

Lecture 6

Logistic Regression for Text Classification

CS 490A, Fall 2021

https://people.cs.umass.edu/~brenocon/cs490a_f21/

9/21/21 Slides

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

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)



BOW linear model for text classif.

- Problem: classify doc d into one of $k=1..K$ classes
- Parameters: For each class $k \in 1..K$ and word type w , there is a word weight
- Representation: bag-of-words vector of doc d 's word counts

$\beta_{w,k} \in \mathbb{R}$
 "association b/w w and k "

$\vec{\beta}_k \in \mathbb{R}^V$

$x \in \mathbb{R}^{|V|}$
 $x_w = \#(\text{word } w \text{ in doc})$

- Prediction rule: choose class y with highest score
- $\forall k:$

$$\text{score}_k = \vec{\beta}_k^T \vec{x} = \sum_{w=1}^V \beta_{k,w} x_w$$

$$\hat{y} = \arg \max_{k \in 1..K} \{ \text{score}_k \}$$

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

Lexical Counting

$$n_{pos} = \#(\text{words in POS dict})$$

$$n_{neg} = \#(\text{words in NEG dict})$$

$$\hat{y} = \begin{cases} \text{pos} & \text{if } n_{pos} > n_{neg} \\ \text{neg} & \text{if } n_{pos} \leq n_{neg} \end{cases}$$

as Linear Model

$K=2$ classes

$$\beta_{pos, w} = \begin{cases} 1 & \text{if } w \in \text{POS dict} \\ 0 & \text{otherwise} \end{cases}$$

$$\beta_{neg, w} = \begin{cases} 1 & \text{if } w \in \text{NEG dict} \\ 0 & \text{o/w} \end{cases}$$

$\text{score}_{pos} = \text{count of pos words}$

Naive Bayes

$$\log p(y|x) \propto \left[\sum_{i=1}^{N_{\text{tok}}} \log p(w_i|y) \right] + \log p(y) \quad \left[\text{Token form} \right]$$

$$\propto \left[\sum_{w=1}^V \underbrace{x_w}_{\substack{\downarrow \\ \vec{x}}} \underbrace{\log p(w|y)}_{\beta_{w,y}} \right] + \underbrace{\log p(y)}_{\text{bias}_y}$$



- Input document **d** (a string...)
- Engineer a feature function, $f(d)$, to generate feature vector **x**

$f(d)$ \longrightarrow **x**

$f(d) =$ $\left(\begin{array}{l} \text{Count of "happy"}, \\ \text{(Count of "happy") / (Length of doc),} \\ \log(1 + \text{count of "happy"}), \\ \text{Count of "not happy"}, \\ \text{Count of words in my pre-specified word list,} \\ \text{"positive words according to my favorite} \\ \text{psychological theory"}, \\ \text{Count of several different "happy" emoticons} \\ \text{Count of "of the"}, \\ \text{Length of document,} \\ \dots \end{array} \right)$

Typically these use feature templates:
Generate many features at once

for each word w :

- $\${w}_count$
- $\${w}_is_count_nonzero$
- $NOT_ \${w}_bigram_count$

- Not just word counts. Anything that might be useful!
- Feature engineering: 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.

- Speed & size? Lang?

- Punct.!

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: SLP3]

Classification: ^{Binary}LogReg (I)

- compute **features** (xs)

$$\mathcal{X}_i = (\text{count } \underline{\text{"nigerian"}}, \text{count } \underline{\text{"prince"}}, \text{count } \underline{\text{"nigerian prince"}})$$

- given **weights** (betas)

$$\beta = (\underline{-1.0}, \underline{-1.0}, \underline{4.0})$$

$$N_{\text{feat}} = 3$$

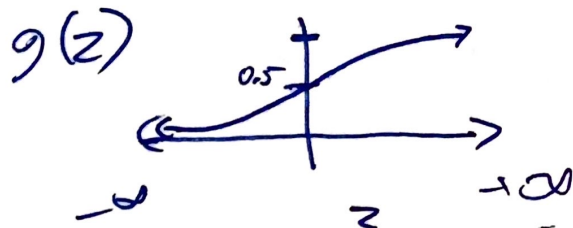
Classification: LogReg (II)

- Compute the **dot product**

$$z = \beta^T x$$

- Compute the **logistic function** for the label probability

$$P(y=1|x) = g(z) = \frac{1}{1 + e^{-z}} = \frac{e^z}{1 + e^z}$$



Logistic function

LogReg Exercise

features: (count “nigerian”, count “prince”, count “nigerian prince”)

$$\mathbf{x} = (1, 1, 1)$$

$$\boldsymbol{\beta} = (-1.0, -1.0, 4.0)$$

$$z = \boldsymbol{\beta}^T \mathbf{x} = -1 - 1 + 4 = 2$$

$$\underline{P(y=1 \mid \mathbf{x})} = \frac{1}{1+e^{-z}} = \frac{1}{1+e^{-2}} = .88$$

Learning Weights

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

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

$$= \arg \max_{\beta} \sum_{d=1}^{N_{doc}} \underbrace{\log P(y_d | x_d; \beta)}_{\log \sigma(\beta^T x_d)}$$

Learning Weights

We know:

$$g(z) = \frac{1}{1 + e^{-z}} \quad \underbrace{P(y = 1 \mid x)}_{\tau} = g \left(\underbrace{\sum_{j=1}^{\text{Nfeat}} \beta_j x_{ij}} \right)$$

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

Learning Weights

$$\begin{pmatrix} x_{00} & x_{01} & \dots & x_{0m} \\ x_{10} & x_{11} & \dots & x_{1m} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n0} & x_{n1} & \dots & x_{nm} \end{pmatrix} \begin{matrix} y_0 \\ y_1 \\ \vdots \\ y_n \end{matrix}$$

The diagram illustrates a dataset matrix where the first four columns represent features (x_{ij}) and the fifth column represents the class (y_i). A blue box highlights the first row of features ($x_{10}, x_{11}, \dots, x_{1m}$), and an arrow points from this box to the corresponding class value y_1 , indicating the learning process for a specific example.

n examples; xs - features; ys - class

Multiclass Logistic Regression

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

$$p(y = 1|x, \beta) = \frac{\exp(\beta^T x)}{1 + \exp(\beta^T x)} \quad \beta \text{ is a weight vector across the } x \text{ features.}$$

- Multiclass logistic regression: weight vector for each class

$$\vec{\beta}_k = [\beta_{k,1}, \beta_{k,2}, \dots, \beta_{k,v}]$$

- Prediction: dot product for each class

$$z_k = \vec{\beta}_k^T \vec{x} \in \mathbb{R} \quad [z_1, z_2, \dots, z_K]$$

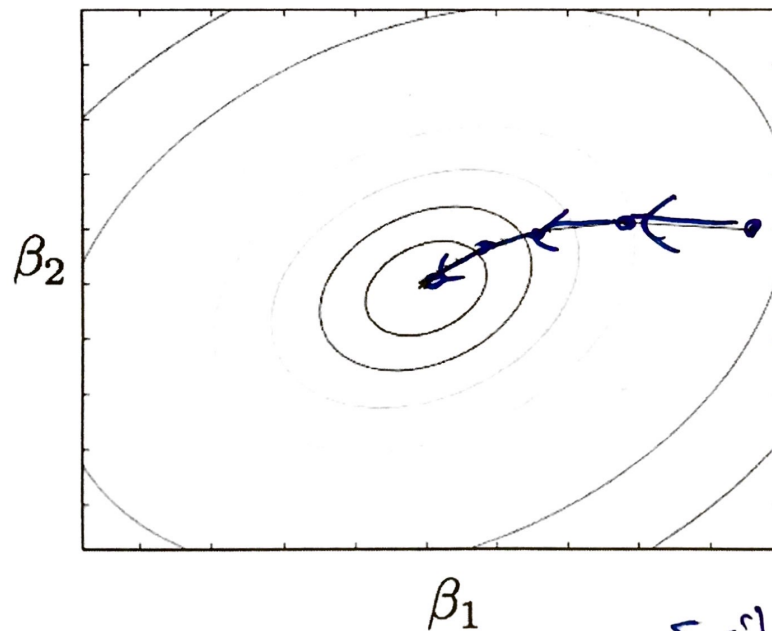
- Predicted probabilities: apply the *softmax function* to normalize

$$P(y=k|x) = P_k = \frac{e^{z_k}}{\sum_{j=1}^K e^{z_j}} \quad [e^{z_1}, e^{z_2}, \dots, e^{z_K}]$$

Gradient ascent learning

- Follow direction of steepest ascent. Iterate

$$\beta^{(new)} = \beta^{(old)} + \eta \frac{\partial \ell}{\partial \beta}$$



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

ℓ : Training set log-likelihood

η : Step size (a.k.a. learning rate)

$\left(\frac{\partial \ell}{\partial \beta_1}, \dots, \frac{\partial \ell}{\partial \beta_J} \right)$: 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.

Scikit-learn "Logistic Regression"

Name: _____

Exercise 3, in-class 9/16/21 – UMass CS 490A

Naïve Bayes example (from J&M exercise 4.2)

Turn in via Gradescope after class or by next Monday

We are given the following short movie reviews, each labeled with a genre (comedy or action):

review	label
fun, couple, love, love	comedy
fast, furious, shoot	action
couple, fly, fast, fun, fun	comedy
furious, shoot, shoot, fun	action
fly, fast, shoot, love	action

Now, we are given a new review:

fast, couple, shoot, fly

Compute the most likely class for this review. Assume a naive Bayes classifier and use add-1 smoothing for the likelihoods.

New V!

$$V=7$$

fun
couple

love

fast

furious

shoot

fly

Name: _____

Exercise 3, in-class 9/16/21 – UMass CS 490A

Naïve Bayes example (from J&M exercise 4.2)

Turn in via Gradescope after class or by next Monday

We are given the following short movie reviews, each labeled with a genre (comedy or action):

review	label
fun, couple, love, love	comedy
fast, furious, shoot	action
couple, fly, fast, fun, fun	comedy
furious, shoot, shoot, fun	action
fly, fast, shoot, love	action

Now, we are given a new review:

fast, couple, shoot, fly

= x

Compute the most likely class for this review. Assume a naive Bayes classifier and use add-1 smoothing for the likelihoods.

$$P(y=c/x) = P(c) \cdot P(\text{fast}|c) \cdot P(\text{couple}|c) \cdot P(\text{sh.}|c) \cdot P(\text{fly}|c)$$

$$= \frac{2}{5} \cdot \frac{1+1}{9+V} \cdot \frac{2+1}{9+V} \cdot \dots$$

$$P(\text{action}|x) = \frac{3}{5} \cdot \dots$$