

Natural Language Processing

Lecture 13—10/6/2015

Jim Martin

Today

- **Multinomial Logistic Regression**
 - Aka log-linear models or maximum entropy (maxent)
 - Components of the model
 - Learning the parameters

10/1/15 Speech and Language Processing - Jurafsky and Martin 2

Logistic Regression Models

- Estimate $P(c|d)$ directly (discriminative model)
 - Features
 - Model w/ weights on features
 - Classification

10/1/15 Speech and Language Processing - Jurafsky and Martin 3

Features

- The kind of features used in NLP-oriented machine learning systems typically involve
 - Binary values
 - Think of a feature as being **on** or **off** rather than as a feature with a value
 - Values that are relative to an object/class pair rather than being a function of the object alone.
 - Typically have lots and lots of features (100,000s of features isn't unusual.)

10/1/15

Speech and Language Processing - Jurafsky and Martin

4

POS Features

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

$$f_4(c, x) = \begin{cases} 1 & \text{if } \text{is_lower_case}(\text{word}_i) \ \& \ c = \text{VB} \\ 0 & \text{otherwise} \end{cases}$$

10/1/15

Speech and Language Processing - Jurafsky and Martin

5

Sentiment Features

$$f_1(c, x) = \begin{cases} 1 & \text{if } \text{"great"} \in x \ \& \ c = + \\ 0 & \text{otherwise} \end{cases}$$

$$f_2(c, x) = \begin{cases} 1 & \text{if } \text{"second-rate"} \in x \ \& \ c = - \\ 0 & \text{otherwise} \end{cases}$$

$$f_3(c, x) = \begin{cases} 1 & \text{if } \text{"no"} \in x \ \& \ c = - \\ 0 & \text{otherwise} \end{cases}$$

$$f_4(c, x) = \begin{cases} 1 & \text{if } \text{"enjoy"} \in x \ \& \ c = - \\ 0 & \text{otherwise} \end{cases}$$

10/4/15

Speech and Language Processing - Jurafsky and Martin

6

Sentiment Features w/ Weights

1.9 $f_1(c, x) = \begin{cases} 1 & \text{if "great"} \in x \ \& \ c = + \\ 0 & \text{otherwise} \end{cases}$

.9 $f_2(c, x) = \begin{cases} 1 & \text{if "second-rate"} \in x \ \& \ c = - \\ 0 & \text{otherwise} \end{cases}$

.7 $f_3(c, x) = \begin{cases} 1 & \text{if "no"} \in x \ \& \ c = - \\ 0 & \text{otherwise} \end{cases}$

-0.8 $f_4(c, x) = \begin{cases} 1 & \text{if "enjoy"} \in x \ \& \ c = - \\ 0 & \text{otherwise} \end{cases}$

10/4/15 Speech and Language Processing - Jurafsky and Martin 7

Logistic Regression Model

$$p(c|x) = \frac{\exp\left(\sum_{i=1}^N w_i f_i(c, x)\right)}{\sum_{c' \in C} \exp\left(\sum_{i=1}^N w_i f_i(c', x)\right)}$$

10/1/15 Speech and Language Processing - Jurafsky and Martin 8

Logistic Regression Model

... there are virtually **no** surprises, and the writing is **second-rate**. So why did I **enjoy** it so much? For one thing, the cast is **great**.

Diagram illustrating weights for sentiment features:

- 0.9 (red dashed line) points to "no" (circled in red)
- 0.7 (red dashed line) points to "surprises" (circled in red)
- 0.8 (red dashed line) points to "second-rate" (circled in red)
- 1.9 (blue dashed line) points to "great" (circled in blue)
- 1.9 (blue dashed line) points to "enjoy" (circled in blue)

10/4/15 Speech and Language Processing - Jurafsky and Martin 9

Logistic Regression Model

- 0.7
- 0.8
- 0.9 ... there are virtually no surprises, and the writing is second-rate. So why did I enjoy it so much? For one thing, the cast is great 1.9 +

$$f_1(c, x) = \begin{cases} 1 & \text{if "great" } \in x \text{ \& } c = + \\ 0 & \text{otherwise} \end{cases}$$

$$f_2(c, x) = \begin{cases} 1 & \text{if "second-rate" } \in x \text{ \& } c = - \\ 0 & \text{otherwise} \end{cases}$$

$$f_3(c, x) = \begin{cases} 1 & \text{if "no" } \in x \text{ \& } c = - \\ 0 & \text{otherwise} \end{cases}$$

$$f_4(c, x) = \begin{cases} 1 & \text{if "enjoy" } \in x \text{ \& } c = + \\ 0 & \text{otherwise} \end{cases}$$

$$P(+|x) = \frac{e^{1.9}}{e^{1.9} + e^{0.9}e^{-0.7}e^{-0.8}} = .82$$

$$P(-|x) = \frac{e^{-0.9}e^{-0.7}e^{-0.8}}{e^{1.9} + e^{0.9}e^{-0.7}e^{-0.8}} = .18$$

10/4/15 Speech and Language Processing - Jurafsky and Martin 10

Argmax

- If we didn't really care about the probabilities then this

$$p(c|x) = \frac{\exp\left(\sum_{i=1}^N w_i f_i(c, x)\right)}{\sum_{c' \in C} \exp\left(\sum_{i=1}^N w_i f_i(c', x)\right)}$$
- Reduces to the just a comparison of the sum of the weights

$1.9 > .9 + .7 - .8$
 $1.9 > .8$

10/4/15 Speech and Language Processing - Jurafsky and Martin 11

Model

- So back to P(C|D)... given:
 - A training set of documents (D)
 - And an associated set of labels (C)
 - A finite set of features of the documents and classes
 - And a weight vector over the features (w)
 - This produces a proper probability distribution over the classes given an document

$$p(c|x) = \frac{\exp\left(\sum_{i=1}^N w_i f_i(c, x)\right)}{\sum_{c' \in C} \exp\left(\sum_{i=1}^N w_i f_i(c', x)\right)}$$

10/1/15 Speech and Language Processing - Jurafsky a 2

Learning the Weights

- We have a training set of labels and documents (y, x)
- We're going to use a Maximum Likelihood approach...
 - Choose the parameters (weights) that maximize the **probability** of the labels given the observations (features derived from the documents)
 - And we'll do that in log-space, so maximize the log prob of the labels given the data

10/4/15

Speech and Language Processing - Jurafsky and Martin

13

Learning the Weights

- So choose the weights that maximizes

$$L(w) = \sum_j \log P(y^{(j)} | x^{(j)})$$

10/4/15

Speech and Language Processing - Jurafsky and Martin

14

Learning the Weights

- So choose the weights that maximizes

$$= \log \sum_j \frac{\exp \left(\sum_{i=1}^N w_i f_i(y^{(j)}, x^{(j)}) \right)}{\sum_{y' \in Y} \exp \left(\sum_{i=1}^N w_i f_i(y'^{(j)}, x^{(j)}) \right)}$$

10/4/15

Speech and Language Processing - Jurafsky and Martin

15

Learning the Weights

- So choose the weights that maximizes the following in log space

$$\log \sum_j \exp \left(\sum_{i=1}^N w_i f_i(y^{(j)}, x^{(j)}) \right) - \log \sum_j \sum_{y' \in Y} \exp \left(\sum_{i=1}^N w_i f_i(y'^{(j)}, x^{(j)}) \right)$$

- Fortunately, that corresponds neatly to a convex optimization problem for which there are many possible approaches

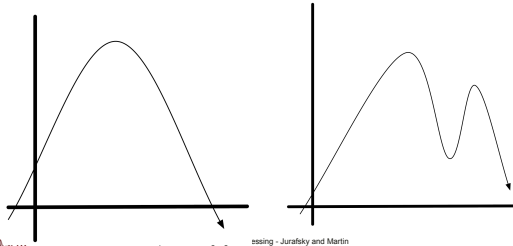
10/4/15

Speech and Language Processing - Jurafsky and Martin

16

Convex Problems

- What does convex/concave mean?
 - And why does it matter



10/4/15

Speech and Language Processing - Jurafsky and Martin

17

Derivatives

$$L'(w) = \sum_j f_k(y^{(j)}, x^{(j)}) - \sum_j \sum_{y' \in Y} P(y'|x^{(j)}) f_k(y'^{(j)}, x^{(j)})$$

$$L'(w) = \sum_j \text{Observed count}(f_k) - \text{Expected count}(f_k)$$

10/4/15

Speech and Language Processing - Jurafsky and Martin

18

Gradient Ascent

- Initialize the weights to zero
 - $w = 0$
- Loop until convergence
 - Calculate derivative for each feature
 - δ_k for each of the k features

$$L'(w) = \sum_j \text{Observed count}(f_k) - \text{Expected count}(f_k)$$

- Set $w = w + \beta * \delta$ for each feature k

10/4/15

Speech and Language Processing - Jurafsky and Martin

19

Optimization

- In practice, that can actually be slow to converge because the algorithm can either be taking steps
 - That are too small and hence take us too long to get where we're going
 - Or too large which leads us to overshoot the target and wander around too much
- Fortunately, you don't have to worry about this. Lots of packages available where you need to specify L and L' and you're done.
 - LBFGs

10/4/15

Speech and Language Processing - Jurafsky and Martin

20

Overfitting

- A problem with the approach as we've described it is that it is overly eager to match the training set as closely as it can
- With thousands or millions of features this can lead to poor performance on new data
- With logistic models this usually manifests itself as extremely large weights on a limited set of features

10/6/15

Speech and Language Processing - Jurafsky and Martin

21

Regularization

- Solution is to add a “penalty term” to the objective function.
- The job of the penalty term is to squash the weights of features that are getting out of control.

$$\hat{w} = \operatorname{argmax}_w \sum_j \log P(y^{(j)} | x^{(j)}) - \alpha R(w)$$

Parameter $\alpha > 0$ learned from held out dev data.

10/6/15

Speech and Language Processing - Jurafsky and Martin

22

Regularization

- Solution is to add a “penalty term” to the objective function.
- The job of the penalty term is to squash the weights of features that are getting out of control.

$$\hat{w} = \operatorname{argmax}_w \sum_j \log P(y^{(j)} | x^{(j)}) - \alpha R(w)$$

$$R(W) = \|W\|_2^2 = \sum_{j=1}^N w_j^2$$

10/6/15

23

Regularization and Learning

- With regularization, the parameter learning has to balance finding models that match the training data well and use small weights.

10/6/15

Speech and Language Processing - Jurafsky and Martin

24

Review

- Finite state methods
- Practical issues in segmentation and tokenization
- N-gram language models
- HMMs
 - HMMs applied to POS tagging
- Logistic regression and text classification

10/4/15

Speech and Language Processing - Jurafsky and Martin

25
