

Master in Artificial Intelligence

Advanced Human Language Technologies Statistical Models of Language

Statistical
Models for
NLP

Prediction
Models

Prediction
Model
Estimation:
MLE

Prediction
Model
Estimation:
Log-Linear &
MaxEnt



UNIVERSITAT POLITÈCNICA DE CATALUNYA
BARCELONATECH

Facultat d'Informàtica de Barcelona



Outline

1 Statistical Models for NLP

- Why modeling
- Prediction & Similarity Models

2 Prediction Models

- Overview

3 Prediction Model Estimation: MLE

- Overview
- Smoothing & Estimator Combination

4 Prediction Model Estimation: Log-Linear & MaxEnt

- Introduction
- Log-Linear Models
- Maximum Entropy Models
- Examples
- Summary

Outline

- 1 Statistical Models for NLP
 - Why modeling
 - Prediction & Similarity Models
- 2 Prediction Models
 - Overview
- 3 Prediction Model Estimation: MLE
 - Overview
 - Smoothing & Estimator Combination
- 4 Prediction Model Estimation: Log-Linear & MaxEnt
 - Introduction
 - Log-Linear Models
 - Maximum Entropy Models
 - Examples
 - Summary

Statistical
Models for
NLP

Why modeling

Prediction
Models

Prediction
Model
Estimation:
MLE

Prediction
Model
Estimation:
Log-Linear &
MaxEnt

We model to make predictions



Statistical
Models for
NLP

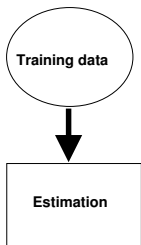
Why modeling

Prediction
Models

Prediction
Model
Estimation:
MLE

Prediction
Model
Estimation:
Log-Linear &
MaxEnt

We model to make predictions



Statistical
Models for
NLP

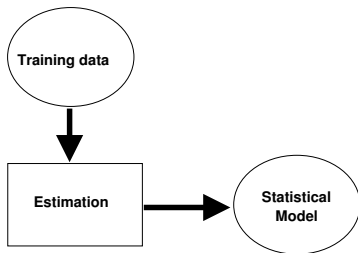
Why modeling

Prediction
Models

Prediction
Model
Estimation:
MLE

Prediction
Model
Estimation:
Log-Linear &
MaxEnt

We model to make predictions



Statistical
Models for
NLP

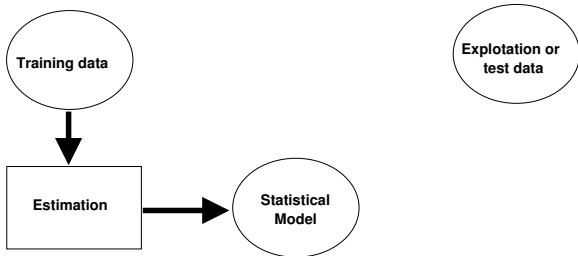
Why modeling

Prediction
Models

Prediction
Model
Estimation:
MLE

Prediction
Model
Estimation:
Log-Linear &
MaxEnt

We model to make predictions



Statistical
Models for
NLP

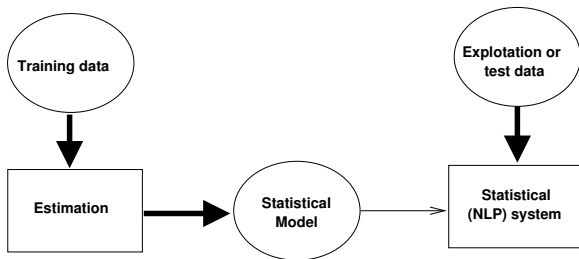
Why modeling

Prediction
Models

Prediction
Model
Estimation:
MLE

Prediction
Model
Estimation:
Log-Linear &
MaxEnt

We model to make predictions



Statistical
Models for
NLP

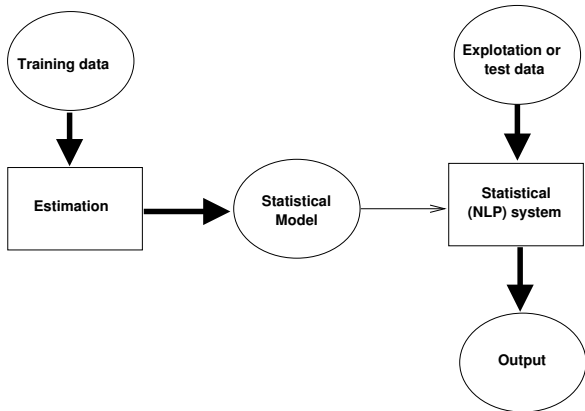
Why modeling

Prediction
Models

Prediction
Model
Estimation:
MLE

Prediction
Model
Estimation:
Log-Linear &
MaxEnt

We model to make predictions



Statistical
Models for
NLP

Why modeling

Prediction
Models

Prediction
Model
Estimation:
MLE

Prediction
Model
Estimation:
Log-Linear &
MaxEnt

Outline

1 Statistical Models for NLP

- Why modeling
- Prediction & Similarity Models

2 Prediction Models

- Overview

3 Prediction Model Estimation: MLE

- Overview
- Smoothing & Estimator Combination

4 Prediction Model Estimation: Log-Linear & MaxEnt

- Introduction
- Log-Linear Models
- Maximum Entropy Models
- Examples
- Summary

Statistical
Models for
NLP

Prediction &
Similarity Models

Prediction
Models

Prediction
Model
Estimation:
MLE

Prediction
Model
Estimation:
Log-Linear &
MaxEnt

Prediction Models & Similarity Models

Statistical
Models for
NLP

Prediction &
Similarity Models

Prediction
Models

Prediction
Model
Estimation:
MLE

Prediction
Model
Estimation:
Log-Linear &
MaxEnt

- **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). [

Outline

- 1 Statistical Models for NLP
 - Why modeling
 - Prediction & Similarity Models
- 2 Prediction Models
 - Overview
- 3 Prediction Model Estimation: MLE
 - Overview
 - Smoothing & Estimator Combination
- 4 Prediction Model Estimation: Log-Linear & MaxEnt
 - Introduction
 - Log-Linear Models
 - Maximum Entropy Models
 - Examples
 - Summary

Statistical
Models for
NLP

Prediction
Models

Prediction
Model
Estimation:
MLE

Prediction
Model
Estimation:
Log-Linear &
MaxEnt

Outline

- 1 Statistical Models for NLP
 - Why modeling
 - Prediction & Similarity Models
- 2 Prediction Models
 - Overview
- 3 Prediction Model Estimation: MLE
 - Overview
 - Smoothing & Estimator Combination
- 4 Prediction Model Estimation: Log-Linear & MaxEnt
 - Introduction
 - Log-Linear Models
 - Maximum Entropy Models
 - Examples
 - Summary

Statistical
Models for
NLP

Prediction
Models
Overview

Prediction
Model
Estimation:
MLE

Prediction
Model
Estimation:
Log-Linear &
MaxEnt

Prediction Models

Example: Noisy Channel Model (Shannon 48)



NLP Applications

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

Given \mathbf{o} , we want to find the most likely \mathbf{i}

$$\operatorname{argmax}_{\mathbf{i}} P(\mathbf{i} | \mathbf{o}) = \operatorname{argmax}_{\mathbf{i}} P(\mathbf{o}, \mathbf{i}) = \operatorname{argmax}_{\mathbf{i}} P(\mathbf{i})P(\mathbf{o} | \mathbf{i})$$

Using Prediction Models

- **Estimation:** Using data to infer information about distributions. A.k.a. *learning*.
- **Classification:** Make predictions based on past behaviour
- In general, ML models estimate (i.e. *learn*) conditional probability distributions $P(\text{target}|\text{features})$, e.g.:
 - language identification given word or subword features.
 - document category given words in it.
 - word PoS given context information.
 - ...
- Many NLP tasks require a posterior search step to find the best combination of predictions.

Outline

- 1 Statistical Models for NLP
 - Why modeling
 - Prediction & Similarity Models
- 2 Prediction Models
 - Overview
- 3 Prediction Model Estimation: MLE**
 - Overview
 - Smoothing & Estimator Combination
- 4 Prediction Model Estimation: Log-Linear & MaxEnt
 - Introduction
 - Log-Linear Models
 - Maximum Entropy Models
 - Examples
 - Summary

Statistical
Models for
NLP

Prediction
Models

Prediction
Model
Estimation:
MLE

Prediction
Model
Estimation:
Log-Linear &
MaxEnt

Outline

- 1 Statistical Models for NLP
 - Why modeling
 - Prediction & Similarity Models
- 2 Prediction Models
 - Overview
- 3 Prediction Model Estimation: MLE
 - Overview
 - Smoothing & Estimator Combination
- 4 Prediction Model Estimation: Log-Linear & MaxEnt
 - Introduction
 - Log-Linear Models
 - Maximum Entropy Models
 - Examples
 - Summary

Statistical
Models for
NLP

Prediction
Models

Prediction
Model
Estimation:
MLE

Overview

Prediction
Model
Estimation:
Log-Linear &
MaxEnt

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 σ^2
- Zipf's Laws. Data sparseness. Smoothing techniques.

$P(x, y)$	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

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. $P(w_i | w_{i-n+1} \dots w_{i-1})$
- 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	$20,000^2 = 4 \times 10^8$
trigram	$20,000^3 = 8 \times 10^{12}$
four-gram	$20,000^4 = 1.6 \times 10^{17}$

Language model sizes for a 20,000 words vocabulary

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:

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

■ MLE (Maximum Likelihood Estimation)

$$P_{MLE}(w_1 \dots w_n) = \frac{C(w_1 \dots w_n)}{N}$$

$$P_{MLE}(w_n | w_1 \dots w_{n-1}) = \frac{C(w_1 \dots w_n)}{C(w_1 \dots w_{n-1})}$$

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

Brief Parenthesis: Zipf's Laws

Zipf's Laws (1929)

- Word frequency is inversely proportional to its rank (speaker/hearer minimum effort) $f \sim 1/r$
- Number of senses is proportional to frequency root $m \sim \sqrt{f}$
- 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 occur only once

Outline

- 1 Statistical Models for NLP
 - Why modeling
 - Prediction & Similarity Models
- 2 Prediction Models
 - Overview
- 3 Prediction Model Estimation: MLE
 - Overview
 - Smoothing & Estimator Combination
 - Adding counts
 - Discounting counts
 - Combining Estimators
- 4 Prediction Model Estimation: Log-Linear & MaxEnt
 - Introduction
 - Log-Linear Models
 - Maximum Entropy Models
 - Examples
 - Summary

Statistical
Models for
NLP

Prediction
Models

Prediction
Model
Estimation:
MLE

Smoothing &
Estimator
Combination

Prediction
Model
Estimation:
Log-Linear &
MaxEnt

Outline

- 1 Statistical Models for NLP
 - Why modeling
 - Prediction & Similarity Models
- 2 Prediction Models
 - Overview
- 3 Prediction Model Estimation: MLE**
 - Overview
 - Smoothing & Estimator Combination**
 - Adding counts
 - Discounting counts
 - Combining Estimators
- 4 Prediction Model Estimation: Log-Linear & MaxEnt
 - Introduction
 - Log-Linear Models
 - Maximum Entropy Models
 - Examples
 - Summary

Statistical
Models for
NLP

Prediction
Models

Prediction
Model
Estimation:
MLE

Adding counts

Prediction
Model
Estimation:
Log-Linear &
MaxEnt

Laplace's Law

LAPLACE'S LAW (adding one count)

General rule:

$$P_{\text{LAP}}(X = x) = \frac{C(X = x) + 1}{N + B}$$

N: number of observations of X
B: number of potentially observable values for X

N-gram probability:

$$P_{\text{LAP}}(w_1 \dots w_n) = \frac{C(w_1 \dots w_n) + 1}{N + B}$$

N: total number of n-gram observations
B: number of potentially observable different n-grams

N-gram conditional probability:

$$P_{\text{LAP}}(w_n | w_1 \dots w_{n-1}) = \frac{C(w_1 \dots w_n) + 1}{C(w_1 \dots w_{n-1}) + B}$$

B: number of potentially observable w_n values

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

Lidstone's Law

LIDSTONE'S LAW (adding λ counts, with $\lambda < 1$)

General rule:

$$P_{\text{LAP}}(X = x) = \frac{C(X = x) + \lambda}{N + B\lambda}$$

N : number of observations of X
 B : number of potentially observable values for X

N-gram probability:

$$P_{\text{LAP}}(w_1 \dots w_n) = \frac{C(w_1 \dots w_n) + \lambda}{N + B\lambda}$$

N : total number of n -gram observations
 B : number of potentially observable different n -grams

N-gram conditional probability:

$$P_{\text{LAP}}(w_n | w_1 \dots w_{n-1}) = \frac{C(w_1 \dots w_n) + \lambda}{C(w_1 \dots w_{n-1}) + B\lambda}$$

B : number of potentially observable w_n values

Equivalent to linear interpolation between MLE and uniform prior:

$$\text{with } \mu = N/(N + B\lambda); P_{\text{LID}}(X = x) = \mu \frac{C(X = x)}{N} + (1 - \mu) \frac{1}{B}$$

Outline

- 1 Statistical Models for NLP
 - Why modeling
 - Prediction & Similarity Models
- 2 Prediction Models
 - Overview
- 3 Prediction Model Estimation: MLE
 - Overview
 - Smoothing & Estimator Combination
 - Adding counts
 - Discounting counts
 - Combining Estimators
- 4 Prediction Model Estimation: Log-Linear & MaxEnt
 - Introduction
 - Log-Linear Models
 - Maximum Entropy Models
 - Examples
 - Summary

Statistical
Models for
NLP

Prediction
Models

Prediction
Model
Estimation:
MLE

Discounting counts

Prediction
Model
Estimation:
Log-Linear &
MaxEnt

Absolute Discounting

ABSOLUTE DISCOUNTING (discount δ counts, with $0 < \delta < 1$)

General rule:

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

N_0 : number of possible values for X observed 0 times

N-gram probability:

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

N-gram conditional probability:

$$P_{\text{ABS}}(w_n | w_1 \dots w_{n-1}) = \begin{cases} \frac{C(w_1 \dots w_n)-\delta}{C(w_1 \dots w_{n-1})} & \text{if } C(w_1 \dots w_n) > 0 \\ \frac{(B-N_0)\delta/N_0}{C(w_1 \dots w_{n-1})} & \text{otherwise} \end{cases}$$

N_0 : number of possible values for w_n observed 0 times

Linear Discounting

LINEAR DISCOUNTING (discount a proportion α of counts)

General rule:

$$P_{\text{LIN}}(X = x) = \begin{cases} (1 - \alpha) \frac{C(X=x)}{N} & \text{if } C(X = x) > 0 \\ \alpha/N_0 & \text{otherwise} \end{cases}$$

N_0 : number of possible values for X observed 0 times

N-gram probability:

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

N_0 : number of possible n -grams observed 0 times

N-gram conditional probability:

$$P_{\text{LIN}}(w_n | w_1 \dots w_{n-1}) = \begin{cases} (1 - \alpha) \frac{C(w_1 \dots w_n)}{C(w_1 \dots w_{n-1})} & \text{if } C(w_1 \dots w_n) > 0 \\ \alpha/N_0 & \text{otherwise} \end{cases}$$

N_0 : number of possible values for w_n observed 0 times

Outline

- 1 Statistical Models for NLP
 - Why modeling
 - Prediction & Similarity Models
- 2 Prediction Models
 - Overview
- 3 Prediction Model Estimation: MLE**
 - Overview
 - Smoothing & Estimator Combination**
 - Adding counts
 - Discounting counts
 - Combining Estimators**
- 4 Prediction Model Estimation: Log-Linear & MaxEnt
 - Introduction
 - Log-Linear Models
 - Maximum Entropy Models
 - Examples
 - Summary

Statistical
Models for
NLP

Prediction
Models

Prediction
Model
Estimation:
MLE

Combining
Estimators

Prediction
Model
Estimation:
Log-Linear &
MaxEnt

Combining Estimators

COMBINING ESTIMATORS FOR CONDITIONAL N-GRAM PROBABILITIES

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:

$$P_{BO}(w_n | h_i) = \begin{cases} \alpha_{h_i} \frac{C(h_i, w_n)}{C(h_i)} & \text{if } C(h_i, w_n) > k \\ \delta_{h_{i-1}} P_{BO}(w_n | h_{i-1}) & \text{otherwise} \end{cases}$$

$h_i = w_{n-i} \dots w_{n-1}$: recent history of length i .

α_{h_i} : remaining proportion after discount.

$\delta_{h_{i-1}}$: mass assigned to back-off distribution.

- Constant: $\alpha_{h_i} = 1 - \delta_{h_{i-1}} = K; \forall h$
- Katz back-off (based on Good-Turing estimation)

Outline

- 1 Statistical Models for NLP
 - Why modeling
 - Prediction & Similarity Models
- 2 Prediction Models
 - Overview
- 3 Prediction Model Estimation: MLE
 - Overview
 - Smoothing & Estimator Combination
- 4 Prediction Model Estimation: Log-Linear & MaxEnt
 - Introduction
 - Log-Linear Models
 - Maximum Entropy Models
 - Examples
 - Summary

Statistical
Models for
NLP

Prediction
Models

Prediction
Model
Estimation:
MLE

Prediction
Model
Estimation:
Log-Linear &
MaxEnt

Outline

- 1 Statistical Models for NLP
 - Why modeling
 - Prediction & Similarity Models
- 2 Prediction Models
 - Overview
- 3 Prediction Model Estimation: MLE
 - Overview
 - Smoothing & Estimator Combination
- 4 Prediction Model Estimation: Log-Linear & MaxEnt
 - Introduction
 - Log-Linear Models
 - Maximum Entropy Models
 - Examples
 - Summary

Statistical
Models for
NLP

Prediction
Models

Prediction
Model
Estimation:
MLE

Prediction
Model
Estimation:
Log-Linear &
MaxEnt

Introduction

Example: Identifying Sentence Boundaries

EXAMPLE

*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

- Object to classify: punctuation sign + context
 $x = \langle \text{sign, prefix, suffix, previous, next} \rangle$

- Assume access to *annotated* data:

y	sign	prefix	suffix	prev	next
no	.	D	C.	Washington,	The
yes	.	D.C		Washington,	The
no	.	Mr		2010.	Wayne

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

Outline

- 1 Statistical Models for NLP
 - Why modeling
 - Prediction & Similarity Models
- 2 Prediction Models
 - Overview
- 3 Prediction Model Estimation: MLE
 - Overview
 - Smoothing & Estimator Combination
- 4 Prediction Model Estimation: Log-Linear & MaxEnt
 - Introduction
 - **Log-Linear Models**
 - Maximum Entropy Models
 - Examples
 - Summary

Statistical
Models for
NLP

Prediction
Models

Prediction
Model
Estimation:
MLE

Prediction
Model
Estimation:
Log-Linear &
MaxEnt

Log-Linear Models

Log-Linear Models

Log-Linear models take the form:

$$P(y | x; \mathbf{w}) = \frac{\exp(\mathbf{w} \cdot \mathbf{f}(x, y))}{Z(x)} = \frac{\exp(\mathbf{w} \cdot \mathbf{f}(x, y))}{\sum_y \exp(\mathbf{w} \cdot \mathbf{f}(x, y))}$$

where

- $\mathbf{f}(x, y) \in \mathbb{R}^d$ is a feature vector representing a *context* x and a *label* y
- $\mathbf{w} \in \mathbb{R}^d$ is a vector containing the *parameters* of the model
- $\mathbf{w} \cdot \mathbf{f}(x, y) = \sum_{i=1}^d w_i f_i(x, y)$ is a *score* for x and y
- $Z(x) = \sum_y \exp(\mathbf{w} \cdot \mathbf{f}(x, y))$ is a normalizer (sum of scores for all possible values for y); a.k.a *partition function*

Features, Indicator Features

- $\mathbf{f}(x, y) \in \mathbb{R}^d$ is a vector of d features encoding some information about x and y

$$\mathbf{f}(x, y) = (f_1(x, y), \dots, f_k(x, y), \dots, f_d(x, y))$$

- What's in a feature $f_k(x, y)$?
 - *Anything* we can compute using x and y (and *suitable* for the task at hand)
 - *Anything* informative for (or against) x belonging to class y .
 - Usually, they are **indicator features**: binary-valued features looking at a (simple) property.

$$f_1(x, y) = \begin{cases} 1 & \text{if } \text{prefix}(x) = Mr \text{ and } y = \text{no} \\ 0 & \text{otherwise} \end{cases}$$

$$f_2(x, y) = \begin{cases} 1 & \text{if } \text{uppercase}(\text{next}(x)) \text{ and } y = \text{yes} \\ 0 & \text{otherwise} \end{cases}$$

Log-linear Models: Name

- Let's take the **log** of the conditional probability:

$$\begin{aligned}\log P(y | x; \mathbf{w}) &= \log \frac{\exp(\mathbf{w} \cdot \mathbf{f}(x, y))}{\sum_y \exp(\mathbf{w} \cdot \mathbf{f}(x, y))} \\ &= \mathbf{w} \cdot \mathbf{f}(x, y) - \log \sum_y \exp(\mathbf{w} \cdot \mathbf{f}(x, y)) \\ &= \mathbf{w} \cdot \mathbf{f}(x, y) - \log Z(x)\end{aligned}$$

- $\log Z(x)$ is a constant for a fixed x
- In the **log** space, computations are **linear**

Log-linear Models: Making Predictions

- Given x , what y in $\{1, \dots, L\}$ is most appropriate?

$$\text{best_label}(x) = \underset{y \in \{1, \dots, L\}}{\operatorname{argmax}} P(y | x; \mathbf{w})$$

Log-linear Models: Making Predictions

- Given x , what y in $\{1, \dots, L\}$ is most appropriate?

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

Statistical
Models for
NLP

Prediction
Models

Prediction
Model
Estimation:
MLE

Prediction
Model
Estimation:
Log-Linear &
MaxEnt

Log-Linear Models

Log-linear Models: Making Predictions

- Given x , what y in $\{1, \dots, L\}$ is most appropriate?

$$\begin{aligned}\text{best_label}(x) &= \operatorname{argmax}_{y \in \{1, \dots, L\}} P(y | x; \mathbf{w}) \\ &= \operatorname{argmax}_{y \in \{1, \dots, L\}} \frac{\exp(\mathbf{w} \cdot \mathbf{f}(x, y))}{Z(x)} \\ &= \operatorname{argmax}_{y \in \{1, \dots, L\}} \exp(\mathbf{w} \cdot \mathbf{f}(x, y))\end{aligned}$$

Log-linear Models: Making Predictions

- Given x , what y in $\{1, \dots, L\}$ is most appropriate?

$$\begin{aligned}\text{best_label}(x) &= \operatorname{argmax}_{y \in \{1, \dots, L\}} P(y | x; \mathbf{w}) \\ &= \operatorname{argmax}_{y \in \{1, \dots, L\}} \frac{\exp(\mathbf{w} \cdot \mathbf{f}(x, y))}{Z(x)} \\ &= \operatorname{argmax}_{y \in \{1, \dots, L\}} \exp(\mathbf{w} \cdot \mathbf{f}(x, y)) \\ &= \operatorname{argmax}_{y \in \{1, \dots, L\}} \mathbf{w} \cdot \mathbf{f}(x, y)\end{aligned}$$

Log-linear Models: Making Predictions

- Given x , what y in $\{1, \dots, L\}$ is most appropriate?

$$\begin{aligned}\text{best_label}(x) &= \operatorname{argmax}_{y \in \{1, \dots, L\}} P(y | x; \mathbf{w}) \\ &= \operatorname{argmax}_{y \in \{1, \dots, L\}} \frac{\exp(\mathbf{w} \cdot \mathbf{f}(x, y))}{Z(x)} \\ &= \operatorname{argmax}_{y \in \{1, \dots, L\}} \exp(\mathbf{w} \cdot \mathbf{f}(x, y)) \\ &= \operatorname{argmax}_{y \in \{1, \dots, L\}} \mathbf{w} \cdot \mathbf{f}(x, y)\end{aligned}$$

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

Log-linear Models: Computing Probabilities

- Sometimes we will be interested in computing $P(y | x)$, not just the argmax.

$P(y | x)$ can be used as a measure of confidence in the prediction, e.g.:

$$P(\text{yes} | x) = 0.51 \quad \text{vs.} \quad P(\text{yes} | x) = 0.99$$

- Since: $P(y | x; \mathbf{w}) = \frac{\exp(\mathbf{w} \cdot \mathbf{f}(x, y))}{Z(x)}$

we need to compute: $Z(x) = \sum_{y=\{1, \dots, L\}} \exp(\mathbf{w} \cdot \mathbf{f}(x, y))$

- Feasible as long as L is not too large

Parameter Estimation in Log-linear Models

- How to estimate model parameters \mathbf{w} given a training set:

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

- We define the conditional log-likelihood of the data:

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

- $L(\mathbf{w})$ measures how well \mathbf{w} explains the data. A good value for \mathbf{w} will give a high value for $P(\mathbf{y}^{(k)} | \mathbf{x}^{(k)}; \mathbf{w})$ for all $k = 1 \dots m$.
- We want \mathbf{w} that *maximizes* $L(\mathbf{w})$

Parameter Estimation in Log-Linear Models

Estimating model parameters is an optimization problem, aiming to find:

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

But low-frequency features may end up having large weights (i.e. overfitting).

Thus, we need a **regularization** factor that penalizes solutions with a large norm (similar to norm-minimization in SVM), redefining $L(\mathbf{w})$ as:

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

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

Parameter Estimation in Log-Linear Models

So we want to find:

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

- In general there is no analytical solution to this optimization.
- ... but it is a **convex** function \Rightarrow numerical optimization iterative techniques, i.e. gradient-based optimization, may be used.
- Very fast algorithms exist (e.g. LBFGS).

Parameter Estimation in Log-Linear Models :

Gradient step

- Initialize $\mathbf{w} = \mathbf{0}$

- Repeat

- Compute gradient $\boldsymbol{\delta} = (\delta_1, \dots, \delta_d)$, where:

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

- Compute step size

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

- Move \mathbf{w} in the direction of the gradient

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

- until convergence ($\|\boldsymbol{\delta}\| < \epsilon$)

Log-linear Models: Computing the Gradient

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

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

Outline

- 1 Statistical Models for NLP
 - Why modeling
 - Prediction & Similarity Models
- 2 Prediction Models
 - Overview
- 3 Prediction Model Estimation: MLE
 - Overview
 - Smoothing & Estimator Combination
- 4 Prediction Model Estimation: Log-Linear & MaxEnt
 - Introduction
 - Log-Linear Models
 - **Maximum Entropy Models**
 - Examples
 - Summary

Statistical
Models for
NLP

Prediction
Models

Prediction
Model
Estimation:
MLE

Prediction
Model
Estimation:
Log-Linear &
MaxEnt

Maximum Entropy
Models

Example

Example: Maximum entropy model for translating English prepositions *in* and *on* to French.

No observations (no constraints)

$P(x, y)$	dans	en	à	sur	au-cours-de	pendant	selon	
in	0.07	0.07	0.07	0.07	0.07	0.07	0.07	
on	0.07	0.07	0.07	0.07	0.07	0.07	0.07	
total								1.0

Example

Example: Maximum entropy model for translating English prepositions *in* and *on* to French.

Observations (constraints):

$$p(\text{en} \vee \grave{\text{a}}) = 0.6$$

$P(x, y)$	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		0.6						1.0

Example

Example: Maximum entropy model for translating English prepositions *in* and *on* to French.

Observations (constraints):

$$p(\text{en} \vee \grave{\text{a}}) = 0.6; \quad p((\text{en} \vee \grave{\text{a}}) \wedge \text{in}) = 0.4$$


$P(x, y)$	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		0.6						1.0

Example

Example: Maximum entropy model for translating English prepositions *in* and *on* to French.

Observations (constraints):

$$p(\text{en} \vee \grave{\text{a}}) = 0.6; \quad p((\text{en} \vee \grave{\text{a}}) \wedge \text{in}) = 0.4; \quad p(\text{in}) = 0.5$$


$P(x, y)$	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

Example

Example: Maximum entropy model for translating English prepositions *in* and *on* to French.

Observations (constraints):

$$p(\text{en} \vee \grave{\text{a}}) = 0.6; \quad p((\text{en} \vee \grave{\text{a}}) \wedge \text{in}) = 0.4; \quad p(\text{in}) = 0.5; \quad p(\text{sur}) = 0.1$$


$P(x, y)$	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				0.1				1.0

Example

Example: Maximum entropy model for translating English prepositions *in* and *on* to French.

Observations (constraints):

$$p(\text{en} \vee \grave{\text{a}}) = 0.6; \quad p((\text{en} \vee \grave{\text{a}}) \wedge \text{in}) = 0.4; \quad p(\text{in}) = 0.5; \quad p(\text{sur}) = 0.1$$

$P(x, y)$	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				0.1				1.0

Not so easy...

Maximum entropy Models

Statistical
Models for
NLP

Prediction
Models

Prediction
Model
Estimation:
MLE

Prediction
Model
Estimation:
Log-Linear &
MaxEnt

Maximum Entropy
Models

MaxEnt models: dual formulation of **Log-Linear models**.

ME principles:

- Do not assume anything about non-observed events.
- Find the most uniform (maximum entropy, less informed) probability distribution that matches the observations.

ME Modeling

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

$$f_i : \mathcal{E} \rightarrow \{0, 1\}$$

- The desired model p must satisfy the constraints:
The expectation predicted by model p for any feature f_i must match the observed expectation for f_i
i.e.:

$$\begin{aligned} E_p(f_i) &= E_{\tilde{p}}(f_i) \quad \forall i \\ \sum_{x \in \mathcal{E}} p(x) f_i(x) &= \sum_{x \in \mathcal{E}} \tilde{p}(x) f_i(x) \quad \forall i \end{aligned}$$

Probability Model

- There is an infinite set \mathcal{P} of probability models consistent with observations:

$$\mathcal{P} = \{p \mid E_p(f_i) = E_{\tilde{p}}(f_i), \quad \forall i\}$$

- Maximum entropy model

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

Parameter estimation

- ME models are exponential models, same as log-linear models:

$$P(\mathbf{y} | \mathbf{x}) = \frac{\exp(\mathbf{w} \cdot \mathbf{f}(\mathbf{x}, \mathbf{y}))}{\sum_{\mathbf{y} \in \mathcal{L}} \exp(\mathbf{w} \cdot \mathbf{f}(\mathbf{x}, \mathbf{y}))}$$

- Each model parameter weights the influence of a feature.
- Same convex optimization algorithms are used (e.g. **LM-BFGS**,)
- Optimized cost function is different, but optimum corresponds to the same \mathbf{w} than a log-linear model.

Outline

- 1 Statistical Models for NLP
 - Why modeling
 - Prediction & Similarity Models
- 2 Prediction Models
 - Overview
- 3 Prediction Model Estimation: MLE
 - Overview
 - Smoothing & Estimator Combination
- 4 Prediction Model Estimation: Log-Linear & MaxEnt
 - Introduction
 - Log-Linear Models
 - Maximum Entropy Models
 - **Examples**
 - Summary

Statistical
Models for
NLP

Prediction
Models

Prediction
Model
Estimation:
MLE

Prediction
Model
Estimation:
Log-Linear &
MaxEnt

Examples

Example: Identifying Sentence Boundaries

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

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?

- Candidate characters: . ! ?
- Candidate tokens: tokens containing candidate characters (e.g. *D.C.*, *2010.*, *young.*, *Mr.*, *Ph.D.*, *98.5%!*, *What?*)
- Given a candidate character in a *context* decide whether the token ends a sentence.

Identifying Sentence Boundaries: Formulation

- **Object to classify:** punctuation sign + context
 $x = \langle \text{sign, prefix, suffix, previous, next} \rangle$
- **Class:** $y \in \{\text{yes, no}\}$
- **Probabilistic approach:**
 - **Goal:** estimate $P(\text{yes} | x)$ and $P(\text{no} | x)$
 - $P(\text{yes} | x) + P(\text{no} | x) = 1$
 - Predict yes if $P(\text{yes} | x) > 0.5$

Identifying Sentence Boundaries: Data set

- Assume access to *annotated* data:

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?

- Data set

y	sign	prefix	suffix	prev	next
no	.	D	C.	Washington,	The
yes	.	D.C		Washington,	The
no	.	D	C.	Washington	in
no	.	D.C		Washington	in
yes	.	2010		in	Mr.
no	.	Mr		2010.	Wayne
yes	.	young		is	Mr.
no	.	Mr		young.	Wayne
no	.	Ph	D.	a	I
yes	.	Ph.D		a	I
no	.	98	5%!	got	What?
yes	!	98.5%!		got	What?
yes	?	What		98.5%!	<eof>

Identifying Sentence Boundaries: Feature templates

- Feature templates:
 - 1 The **prefix**: part of the token before the **sign**.
 - 2 The **suffix**: part of the token after the **sign**.
 - 3 The **previous** token.
 - 4 The **next** token.
 - 5 Whether **prefix** or **suffix** are in ABBREVIATIONS
 - 6 Whether **previous** or **next** are in ABBREVIATIONS
- ABBREVIATIONS: list of all tokens in training data that contain **sign** and are *not* sentence boundaries.
- Actual features are generated by applying *each template* to *each training example*.

Identifying Sentence Boundaries: Feature generation

FEATURE TEMPLATES

1 The **prefix**

3 The **previous** token

5 Whether **prefix** or **suffix** are in ABBREVIATIONS

2 The **suffix**

4 The **next** token

6 Whether **previous** or **next** are in ABBREVIATIONS

TRAINING DATA EXAMPLE 1

$y = \text{no}; \quad x = \langle \text{sign}=. \text{ pref}=\text{Mr} \text{ suff}=\text{ prev}=\text{2010. next}=\text{Wayne} \rangle$

GENERATED FEATURES

$$f_{1_Mr_no}(x, y) = \begin{cases} 1 & \text{if } \text{pref}(x) = \text{Mr} \\ & \text{and } y = \text{no} \\ 0 & \text{otherwise} \end{cases}$$

$$f_{2_null_no}(x, y) = \begin{cases} 1 & \text{if } \text{suff}(x) = \text{NULL} \\ & \text{and } y = \text{no} \\ 0 & \text{otherwise} \end{cases}$$

$$f_{3_2010_no}(x, y) = \begin{cases} 1 & \text{if } \text{prev}(x) = \text{2010.} \\ & \text{and } y = \text{no} \\ 0 & \text{otherwise} \end{cases}$$

$$f_{4_Wayne_no}(x, y) = \begin{cases} 1 & \text{if } \text{next}(x) = \text{Wayne} \\ & \text{and } y = \text{no} \\ 0 & \text{otherwise} \end{cases}$$

$$f_{5_no}(x, y) = \begin{cases} 1 & \text{if } (\text{abbr}(\text{pref}(x)) \\ & \text{or } \text{abbr}(\text{suff}(x))) \\ & \text{and } y = \text{no} \\ 0 & \text{otherwise} \end{cases}$$

$$f_{6_no}(x, y) = \begin{cases} 1 & \text{if } (\text{abbr}(\text{prev}(x)) \\ & \text{or } \text{abbr}(\text{next}(x))) \\ & \text{and } y = \text{no} \\ 0 & \text{otherwise} \end{cases}$$

Statistical
Models for
NLP

Prediction
Models

Prediction
Model
Estimation:
MLE

Prediction
Model
Estimation:
Log-Linear &
MaxEnt

Examples

Identifying Sentence Boundaries: Feature generation

FEATURE TEMPLATES

1 The **prefix**

3 The **previous** token

5 Whether **prefix** or **suffix** are in ABBREVIATIONS

2 The **suffix**

4 The **next** token

6 Whether **previous** or **next** are in ABBREVIATIONS

TRAINING DATA EXAMPLE 2

$y = \text{yes}$ $x = \langle \text{punc}=. \text{ pref}=\text{D.C} \text{ suff}=\text{ } \text{ prev}=\text{,} \text{ next}=\text{The} \rangle$

GENERATED FEATURES

$$f_{1_{\text{D.C.}, \text{yes}}}(x, y) = \begin{cases} 1 & \text{if } \text{pref}(x) = \text{D.C} \\ & \text{and } y = \text{yes} \\ 0 & \text{otherwise} \end{cases}$$

$$f_{2_{\text{null}, \text{yes}}}(x, y) = \begin{cases} 1 & \text{if } \text{suff}(x) = \text{NULL} \\ & \text{and } y = \text{yes} \\ 0 & \text{otherwise} \end{cases}$$

$$f_{3_{\text{,}, \text{yes}}}(x, y) = \begin{cases} 1 & \text{if } \text{prev}(x) = \text{,} \\ & \text{and } y = \text{yes} \\ 0 & \text{otherwise} \end{cases}$$

$$f_{4_{\text{The}, \text{yes}}}(x, y) = \begin{cases} 1 & \text{if } \text{next}(x) = \text{The} \\ & \text{and } y = \text{yes} \\ 0 & \text{otherwise} \end{cases}$$

$$f_{5_{\text{yes}}}(x, y) = \begin{cases} 1 & \text{if } (\text{abbr}(\text{pref}(x)) \\ & \text{or } \text{abbr}(\text{suff}(x))) \\ & \text{and } y = \text{yes} \\ 0 & \text{otherwise} \end{cases}$$

$$f_{6_{\text{yes}}}(x, y) = \begin{cases} 1 & \text{if } (\text{abbr}(\text{prev}(x)) \\ & \text{or } \text{abbr}(\text{next}(x))) \\ & \text{and } y = \text{yes} \\ 0 & \text{otherwise} \end{cases}$$

Statistical
Models for
NLP

Prediction
Models

Prediction
Model
Estimation:
MLE

Prediction
Model
Estimation:
Log-Linear &
MaxEnt

Examples

Example: Text Categorization

Goal: given a text, classify it according to a set of predefined classes

British hurdler Sarah Claxton is confident she can win her first major medal at next month's European Indoor Championships in Madrid.

- Candidate classes: business, entertainment, politics, sport, tech.
- Given a document decide which category (or categories) it belongs to.

Text Categorization: Formulation

- **Object to classify:** document (as a set of words)
 $x = \langle \text{word}_1, \text{word}_2, \dots, \text{word}_n \rangle$
- **Class:** $y \in L = \{\text{biz}, \text{ent}, \text{pol}, \text{spo}, \text{tech}\}$
- **Probabilistic approach:**
 - **Goal:** estimate $P(y | x) \quad \forall y \in L$
 - $\sum_{y \in L} P(y | x) = 1$
 - **Predict as output class:**
 - Class with highest probability: $\text{argmax}_y P(y | x)$
 - Any class with probability over a threshold:
 $\{y | P(y | x) > k\}$

Text Categorization: Data set

- Assume access to *annotated* data:

spo	British hurdler Sarah Claxton is confident she can win her first major medal at next month's European Indoor Championships in Madrid.
pol	The Labour Party will hold its 2006 autumn conference in Manchester and not Blackpool, it has been confirmed.
tech	Microsoft has warned PC users to update their systems with the latest security fixes for flaws in Windows programs.

Text Categorization: Feature templates

- Feature templates:
 - 1 The occurrence of a **word** in the document
- Actual features are generated by applying the *template* to *each training example*.

Text Categorization: Feature generation

FEATURE TEMPLATES

- 1 A **word** occurring in the document.

TRAINING DATA EXAMPLE 1

$y = \text{spo}$

$x = \{\text{British hurdler Sarah Claxton confident win first major medal next month European Indoor Championships Madrid}\}$

GENERATED FEATURES

$$f_{1_\text{British_spo}}(x, y) = \begin{cases} 1 & \text{if British} \in x \\ & \text{and } y = \text{spo} \\ 0 & \text{otherwise} \end{cases}$$

$$f_{1_\text{hurdler_spo}}(x, y) = \begin{cases} 1 & \text{if hurdler} \in x \\ & \text{and } y = \text{spo} \\ 0 & \text{otherwise} \end{cases}$$

...

$$f_{1_win_spo}(x, y) = \begin{cases} 1 & \text{if win} \in x \\ & \text{and } y = \text{spo} \\ 0 & \text{otherwise} \end{cases}$$

...

$$f_{1_medal_spo}(x, y) = \begin{cases} 1 & \text{if medal} \in x \\ & \text{and } y = \text{spo} \\ 0 & \text{otherwise} \end{cases}$$

$$f_{1_next_spo}(x, y) = \begin{cases} 1 & \text{if next} \in x \\ & \text{and } y = \text{spo} \\ 0 & \text{otherwise} \end{cases}$$

...

$$f_{1_Madrid_spo}(x, y) = \begin{cases} 1 & \text{if Madrid} \in x \\ & \text{and } y = \text{spo} \\ 0 & \text{otherwise} \end{cases}$$

Text Categorization: Feature generation

FEATURE TEMPLATES

- 1 A **word** occurring in the document.

TRAINING DATA EXAMPLE 2

$y = \text{tech}$

$x = \{\text{Microsoft warned PC users update systems latest security fixes flaws Windows programs}\}$

GENERATED FEATURES

$$f_{1_Microsoft_tech}(x, y) = \begin{cases} 1 & \text{if Microsoft} \in x \\ & \text{and } y = \text{tech} \\ 0 & \text{otherwise} \end{cases}$$

$$f_{1_latest_tech}(x, y) = \begin{cases} 1 & \text{if latest} \in x \\ & \text{and } y = \text{tech} \\ 0 & \text{otherwise} \end{cases}$$

$$f_{1_warned_tech}(x, y) = \begin{cases} 1 & \text{if warned} \in x \\ & \text{and } y = \text{tech} \\ 0 & \text{otherwise} \end{cases}$$

$$f_{1_security_tech}(x, y) = \begin{cases} 1 & \text{if security} \in x \\ & \text{and } y = \text{tech} \\ 0 & \text{otherwise} \end{cases}$$

...

...

$$f_{1_systems_tech}(x, y) = \begin{cases} 1 & \text{if systems} \in x \\ & \text{and } y = \text{tech} \\ 0 & \text{otherwise} \end{cases}$$

$$f_{1_programs_tech}(x, y) = \begin{cases} 1 & \text{if programs} \in x \\ & \text{and } y = \text{tech} \\ 0 & \text{otherwise} \end{cases}$$

...

Outline

- 1 Statistical Models for NLP
 - Why modeling
 - Prediction & Similarity Models
- 2 Prediction Models
 - Overview
- 3 Prediction Model Estimation: MLE
 - Overview
 - Smoothing & Estimator Combination
- 4 Prediction Model Estimation: Log-Linear & MaxEnt
 - Introduction
 - Log-Linear Models
 - Maximum Entropy Models
 - Examples
 - **Summary**

Statistical
Models for
NLP

Prediction
Models

Prediction
Model
Estimation:
MLE

Prediction
Model
Estimation:
Log-Linear &
MaxEnt

Summary

Log-linear Models 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).