

Introduction

Statistical
Models for
NLP

Maximum
Likelihood
Estimation
(MLE)

Maximum
Entropy
Modeling

References

Model Estimation: Maximum Likelihood vs. Maximum Entropy

DMKM - Universitat Politècnica de Catalunya

1 Introduction

2 Statistical Models for NLP

- Overview
- Prediction & Similarity Models
- Statistical Inference of Models for NLP

3 Maximum Likelihood Estimation (MLE)

- Overview
- Smoothing & Estimator Combination

4 Maximum Entropy Modeling

- Overview
- Building ME Models
- Application to NLP

5 References

Introduction

Statistical
Models for
NLPMaximum
Likelihood
Estimation
(MLE)Maximum
Entropy
Modeling

References

- Random variable: Function on a stochastic process.
 $X : \Omega \longrightarrow \mathcal{R}$
- Continuous and discrete random variables.
- Probability mass (or density) function, Frequency function:
 $p(x) = P(X = x)$.
Discrete R.V.: $\sum_x p(x) = 1$
Continuous R.V: $\int_{-\infty}^{\infty} p(x)dx = 1$
- Distribution function: $F(x) = P(X \leq x)$
- Expectation and variance, standard deviation
 $E(X) = \mu = \sum_x xp(x)$
 $VAR(X) = \sigma^2 = E((X - E(X))^2) = \sum_x (x - \mu)^2 p(x)$

- Joint probability mass function: $p(x, y)$
- Marginal distribution:

$$p_X(x) = \sum_y p(x, y) \quad p_{X|Y}(x | y) = \frac{p(x, y)}{p_Y(y)}$$

$$p_Y(y) = \sum_x p(x, y)$$

Simplified Polynesian. Sequences of C-V syllables: Two random variables C,V

P(C,V)	p	t	k	
a	1/16	3/8	1/16	1/2
i	1/16	3/16	0	1/4
u	0	3/16	1/16	1/4
	1/8	3/4	1/8	

$$P(p | i) = ?$$

$$P(a | t \vee k) = ?$$

$$P(a \vee i | p) = ?$$

Introduction

Statistical
Models for
NLPMaximum
Likelihood
Estimation
(MLE)Maximum
Entropy
Modeling

References

Introduction

Statistical
Models for
NLP

Maximum
Likelihood
Estimation
(MLE)

Maximum
Entropy
Modeling

References

- Random samples

- Sample variables:

Sample mean: $\bar{\mu}_n = \frac{1}{n} \sum_{i=1}^n x_i$

Sample variance: $s_n^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{\mu}_n)^2$.

- Law of Large Numbers: as n increases, $\bar{\mu}_n$ and s_n^2 converge to μ and σ^2
- Estimators: Sample variables used to estimate real parameters.

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
- Data sparseness problem. Smoothing techniques.

$P(a, b)$	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

Maximum Entropy Estimation (MEE)

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

Observations:

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

$P(a, b)$	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		$\underbrace{\hspace{1.5cm}}$ 0.6						1.0

Maximum Entropy Estimation (MEE)

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

Observations:

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

$P(a, b)$	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		<div style="text-align: center;"> $\underbrace{\hspace{1.5cm}}$ 0.6 </div>						1.0

Maximum Entropy Estimation (MEE)

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

Observations:

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

$P(a, b)$	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		<div style="text-align: center;"> $\underbrace{\hspace{1.5cm}}$ 0.6 </div>						1.0

1 Introduction

2 Statistical Models for NLP

- Overview
- Prediction & Similarity Models
- Statistical Inference of Models for NLP

3 Maximum Likelihood Estimation (MLE)

- Overview
- Smoothing & Estimator Combination

4 Maximum Entropy Modeling

- Overview
- Building ME Models
- Application to NLP

5 References

1 Introduction

2 Statistical Models for NLP

- Overview
- Prediction & Similarity Models
- Statistical Inference of Models for NLP


3 Maximum Likelihood Estimation (MLE)

- Overview
- Smoothing & Estimator Combination

4 Maximum Entropy Modeling

- Overview
- Building ME Models
- Application to NLP

5 References



Training data

Introduction

Statistical
Models for
NLP

Overview

Maximum
Likelihood
Estimation
(MLE)

Maximum
Entropy
Modeling

References

Introduction

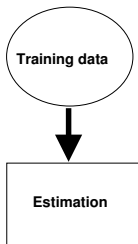
Statistical Models for NLP

Overview

Maximum Likelihood Estimation (MLE)

Maximum Entropy Modeling

References



Introduction

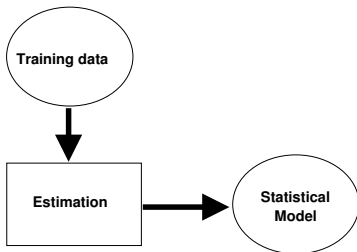
Statistical Models for NLP

Overview

Maximum Likelihood Estimation (MLE)

Maximum Entropy Modeling

References



Introduction

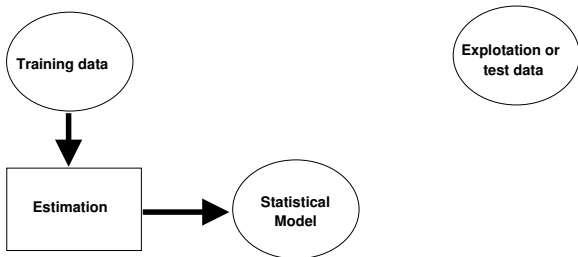
Statistical Models for NLP

Overview

Maximum Likelihood Estimation (MLE)

Maximum Entropy Modeling

References



Introduction

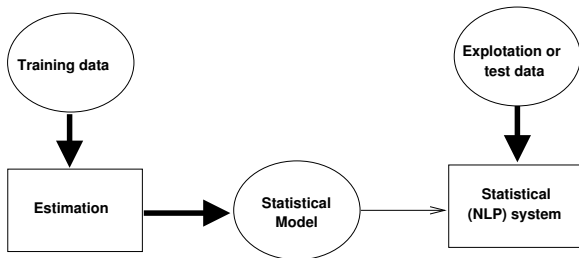
Statistical Models for NLP

Overview

Maximum Likelihood Estimation (MLE)

Maximum Entropy Modeling

References



Introduction

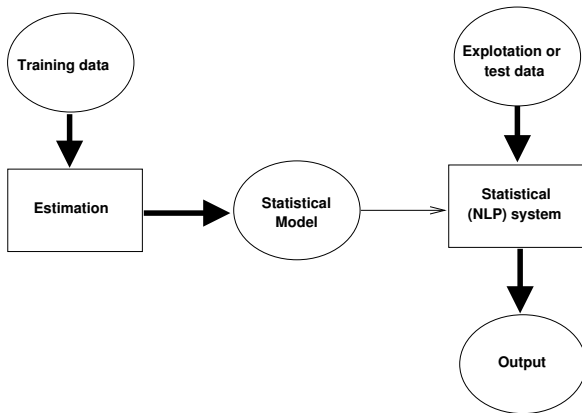
Statistical Models for NLP

Overview

Maximum Likelihood Estimation (MLE)

Maximum Entropy Modeling

References



1 Introduction

2 Statistical Models for NLP

- Overview
- Prediction & Similarity Models
- Statistical Inference of Models for NLP

3 Maximum Likelihood Estimation (MLE)

- Overview
- Smoothing & Estimator Combination

4 Maximum Entropy Modeling

- Overview
- Building ME Models
- Application to NLP

5 References

Introduction

Statistical Models for NLP

Prediction & Similarity Models

Maximum Likelihood Estimation (MLE)

Maximum Entropy Modeling

References

- Prediction Models: Able to *predict* probabilities of future events, knowing past and present.
- Similarity Models: Able to compute *similarities* between objects (may predict, too).
 - Compare feature-vector/feature-set represented objects.
 - Compare distribution-vector represented objects
 - Used to group objects (clustering, data analysis, pattern discovery, ...)
 - If objects are “present and past” situations, computing similarities may be used as a prediction (memory-based ML techniques).

Example: Document representation

- Documents are represented as vectors in a high dimensional \mathbb{R}^N space.
- Dimensions are word forms, lemmas, NEs, ...
- Values may be either binary or real-valued (count, frequency, ...)

$$\vec{x} = \begin{bmatrix} x_1 \\ \vdots \\ x_N \end{bmatrix} \quad \vec{x}^T = [x_1 \dots x_N] \quad |\vec{x}| = \sqrt{\sum_{i=1}^N x_i^2}$$

Introduction

Statistical
Models for
NLP

Prediction &
Similarity
Models

Maximum
Likelihood
Estimation
(MLE)

Maximum
Entropy
Modeling

References

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

Introduction

Statistical Models for NLP

Prediction & Similarity Models

Maximum Likelihood Estimation (MLE)

Maximum Entropy Modeling

References

Introduction

Statistical Models for NLP

Statistical Inference of Models for NLP

Maximum Likelihood Estimation (MLE)

Maximum Entropy Modeling

References

1 Introduction

2 Statistical Models for NLP

- Overview
- Prediction & Similarity Models
- Statistical Inference of Models for NLP

3 Maximum Likelihood Estimation (MLE)

- Overview
- Smoothing & Estimator Combination

4 Maximum Entropy Modeling

- Overview
- Building ME Models
- Application to NLP

5 References

Introduction

Statistical Models for NLP

Statistical Inference of Models for NLP

Maximum Likelihood Estimation (MLE)

Maximum Entropy Modeling

References

- Using data to infer information about distributions
 - Parametric / non-parametric estimation
 - Finding good estimators: MLE, MEE, ...
- Example: Language Modeling (Shannon game), N-gram models.
- Predictions based on past behaviour
 - Target / classification features → Independence assumptions
 - Equivalence classes (bins).
Granularity: discrimination vs. statistical reliability

Introduction

Statistical Models for NLP

Statistical Inference of Models for NLP

Maximum Likelihood Estimation (MLE)

Maximum Entropy Modeling

References

- Predicting the next word in a sequence, given the *history* or *context*. $P(w_n \mid w_1 \dots w_{n-1})$
- Markov assumption: Only *local* context (of size $n - 1$) is taken into account. $P(w_i \mid w_{i-n+1} \dots w_{i-1})$
- 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

1 Introduction

2 Statistical Models for NLP

- Overview
- Prediction & Similarity Models
- Statistical Inference of Models for NLP

3 Maximum Likelihood Estimation (MLE)

- Overview
- Smoothing & Estimator Combination

4 Maximum Entropy Modeling

- Overview
- Building ME Models
- Application to NLP

5 References

1 Introduction

2 Statistical Models for NLP

- Overview
- Prediction & Similarity Models
- Statistical Inference of Models for NLP

3 Maximum Likelihood Estimation (MLE)

- Overview
- Smoothing & Estimator Combination

4 Maximum Entropy Modeling

- Overview
- Building ME Models
- Application to NLP

5 References

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:

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

■ MLE (Maximum Likelihood Estimation)

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

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

- No probability mass for unseen events
- Unsuitable for NLP
- Data sparseness, Zipf's Law

Introduction

Statistical
Models for
NLP

Maximum
Likelihood
Estimation
(MLE)

Overview

Maximum
Entropy
Modeling

References

1 Introduction

2 Statistical Models for NLP

- Overview
- Prediction & Similarity Models
- Statistical Inference of Models for NLP

3 Maximum Likelihood Estimation (MLE)

- Overview
- Smoothing & Estimator Combination

4 Maximum Entropy Modeling

- Overview
- Building ME Models
- Application to NLP

5 References

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

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

- N_k : Number of classes (n-grams) observed k times.
- N_k^A : Number of classes (n-grams) observed k times on data subset A .
- B : Number of equivalence classes or bins (number of potentially observable n-grams).

■ 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

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

- Usually $\lambda = 0.5$, *Expected Likelihood Estimation*.
- Equivalent to linear interpolation between MLE and uniform prior, with $\mu = N/(N + B\lambda)$,

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

■ Absolute Discounting

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

■ Linear Discounting

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

Introduction

Statistical
Models for
NLP

Maximum
Likelihood
Estimation
(MLE)

Smoothing &
Estimator
Combination

Maximum
Entropy
Modeling

References

- *Notation:* γ stands for $w_1 \dots w_n$.
- Divide the train corpus in two subsets, A and B.
- Define: $T_r^{AB} = \sum_{\gamma: C_A(\gamma)=r} C_B(\gamma)$
- **Held Out Estimator**

$$P_{HO}(w_1 \dots w_n) = \frac{T_{C_A(\gamma)}^{AB}}{N_{C_A(\gamma)}^A} \times \frac{1}{N}$$

- **Cross Validation** (deleted estimation)

$$P_{DEL}(w_1 \dots w_n) = \frac{T_{C_A(\gamma)}^{AB} + T_{C_B(\gamma)}^{BA}}{N_{C_A(\gamma)}^A + N_{C_B(\gamma)}^B} \times \frac{1}{N}$$

- **Cross Validation** (Leave-one-out)

■ Simple Linear Interpolation

$$\begin{aligned} P_{LI}(w_n \mid w_{n-2}, w_{n-1}) &= \\ &= \lambda_1 P_1(w_n) + \lambda_2 P_2(w_n \mid w_{n-1}) + \lambda_3 P_3(w_n \mid w_{n-2}, w_{n-1}) \end{aligned}$$

■ General Linear Interpolation

$$P_{LI}(w_n \mid h) = \sum_{i=1}^k \lambda_i(h) P_i(w \mid h_i)$$

■ Katz's Backing-off

$$P_{BO}(w_i \mid w_{i-n+1} \dots w_{i-1}) = \begin{cases} (1 - d_{w_{i-n+1} \dots w_{i-1}}) \frac{C(w_{i-n+1} \dots w_i)}{C(w_{i-n+1} \dots w_{i-1})} & \text{if } C(w_{i-n+1} \dots w_i) > k \\ \alpha_{w_{i-n+1} \dots w_{i-1}} P_{BO}(w_i \mid w_{i-n+2} \dots w_{i-1}) & \text{otherwise} \end{cases}$$

Introduction

Statistical
Models for
NLP

Maximum
Likelihood
Estimation
(MLE)

Smoothing &
Estimator
Combination

Maximum
Entropy
Modeling

References

1 Introduction

2 Statistical Models for NLP

- Overview
- Prediction & Similarity Models
- Statistical Inference of Models for NLP

3 Maximum Likelihood Estimation (MLE)

- Overview
- Smoothing & Estimator Combination

4 Maximum Entropy Modeling

- Overview
- Building ME Models
- Application to NLP

5 References

1 Introduction

2 Statistical Models for NLP

- Overview
- Prediction & Similarity Models
- Statistical Inference of Models for NLP

3 Maximum Likelihood Estimation (MLE)

- Overview
- Smoothing & Estimator Combination

4 Maximum Entropy Modeling

- Overview
- Building ME Models
- Application to NLP

5 References

Introduction

Statistical
Models for
NLPMaximum
Likelihood
Estimation
(MLE)Maximum
Entropy
Modeling
Overview

References

- 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:

$p(a, b)$	0	1	
x	?	?	
y	?	?	
total	0.6	1.0	

Observations

$p(a, b)$	0	1	
x	0.5	0.1	
y	0.1	0.3	
total	0.6	1.0	

One possible $p(a, b)$

$p(a, b)$	0	1	
x	0.3	0.2	
y	0.3	0.2	
total	0.6	1.0	

Max. Entropy $p(a, b)$

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

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

- The desired model must satisfy the constraints:

$$E_p(f_i) = E_{\tilde{p}}(f_i) \quad \forall i$$

where:

$$E_p(f_i) = \sum_{x \in \varepsilon} p(x) f_i(x) \quad \text{expectation of model } p.$$

$$E_{\tilde{p}}(f_i) = \sum_{x \in \varepsilon} \tilde{p}(x) f_i(x) \quad \text{observed expectation.}$$

- Example:

$$\varepsilon = \{x, y\} \times \{0, 1\}$$

$p(a, b)$	0	1
x	?	?
y	?	?
total	0.6	1.0

- Observed fact: $p(x, 0) + p(y, 0) = 0.6$
- Encoded as a constraint: $E_p(f_1) = 0.6$

where:

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

- 1 Introduction
- 2 Statistical Models for NLP
 - Overview
 - Prediction & Similarity Models
 - Statistical Inference of Models for NLP
- 3 Maximum Likelihood Estimation (MLE)
 - Overview
 - Smoothing & Estimator Combination
- 4 **Maximum Entropy Modeling**
 - Overview
 - **Building ME Models**
 - Application to NLP
- 5 References

Introduction

Statistical Models for NLP

Maximum Likelihood Estimation (MLE)

Maximum Entropy Modeling

Building ME Models

References

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

$$P = \{p \mid E_p(f_i) = E_{\tilde{p}}(f_i), \forall i = 1 \dots k\}$$

- Maximum entropy model

$$p^* = \operatorname{argmax}_{p \in P} H(p)$$

$$H(p) = - \sum_{x \in \mathcal{E}} p(x) \log p(x)$$

- For NLP applications, we are usually interested in conditional distributions $P(A|B)$, thus:

$$E_{\tilde{p}}(f_j) = \sum_{a,b} \tilde{p}(a,b) f_j(a,b)$$

$$E_p(f_j) = \sum_{a,b} \tilde{p}(b) p(a|b) f_j(a,b)$$

- Maximum entropy model

$$p^* = \operatorname{argmax}_{p \in P} H(p)$$

$$H(p) = H(A|B) = - \sum_{a,b} \tilde{p}(b) p(a|b) \log p(a|b)$$

Example: Maximum entropy model for translating *in* to French

■ No constraints

$P(x)$	dans	en	à	au-cours-de	pendant	
	0.2	0.2	0.2	0.2	0.2	
total						1.0

■ With constraint $p(dans) + p(en) = 0.3$

$P(x)$	dans	en	à	au-cours-de	pendant	
	0.15	0.15	0.233	0.233	0.233	
total	0.3					1.0

■ With constraints $p(dans) + p(en) = 0.3$; $p(en) + p(à) = 0.5$
...Not so easy !

- Exponential models. (Lagrange multipliers optimization)

$$p(a | b) = \frac{1}{Z(b)} \prod_{j=1}^k \alpha_j^{f_j(a,b)} \quad \alpha_j > 0$$

$$Z(b) = \sum_a \prod_{i=1}^k \alpha_i^{f_i(a,b)}$$

- also formulated as

$$p(a | b) = \frac{1}{Z(b)} \exp(\sum_{j=1}^k \lambda_j f_j(a, b))$$

$$\lambda_j = \ln \alpha_j$$

- Each model parameter weights the influence of a feature.
- Optimal parameters (ME model) can be computed with:
 - GIS. Generalized Iterative Scaling (Darroch & Ratcliff 72)
 - IIS. Improved Iterative Scaling (Della Pietra et al. 96)
 - LM-BFGS. Limited Memory BFGS (Malouf 03)

Improved Iterative Scaling (IIS)

Input: Feature functions $f_1 \dots f_n$, empirical distribution $\tilde{p}(a, b)$

Output: λ_i^* parameters for optimal model p^*

Start with $\lambda_i = 0$ for all $i \in \{1 \dots n\}$

Repeat

For each $i \in \{1 \dots n\}$ **do**

let $\Delta\lambda_i$ be the solution to

$$\sum_{a,b} \tilde{p}(b) p(a | b) f_i(a, b) \exp(\Delta\lambda_i \sum_{j=1}^n f_j(a, b)) = \tilde{p}(f_i)$$

$$\lambda_i \leftarrow \lambda_i + \Delta\lambda_i$$

end for

Until all λ_i have converged

Introduction

Statistical
Models for
NLP

Maximum
Likelihood
Estimation
(MLE)

Maximum
Entropy
Modeling
Building ME
Models

References

1 Introduction

2 Statistical Models for NLP

- Overview
- Prediction & Similarity Models
- Statistical Inference of Models for NLP

3 Maximum Likelihood Estimation (MLE)

- Overview
- Smoothing & Estimator Combination

4 Maximum Entropy Modeling

- Overview
- Building ME Models
- Application to NLP

5 References

Introduction

Statistical Models for NLP

Maximum Likelihood Estimation (MLE)

Maximum Entropy Modeling

Application to NLP

References

- Speech processing (Rosenfeld 94)
- Machine Translation (Brown et al 90)
- Morphology (Della Pietra et al. 95)
- Clause boundary detection (Reynar & Ratnaparkhi 97)
- PP-attachment (Ratnaparkhi et al 94)
- PoS Tagging (Ratnaparkhi 96, Black et al 99)
- Partial Parsing (Skut & Brants 98)
- Full Parsing (Ratnaparkhi 97, Ratnaparkhi 99)
- Text Categorization (Nigam et al 99)

- Probabilistic model over $H \times T$

$$h_i = (w_i, w_{i+1}, w_{i+2}, w_{i-1}, w_{i-2}, t_{i-1}, t_{i-2})$$

$$f_j(h_i, t) = \begin{cases} 1 & \text{if } \text{suffix}(w_i) = \text{ing} \wedge t = \text{VBG} \\ 0 & \text{otherwise} \end{cases}$$

- Compute $p^*(h, t)$ using GIS
- Disambiguation algorithm: *beam search*

$$p(t \mid h) = \frac{p(h, t)}{\sum_{t' \in T} p(h, t')}$$

$$p(t_1 \dots t_n \mid w_1 \dots w_n) = \prod_{i=1}^n p(t_i \mid h_i)$$

Introduction

Statistical
Models for
NLP

Maximum
Likelihood
Estimation
(MLE)

Maximum
Entropy
Modeling

Application to
NLP

References

- Probabilistic model over $W \times C$

$$d = (w_1, w_2 \dots w_N)$$

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

- Compute $p^*(c \mid d)$ using IIS
- Disambiguation algorithm: Select class with highest

$$P(c \mid d) = \frac{1}{Z(d)} \exp\left(\sum_i \lambda_i f_i(d, c)\right)$$

■ 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).
- High computational cost of GIS and IIS.
- Overfitting in some cases.

1 Introduction

2 Statistical Models for NLP

- Overview
- Prediction & Similarity Models
- Statistical Inference of Models for NLP

3 Maximum Likelihood Estimation (MLE)

- Overview
- Smoothing & Estimator Combination

4 Maximum Entropy Modeling

- Overview
- Building ME Models
- Application to NLP

5 References

- T. Cover & J. Thomas, **Elements of Information Theory**. John Wiley & Sons, 1991.
- C. Manning & H. Schütze, **Foundations of Statistical Natural Language Processing**. The MIT Press. Cambridge, MA. May 1999.
- D. Jurafsky & J.H. Martin. **Speech and Language Processing: An Introduction to Natural Language Processing, Speech Recognition, and Computational Linguistics**, 2nd edition. Prentice-Hall, 2009.
- A. Berger, S.A. Della Pietra & V.J. Della Pietra, **A Maximum Entropy Approach to Natural Language Processing**. Computational Linguistics, 22(1):39-71, 1996.
- R Malouf, **A comparison of algorithms for maximum entropy parameter estimation**. In Proceedings of the Sixth Conference on Natural Language Learning (CoNLL-2002), Pages 49-55, 2002.
- A. Ratnaparkhi, **Maximum Entropy Models for Natural Language Ambiguity Resolution**. Ph.D Thesis. University of Pennsylvania, 1998.