exp\left(\mathbf{w} \cdot \mathbf{f}(\mathsf{x}, \mathsf{y})\right)}{\sum_{\mathsf{y}} \exp\left(\mathbf{w} \cdot \mathbf{f}(\mathsf{x}, \mathsf{y})\right)}$$

Statistical Models for NLP

Maximum Likelihood Estimation (MLE)

Maximum Entropy Modeling

Log-Linear Models where

- f(x, y) is a feature vector representing x and y
- w are the parameters of the model
- $\blacksquare \ \mathbf{w} \cdot \mathbf{f}(x,y)$  is a score for x and y
- Z(x) = ∑y exp (w ⋅ f(x, y)) is a normalizer (sums over all possible values y for x); it's sometimes called the *partition function*

#### Statistical Models for NLP

Maximum Likelihood Estimation (MLE)

Maximum Entropy Modeling

Log-Linear Models

## Features, Indicator Features

 $\blacksquare~\mathbf{f}(x,y)$  is a vector of d features representing x and y

 $(\ \mathbf{f_1}(x,y),\ldots,\mathbf{f_j}(x,y),\ldots,\mathbf{f_d}(x,y)\ )$ 

- What's in a feature  $f_j(x, y)$ ?
  - Anything we can compute using x and y
  - Anything that is informative for (or against) x belonging to class y
  - Indicator features: binary-valued features looking at a single simple property

$$\begin{split} \mathbf{f}_{j}(c,b) &= \left\{ \begin{array}{ll} 1 & \text{if } \text{prefix}(c) = Mr \text{ and } b = \text{no} \\ 0 & \text{otherwise} \end{array} \right. \\ \mathbf{f}_{k}(c,b) &= \left\{ \begin{array}{ll} 1 & \text{if } \text{uppercase}(\text{next}(c)) \text{ and } b = \text{yes} \\ 0 & \text{otherwise} \end{array} \right. \end{split}$$

### Features, Parameters, Inner Products

$$\mathsf{P}(\mathsf{y} \mid \mathsf{x}; \mathbf{w}) = \frac{\exp\left(\mathbf{w} \cdot \mathbf{f}(\mathsf{x}, \mathsf{y})\right)}{\sum_{\mathsf{y}} \exp\left(\mathbf{w} \cdot \mathbf{f}(\mathsf{x}, \mathsf{y})\right)}$$

Statistical Models for NLP

Maximum Likelihood Estimation (MLE)

Maximum Entropy Modeling

Log-Linear Models •  $f(x, y) \in \mathbb{R}^d$  is a feature vector with d features •  $w \in \mathbb{R}^d$  is a parameter vector, with d parameters

Inner products (a.k.a. dot products)

$$\mathbf{w} \cdot \mathbf{f}(\mathbf{x}, \mathbf{y}) = \sum_{i=1}^{d} \mathbf{w}_{i} \mathbf{f}_{i}(\mathbf{x}, \mathbf{y})$$

## Log-linear Models

$$\mathsf{P}(y \mid x; \mathbf{w}) = \frac{\exp\left(\mathbf{w} \cdot \mathbf{f}(x, y)\right)}{\sum_{y} \exp\left(\mathbf{w} \cdot \mathbf{f}(x, y)\right)}$$

Statistical Models for NLP

Maximum Likelihood Estimation (MLE)

Maximum Entropy Modeling

Log-Linear Models where

- f(x, y) is a feature vector representing x and y
  - Arbitrary features of x and y are allowed
  - They are provided for the application in turn
- w are the parameters of the model
- Two problems:
  - How to make predictions using P(y | x)
  - How to estimate the parameters w?

## Log-linear Models: Name

• Let's take the log of the conditional probability:

Statistical Models for NLP

Maximum Likelihood Estimation (MLE)

Maximum Entropy Modeling

Log-Linear Models

$$\log P(y \mid x; \mathbf{w}) = \log \frac{\exp \left(\mathbf{w} \cdot \mathbf{f}(x, y)\right)}{\sum_{y} \exp \left(\mathbf{w} \cdot \mathbf{f}(x, y)\right)}$$
$$= \mathbf{w} \cdot \mathbf{f}(x, y) - \log \sum_{y} \exp \left(\mathbf{w} \cdot \mathbf{f}(x, y)\right)$$
$$= \mathbf{w} \cdot \mathbf{f}(x, y) - \log Z(x)$$

• Partition function:  $Z(x) = \sum_{y} \exp(\mathbf{w} \cdot \mathbf{f}(x, y))$ 

 $\blacksquare$  log Z(x) is a constant for a fixed x

In the log space, computations are linear

## Log-linear Models: Making Predictions

■ Given x, what y in {1, ..., L} is most appropriate?

$$\begin{aligned} \mathsf{best}(\mathbf{x}) &= \operatorname*{argmax}_{\mathbf{y} \in \{1, \dots, L\}} \mathsf{P}(\mathbf{y} \mid \mathbf{x}; \mathbf{w}) \\ &= \operatorname*{argmax}_{\mathbf{y} \in \{1, \dots, L\}} \frac{\mathsf{exp}\left(\mathbf{w} \cdot \mathbf{f}(\mathbf{x}, \mathbf{y})\right)}{\mathsf{Z}(\mathbf{x})} \end{aligned}$$

Statistical Models for NLP

Maximum Likelihood Estimation (MLE)

Maximum Entropy Modeling

## Log-linear Models: Making Predictions

■ Given x, what y in {1, ..., L} is most appropriate?

 $best(x) = \underset{y \in \{1, \dots, L\}}{\operatorname{argmax}} P(y \mid x; \mathbf{w})$  $= \underset{y \in \{1, \dots, L\}}{\operatorname{argmax}} \frac{\exp(\mathbf{w} \cdot \mathbf{f}(x, y))}{Z(x)}$  $= \underset{y \in \{1, \dots, L\}}{\operatorname{argmax}} \exp(\mathbf{w} \cdot \mathbf{f}(x, y))$  $= \underset{y \in \{1, \dots, L\}}{\operatorname{argmax}} \mathbf{w} \cdot \mathbf{f}(x, y)$ 

Statistical Models for NLP

Maximum Likelihood Estimation (MLE)

Maximum Entropy Modeling

## Log-linear Models: Making Predictions

■ Given x, what y in {1, . . . , L} is most appropriate?

 $best(x) = \underset{y \in \{1, \dots, L\}}{\operatorname{argmax}} P(y \mid x; w)$  $= \underset{y \in \{1, \dots, L\}}{\operatorname{argmax}} \frac{\exp(w \cdot f(x, y))}{Z(x)}$  $= \underset{y \in \{1, \dots, L\}}{\operatorname{argmax}} \exp(w \cdot f(x, y))$  $= \underset{u \in \{1, \dots, L\}}{\operatorname{argmax}} w \cdot f(x, y)$ 

Predictions only require simple inner products (linear)No need to exponentiate!

Statistical Models for NLP

Maximum Likelihood Estimation (MLE)

Maximum Entropy Modeling

## Log-linear Models: Computing Probabilities

Statistical Models for NLP

Maximum Likelihood Estimation (MLE)

Maximum Entropy Modeling

Log-Linear Models

$$\mathsf{P}(y \mid x; \mathbf{w}) = \frac{\mathsf{exp}\left(\mathbf{w} \cdot \mathbf{f}(x, y)\right)}{\mathsf{Z}(x)}$$

- Sometimes we will be interested in computing P(y | x)
   It can be used as a measure of confidence, e.g. P(yes | c) = 0.51 versus P(yes | c) = 0.99
- We need to compute:

$$\mathsf{Z}(x) = \sum_{\mathsf{y} = \{1, \dots, L\}} \exp\left(\mathbf{w} \cdot \mathbf{f}(x, \mathsf{y})\right)$$

Fast as long as L is not too large

## Parameter Estimation in Log-linear Models

How to estimate model parameters w given a training set:

$$\left\{ (x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)}) \right\}$$

Let's define the conditional log-likelihood of the data:

$$L(\mathbf{w}) = \frac{1}{m} \sum_{k=1}^{m} \log \mathsf{P}(\boldsymbol{y}^{(k)} | \boldsymbol{x}^{(k)}; \mathbf{w})$$

- L(w) measures how well w explains the data. A good value for w will give a high value for P(y<sup>(k)</sup>|x<sup>(k)</sup>; w) for all k = 1...m.
- We want w that maximizes L(w)

Statistical Models for NLP

Maximum Likelihood Estimation (MLE)

Maximum Entropy Modeling

## Parameter Estimation in Log-Linear Models

We pose it as an optimization problemFind:

$$\mathbf{w}^* = \operatorname*{argmax}_{\mathbf{w} \in \mathbb{R}^d} L(\mathbf{w})$$

where

- But low-frequency features may end up having large weights (i.e. overfitting)
- We need a regularization factor that penalizes solutions with a large norm (similar to norm-minimization in SVM):

$$L'(\mathbf{w}) = \frac{1}{m} \sum_{k=1}^{m} \log \mathsf{P}(\boldsymbol{y}^{(k)} | \boldsymbol{x}^{(k)}; \mathbf{w}) - \frac{\lambda}{2} ||\mathbf{w}||^2$$

 where λ is a parameter to control the trade-off between fitting the data and model complexity. Tuned experimentally.

Statistical Models for NLP

Maximum Likelihood Estimation (MLE)

Maximum Entropy Modeling

## Parameter Estimation in Log-Linear Models

So we want to find:

 $\mathbf{w}^{*}$ 

Statistical Models for NLP

Maximum Likelihood Estimation (MLE)

Maximum Entropy Modeling

$$= \operatorname{argmax}_{\mathbf{w} \in \mathbb{R}^{d}} L'(\mathbf{w})$$

$$= \operatorname{argmax}_{\mathbf{w} \in \mathbb{R}^{d}} \left( \frac{1}{m} \sum_{k=1}^{m} \log P(\mathbf{y}^{(k)} | \mathbf{x}^{(k)}; \mathbf{w}) - \frac{\lambda}{2} ||\mathbf{w}||^{2} \right)$$

- In general there is no analytical solution to this optimization
- ... but it is a convex function ⇒ We use iterative techniques, i.e. gradient-based optimization
- Very fast algorithms exist (e.g. LBFGS)

## Parameter Estimation in Log-Linear Models : Gradient step

- Initialize w = 0
- Repeat

• Compute gradient  $\delta = (\delta_1, \dots, \delta_d)$ , where:

$$\delta_j = \frac{\partial L'(\mathbf{w})}{\partial \mathbf{w}_j} \quad \forall j = 1 \dots d$$

$$\beta^* = \operatorname*{argmax}_{\beta \in \mathbb{R}} L'(\mathbf{w} + \beta \delta)$$

Move w in the direction of the gradient

$$\mathbf{w} \leftarrow \mathbf{w} + \beta^* \delta$$

• until convergence  $(\|\delta\| < \varepsilon)$ 

Statistical Models for NLP

Maximum Likelihood Estimation (MLE)

Maximum Entropy Modeling

## Log-linear Models: Computing the Gradient

Statistical Models for NLP

Maximum Likelihood Estimation (MLE)

Maximum Entropy Modeling

Log-Linear Models

$$\begin{split} \frac{\partial L'(\mathbf{w})}{\partial \mathbf{w}_j} &= \frac{1}{m} \sum_{k=1}^m \mathbf{f}_j(x^{(k)}, y^{(k)}) \\ &- \sum_{k=1}^m \sum_{y \in \{1, \dots, L\}} \mathsf{P}(y | x^{(k)}; \mathbf{w}) \ \mathbf{f}_j(x^{(k)}, y) \\ &- \lambda \mathbf{w}_j \end{split}$$

- First term: observed mean feature value
- Second term: expected feature value under current w
- In the optimal, observed = expected

# Maximum log-likelihood log-linear models correspond to Maximum Entropy models

## Example: Identifying Sentence Boundaries

Statistical Models for NLP

Maximum Likelihood Estimation (MLE)

Maximum Entropy Modeling

Log-Linear Models The president lives in Washington, D.C. The presidents met in Washington D.C. in 2010. Mr. Wayne is young. Mr. Wayne is a Ph.D. I got 98.5%! What?

Goal: given a text, identify tokens that end a sentence

- Candidate characters: . ! ?
- Candidate tokens: tokens containing candidate characters
- Given a candidate token in a *context* decide whether it ends a sentence or not

## Example: Sentence Boundaries

■ Candidate: punctuation sign + context

c = < sign, prefix, suffix, previous, next >
 Assume access to annotated data:

b	sign	prefix	suffix	prev	next
no		D D D C	С.	Washington,	The
no	•	Mr		2010.	Wayne

- Let's take a probabilistic approach:
  - P(yes | c): conditional probability of c being end of sentence
  - P(no | c): conditional probability of c not being e.o.s.
  - Obviously,  $\mathsf{P}(\texttt{yes} \mid c) + \mathsf{P}(\texttt{no} \mid c) = 1$
  - Predict yes if P(yes | c) > 0.5
- How to model P(yes | c) and P(no | c)?

Statistical Models for NLP

Maximum Likelihood Estimation (MLE)

Maximum Entropy Modeling

### Example system: Identifying Sentence Boundaries (Reynar and Ratnaparkhi '97)

- Candidate: punctuation sign + context
  - c = < sign, prefix, suffix, previous, next >
- **Goal**: estimate P(yes | c) and P(no | c)
- Feature templates:
  - The prefix
  - 2 The suffix
  - 3 The word previous
  - 4 The word next
  - **5** Whether prefix or suffix are in ABBREVIATIONS
    - ABBREVIATIONS: list of all training tokens that contain a

       and are *not* sentence boundaries
  - **6** Whether previous or next are in ABBREVIATIONS
- Actual features are generated by applying each template to each training example

Statistical Models for NLP

Maximum Likelihood Estimation (MLE)

Maximum Entropy Modeling

#### Example system: Identifying Sentence Boundaries (Reynar and Ratnaparkhi '97) FEATURE TEMPLATES 1 The prefix 2 The suffix 3 The word previous 4 The word next Statistical Models for 5 Whether prefix or suffix are in ABBREVIATIONS 6 Whether previous or next are in ABBREVIATIONS Maximum < b=no punc=. pref=Mr suff= prev=2010. next=Wayne > Likelihood Estimation GENERATED FEATURES Maximum $\mathbf{f}_1(\mathbf{c}, \mathbf{b}) = \begin{cases} 1 & \text{if } \mathsf{pref}(\mathbf{c}) = \mathsf{Mr} \\ & \text{and } \mathbf{b} = \mathsf{no} \\ 0 & \text{otherwise} \end{cases} \quad \mathbf{f}_4(\mathbf{c}, \mathbf{b}) = \begin{cases} 1 & \text{if } \mathsf{next}(\mathbf{c}) = \mathsf{Wayne} \\ & \text{and } \mathbf{b} = \mathsf{no} \\ 0 & \text{otherwise} \end{cases}$ Modeling Log-Linear $f_{2}(c,b) = \begin{cases} 1 & \text{if suff}(c) = \texttt{NULL} \\ & \text{and } b = \texttt{no} \\ 0 & \text{otherwise} \end{cases} \quad f_{5}(c,b) = \begin{cases} 1 & \text{if } (\texttt{abbr}(\texttt{pref}(c)) \text{ or } \texttt{abbr}(\texttt{suff}(c))) \\ & \text{and } b = \texttt{no} \\ 0 & \text{otherwise} \end{cases}$ $f_3(c,b) = \begin{cases} 1 & \text{if } prev(c) = 2010. \\ & \text{and } b = no \\ 0 & \text{otherwise} \end{cases} \quad f_6(c,b) = \begin{cases} 1 & \text{if } (abbr(prev(c)) \text{ or } abbr(next(c))) \\ & \text{and } b = no \\ 0 & \text{otherwise} \end{cases}$

NI P

(MLE)

Entropy

Models

#### Example System: Identifying Sentence Boundaries (Reynar and Ratnaparkhi '97)

	training sentences	test accuracy
Statistical Models for	500	96.5%
NLP	1000	97.3%
Maximum Likelihood	2000	97.3%
Estimation	4000	97.6%
(Maximum	8000	97.6%
Entropy	16000	97.8%
l og-l inear	39441	98.0%
Models		

Corpus: Wall Street Journal, English