

Maximum entropy and MEMMs

based on Ch. 6 of Jurafsky & Martin

Grzegorz Chrupała

National Centre for Language Technology
School of Computing
Dublin City University

NCLT Seminar 2007



Outline

- 1 Classification and Sequence labeling
- 2 Linear and Logistic Regression
- 3 Maximum Entropy Modeling
- 4 Maximum Entropy Markov Models

Markov and Shannon



Outline

- 1 Classification and Sequence labeling
- 2 Linear and Logistic Regression
- 3 Maximum Entropy Modeling
- 4 Maximum Entropy Markov Models

Supervised learning

In supervised learning we try to learn a function $h : \mathcal{X} \rightarrow \mathcal{Y}$ where $x \in \mathcal{X}$ are inputs and $y \in \mathcal{Y}$ are outputs.

- Binary classification: $\mathcal{Y} = \{-1, +1\}$
- Multiclass classification: $\mathcal{Y} = \{1, \dots, K\}$ (finite set of labels)
- Regression: $\mathcal{Y} = \mathbb{R}$

The prediction is based on the *feature function* $\Phi : \mathcal{X} \rightarrow \mathcal{F}$ where usually $\mathcal{F} = \mathbb{R}^D$ (D -dimensional vector space)

Sequence labeling

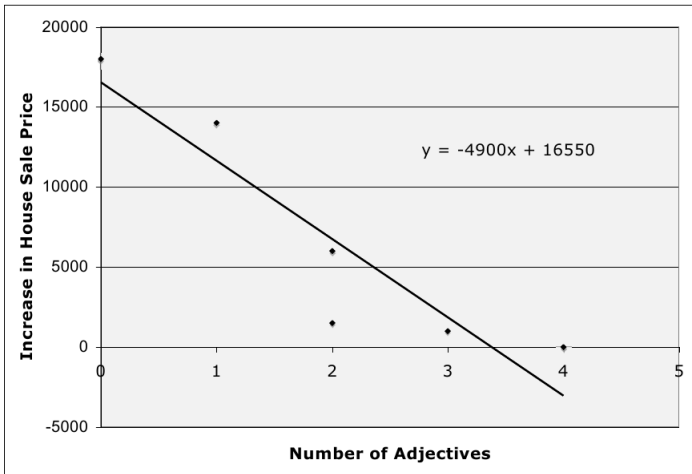
- In sequence labeling the function h to learn is instantiated as $h : \mathcal{W}^n \rightarrow \mathcal{L}^n$ where $w \in \mathcal{W}$ are elements of the sequence and $l \in \mathcal{L}$ are labels.
- Input and output sequences are of equal length
- When this is not the case (i.e. chunking, NER) BIO-encoding can be used
- Maximum Entropy Markov Models combine predictions of a MaxEnt classifier using the Viterbi algorithm to find the best sequence of labels

Outline

- 1 Classification and Sequence labeling
- 2 Linear and Logistic Regression**
- 3 Maximum Entropy Modeling
- 4 Maximum Entropy Markov Models

Linear Regression

- Training data: observations paired with outcomes ($n \in \mathbb{R}$)
- Observations have features (predictors, typically also real numbers)
- The model is a **regression line** $y = ax + b$ which best fits the observations
 - ▶ a is the **slope**
 - ▶ b is the **intercept**
 - ▶ This model has two parameters (or weights)
 - ▶ One feature = x
 - ▶ Example:
 - ★ x = number of vague adjectives in property descriptions
 - ★ y = amount house sold over asking price



Multiple linear regression

- More generally $y = w_0 + \sum_{i=1}^N w_i f_i$, where
 - ▶ y = outcome
 - ▶ w_0 = intercept
 - ▶ $f_1..f_N$ = features vector and $w_1..w_N$ weight vector
- We ignore w_0 by adding a special f_0 feature, then the equation is equivalent to dot product: $y = \mathbf{w} \cdot \mathbf{f}$

Learning linear regression

- Minimize **sum squared error** over the training set of M examples

$$\text{cost}(W) = \sum_{j=0}^M (y_{\text{pred}}^{(j)} - y_{\text{obs}}^{(j)})^2$$

where

$$y_{\text{pred}}^j = \sum_{i=0}^N w_i f_i^{(j)}$$

- Closed-form formula for choosing the best set of weights W is given by:

$$W = (X^T X)^{-1} X^T \vec{y}$$

where the matrix X contains training example features, and \vec{y} is the vector of outcomes.

Logistic regression

- In logistic regression we use the linear model to do classification, i.e. assign probabilities to class labels
- For binary classification, predict $p(y = \text{true}|x)$. But predictions of linear regression model are $\in \mathbb{R}$, whereas $p(y = \text{true}|x) \in [0, 1]$
- Instead predict logit function of the probability:

$$\ln \left(\frac{p(y = \text{true}|x)}{1 - p(y = \text{true}|x)} \right) = \mathbf{w} \cdot \mathbf{f} \quad (1)$$

$$\frac{p(y = \text{true}|x)}{1 - p(y = \text{true}|x)} = e^{\mathbf{w} \cdot \mathbf{f}} \quad (2)$$

- Solving for $p(y = \text{true}|x)$ we obtain:

$$p(y = \text{true}|x) = \frac{e^{\mathbf{w} \cdot \mathbf{f}}}{1 + e^{\mathbf{w} \cdot \mathbf{f}}} \quad (3)$$

$$= \frac{\exp \left(\sum_{i=0}^N w_i f_i \right)}{1 + \exp \left(\sum_{i=0}^N w_i f_i \right)} \quad (4)$$

Logistic regression - classification

- Example x belongs to class *true* if:

$$\frac{p(y = \text{true}|x)}{1 - p(y = \text{true}|x)} > 1 \quad (5)$$

$$e^{\mathbf{w} \cdot \mathbf{f}} > 1 \quad (6)$$

$$\mathbf{w} \cdot \mathbf{f} > 0 \quad (7)$$

$$\sum_{i=0}^N w_i f_i > 0 \quad (8)$$

- The equation $\sum_{i=0}^N w_i f_i = 0$ defines the **hyperplane** in N -dimensional space, with points above this hyperplane belonging to class *true*

Logistic regression - learning

- Conditional likelihood estimation: choose the weights which make the probability of the observed values y be the highest, given the observations x
- For the training set with M examples:

$$\hat{\mathbf{w}} = \operatorname{argmax}_w \prod_{i=0}^M P(y^{(i)} | x^{(i)})$$

- A problem in convex optimization (not covered here)
 - ▶ L-BFGS (Limited-memory Broyden-Fletcher-Goldfarb-Shanno method)
 - ▶ gradient ascent
 - ▶ conjugate gradient
 - ▶ iterative scaling algorithms

Outline

- 1 Classification and Sequence labeling
- 2 Linear and Logistic Regression
- 3 Maximum Entropy Modeling**
- 4 Maximum Entropy Markov Models

Maximum Entropy model

- Logistic regression with more than two classes = **multinomial logistic regression**
- Also known as Maximum Entropy (MaxEnt)
- The MaxEnt equation generalizes (4) above:

$$p(c|x) = \frac{\exp\left(\sum_{i=0}^N w_{ci} f_i\right)}{\sum_{c' \in C} \exp\left(\sum_{i=0}^N w_{c'i} f_i\right)} \quad (9)$$

- The denominator is the normalization factor usually called Z used to make the score into a proper probability distribution

$$p(c|x) = \frac{1}{Z} \exp \sum_{i=0}^N w_{ci} f_i$$

MaxEnt features

- In Maxent modeling normally binary features are used
- Features **depend on classes**: $f_i(c, x) \in \{0, 1\}$
- Those are **indicator features**
- Example x :
Secretariat/NNP is/BEZ expected/VBN to/TO race/VB tomorrow
- Example features:

$$f_1(c, x) = \begin{cases} 1 & \text{if } word_i = race \ \& \ c = NN \\ 0 & \text{otherwise} \end{cases}$$

$$f_2(c, x) = \begin{cases} 1 & \text{if } t_{i-1} = TO \ \& \ c = VB \\ 0 & \text{otherwise} \end{cases}$$

$$f_3(c, x) = \begin{cases} 1 & \text{if } suffix(word_i) = ing \ \& \ c = VBG \\ 0 & \text{otherwise} \end{cases}$$

Binarizing features

- Example x :

Secretariat/NNP is/BEZ expected/VBN to/TO race/VB tomorrow

- Vector of symbolic features of x :

word _{i}	suf	tag _{$i-1$}	is-case(w_i)
race	ace	TO	TRUE

- Class-dependent indicator features of x :

	word _{i} =race	suf=ing	suf=ace	tag _{$i-1$} =TO	tag _{$i-1$} =DT	is-lower(w_i)=TRUE	...
JJ	0	0	0	0	0	0	
VB	1	0	1	1	0	1	
NN	0	0	0	0	0	0	
...							

Complex features

- Consider **proper nouns**: a word is likely to be a proper noun if it starts with a capital letter **and is not** the first word of the sentence.
- We need a conjunction of two primitive features

$$f_{123}(c, x) = \begin{cases} 1 & \text{if } word_{i-1} \neq \langle s \rangle \ \& \ \text{capitalized}(w_i) \ \& \ c = \text{NNP} \\ 0 & \text{otherwise} \end{cases}$$

- Digression: some ML methods (e.g. SVMs) use kernels which can capture feature interactions
 - ▶ If model can be rewritten in terms of dot products, those can be replaced by kernel functions
 - ▶ Using a kernel is equivalent to mapping the input vectors to higher dimensional feature representation
 - ▶ E.g. quadratic kernel $K(\mathbf{x}, \mathbf{z}) = \langle \mathbf{x} \cdot \mathbf{z} \rangle^2$ maps to:

$$\Phi(\mathbf{x}) = (x_1x_1, \dots, x_1x_n, x_2x_1, \dots, x_2x_n, \dots, x_nx_1, \dots, x_nx_n)$$

- MaxEnt does not do that!

Entia non sunt multiplicanda praeter necessitatem



ockham wIELDING razor

Entropy

- Out of all possible models, choose the simplest one consistent with the data (Occam's razor)
- Entropy of the distribution of discrete random variable X :

$$H(X) = - \sum_x P(X = x) \log_2 P(X = x)$$

- The uniform distribution has the **highest entropy**
- Finding the maximum entropy distribution in the set C of possible distributions

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

- Berger et al. (1996) showed that solving this optimization problem is equivalent to finding the multinomial logistic regression model whose weights maximize the likelihood of the training data.

Outline

- 1 Classification and Sequence labeling
- 2 Linear and Logistic Regression
- 3 Maximum Entropy Modeling
- 4 Maximum Entropy Markov Models**

HMMs and MEMMs

- HMM POS tagging model:

$$\hat{T} = \operatorname{argmax}_T P(T|W) \quad (10)$$

$$= \operatorname{argmax}_T P(W|T)P(T) \quad (11)$$

$$= \operatorname{argmax}_T \prod_i P(\text{word}_i|\text{tag}_i) \prod_i P(\text{tag}_i|\text{tag}_{i-1}) \quad (12)$$

- MEMM POS tagging model:

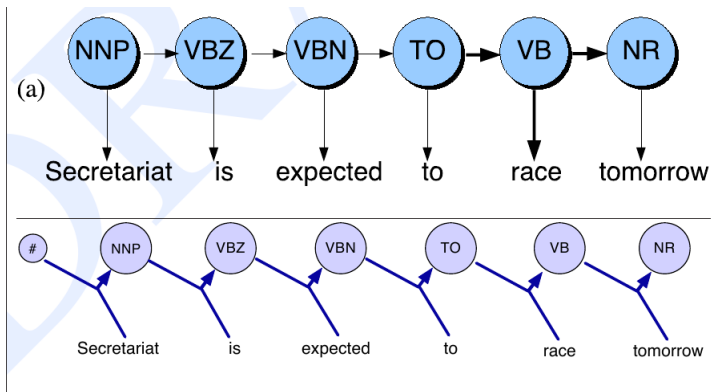
$$\hat{T} = \operatorname{argmax}_T P(T|W) \quad (13)$$

$$= \operatorname{argmax}_T \prod_i P(\text{tag}_i|\text{word}_i, \text{tag}_{i-1}) \quad (14)$$

- Maximum entropy model gives conditional probabilities



Conditioning probabilities in a HMM and a MEMM



Viterbi in MEMMs

- Decoding works almost the same as in HMM
- Except entries in the DP table are values of $P(t_j|t_{j-1}, word_j)$
- Recursive step: Viterbi value of time t for state j :

$$v_t(j) = \max_{i=1}^N v_{t-1}(i)P(s_j|s_i, o_t) \quad 1 \leq j \leq N, 1 < t \leq T \quad (15)$$

More

- Accessible intro to MaxEnt: A simple introduction to maximum entropy models for natural language processing, Ratnaparkhi (1997)
- More complete: A maximum entropy approach to natural language processing, Berger et al. (1996) in CL
- MEMM paper: Maximum Entropy Markov Models for Information Extraction and Segmentation, McCallum (2000)
- Other sequence labeling methods
 - ▶ Conditional Random Fields (also based on MaxEnt): Lafferty et al (2001)
 - ▶ Max-Margin Markov Networks, Taskar et al. (2003)
 - ▶ SVM^{struct} Tsochantaridis et al. (2005)