Lecture 6 Logistic Regression for Text Classification

CS 490A, Fall 2021 https://people.cs.umass.edu/~brenocon/cs490a_f21/

Laure Thompson and Brendan O'Connor

College of Information and Computer Sciences University of Massachusetts Amherst

[Many slides from Ari Kobren]

- HW1!
- Zoom OH demand?
- Accept exercises late? Yes, for check-minus: see the "Grading and Policies" page. Please let us know if you can't upload to Gradescope.

BOW linear model for text classif.

• Problem: classify doc d into one of k=1..K classes

• Parameters: For each class k, and word type w, there is a word weight

Prediction rule: choose class y with highest score

Have we seen a text classification model that can be seen as an instance of this??

 Representation: bag-of-words vector of doc d's word counts

Linear classification models

- The foundational model for machine learning-based NLP!
- Examples
 - The humble "keyword count" classifier (no ML)
 - Naive Bayes ("generative" ML)

- Today: Logistic Regression
 - "discriminative" ML
 - allows for features
 - used within more complex models (neural networks)



- Not just word counts. Anything that might be useful!
- <u>Feature engineering</u>: when you spend a lot of trying and testing new features. Very important!! This is a place to put linguistics in, or just common sense about your data.

Negation

Das, Sanjiv and Mike Chen. 2001. Yahoo! for Amazon: Extracting market sentiment from stock message boards. In Proceedings of the Asia Pacific Finance Association Annual Conference (APFA). Bo Pang, Lillian Lee, and Shivakumar Vaithyanathan. 2002. Thumbs up? Sentiment Classification using Machine Learning Techniques. EMNLP-2002, 79-86.

Add NOT to every word between negation and following punctuation:

didn't like this movie , but I didn't NOT like NOT this NOT movie but I

| Slide: <u>SLP3</u>|

Classification: LogReg (I) First, we'll discuss **how LogReg works**.

Then, why it's set up the way that it is.

Application: spam filtering

Classification: LogReg (I) • compute features (xs)

 $x_i = (\text{count "nigerian", count "prince", count "nigerian prince")}$

• given weights (betas)

 $\beta = (-1.0, -1.0, 4.0)$



Classification: LogReg (II)

• Compute the **dot product**



• Compute the logistic function for the label probability $P(z) = \frac{e^z}{e^z+1} = \frac{1}{1+e^{-z}}$

LogReg Exercise

features: (count "nigerian", count "prince", count "nigerian prince")



Classification: Dot Product



Why the logistic function?

 $P(z) = \frac{e^z}{e^z + 1} = \frac{1}{1 + e^z}$

$$\frac{1}{1+e^{-z}}$$

Multiclass Logistic Regression

• Binary logreg: let x be a feature vector for the doc, and y either 0 or 1

$$p(y = 1|x, \beta) = \frac{\exp(\beta^{\mathsf{T}}x)}{1 + \exp(\beta^{\mathsf{T}}x)}$$

Multiclass logistic regression: weight vector for each class

Prediction: dot product for each class

• Predicted probabilities: apply the softmax function to normalize

is a weight vector across the x features. ß

NB as Log-Linear Model What are the features in Naive Bayes?

What are the weights in Naive Bayes?

NB as Log-Linear Model

In both NB and LogReg we compute the dot product!

NB vs. LogReg Both compute the dot product

• NB: sum of log probs; LogReg: logistic fun.

Learning Weights NB: learn conditional probabilities separately via counting

LogReg: learn weights jointly

Learning Weights given: a set of feature vectors and labels

• goal: learn the weights.

Learning V $x_{00} = x_{01}$ $x_{10} \quad x_{11}$. . . · · · · $x_{n0} \quad x_{n1}$ n examples; xs - fea

Neights	
x_{0m}	y_0
x_{1m}	y_1
• • •	•
x_{nm}	y_n
tures: vs -	class

Learning Weights $P(z) = \frac{e^{z}}{e^{z} + 1} = \frac{1}{1 + e^{-z}}$ $g(z) = \frac{1}{1 + e^{-z}}$ $P(y = 1 \mid x) = g\left(\sum_{j=1}^{\infty} \beta_{j} x_{ij}\right)$ $P(z) = \frac{e^{z}}{e^{z} + 1} = \frac{1}{1 + e^{-z}}$

We know:



So let's try to maximize probability of the entire dataset - maximum likelihood estimation

Learning Weights So let's try to maximize probability of the entire dataset - maximum likelihood estimation

 $\beta^{MLE} = \arg\max_{\beta} \log P(y_0, \dots, y_n | \mathbf{x_0}, \dots, \mathbf{x_n}; \beta)$

Gradient ascent learning

• Follow direction of steepest ascent. Iterate





$$\ell(\beta) = \sum_{i=0}^{|X|} \log P(y_i | \mathbf{x_i}; \beta)$$

 ℓ : Training set log-likelihood η : Step size (a.k.a. learning rate)

 $\left(\frac{\partial \ell}{\partial \beta_1}, ..., \frac{\partial \ell}{\partial \beta_J} \right): \text{Gradient vector}$ (vector of per-element derivatives)

This is a generic optimization technique. Not specific to logistic regression! Finds the maximizer of any function where you can compute the gradient.

Pros & Cons

 LogReg doesn't assume independence better calibrated probabilities \bigcirc

NB is faster to train; less likely to overfit

NB & Log Reg • Both are linear models:

• Training is different: NB: weights trained independently \bigcirc LogReg: weights trained jointly



LogReg: Important Details! Visualizing decision boundary / bias term

- Overfitting / regularization
- Multiclass LogReg

You can use scikit-learn (python) to test it out!

Overfitting and generalization

- Overfitting: your model performs overly optimistically on training set, but generalizes poorly to other data (even from same distribution)
- To diagnose: separate training set vs. test set.
- Why might logreg (or NB) overfit to the training set?

Regularization to address overfitting

• How did we regularize Naive Bayes and language modeling?

• For logistic regression: L2 regularization for training

Regularization tradeoffs

No regularization <-----> Very strong regularization



Regularization

- Just like in language models, there's a danger of overfitting the training data. (For LM's, how did we combat this?)
- One method is <u>count thresholding</u>: throw out features that occur in < L documents (e.g. L=5). This is OK, and makes training faster, but not as good as....
- <u>Regularized logistic regression</u>: add a new term to penalize solutions with large weights. Controls the **bias/variance** tradeoff.

$$\beta^{\text{MLE}} = \arg \max_{\beta} \left[\log p(y_1 ... y_n) \right]$$
$$\beta^{\text{Regul}} = \arg \max_{\beta} \left[\log p(y_1 ... y_n) \right]$$

"Regularizer constant" Strength of penalty



Visualizing a classifier in feature space

Feature vector

Weights/parameters

50% prob where $\beta^{\mathsf{T}} x = 0$

Predict y=1 when $\beta^{\mathsf{T}} x > 0$

Predict y=0 when

 $\beta^{\mathsf{T}} x \leq 0$







x = (1, count "happy", count "hello", ...)