

1. Naive Bayes

2. Logistic Regression Classifiers

CS 490A, Fall 2020

Applications of Natural Language Processing
https://people.cs.umass.edu/~brenocon/cs490a_f20/

Brendan O'Connor

College of Information and Computer Sciences
University of Massachusetts Amherst

[Many slides from Ari Kobren]

Logistic regression

- a.k.a. “Maximum Entropy” a.k.a. “Log-Linear Classifier”
- Kinda like Naive Bayes, but:
 - No Bayes Rule. Just directly define a function for $p(y | \text{text})$. *→ Discriminative model*
 - Learn all feature weights jointly
 - Correlated features don’t “double count”
 - Better to “throw in lots of features”
 - Discriminative training: optimize for label predictions
 - $p(y | \text{text})$, not $p(y, \text{text})$
- Tends to work a bit better than NB - best baseline method
- Good off-the-shelf implementations: use scikit-learn
- Requires regularization (similar to NB pseudocount smoothing)
- Advanced: integrate with neural networks

Features!

Features!

Features!

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

$f(d) \longrightarrow x$

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

Typically these use feature templates:
Generate many features at once

for each word w :

- $\{w\}_{\text{count}}$
- $\{w\}_{\log_1_plus_count}$
- $\{w\}_{\text{with_NOT_before_it_count}}$

- 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.

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

- given **weights** (betas)

β = (-1.0, -1.0, 4.0)

Params
of the model

bigram

Classification: LogReg (II)

- Compute the dot product

$$z = x_i^T \beta = \sum_{j=1}^{N_{\text{feat}}} \beta_j x_{i,j}$$

$z \rightarrow$ high if $y=1$ is likely

- Compute the logistic function for the label probability

$$P(y=1 | x) = \downarrow g(z)$$

$$g(z): \mathbb{R} \rightarrow [0, 1]$$

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

LogReg Exercise

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

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

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

$$P(y=1 | \mathbf{x}) = g(x_1\beta_1 + x_2\beta_2 + x_3\beta_3) = g(-1 - 1 + 4) = g(2) = \frac{1}{1 + e^{-2}} \approx 0.88$$

Classification: Dot Product

"score"
~~z~~

$$z = \sum_{j=1}^{\text{Nfeat}} \beta_j x_{ij}$$

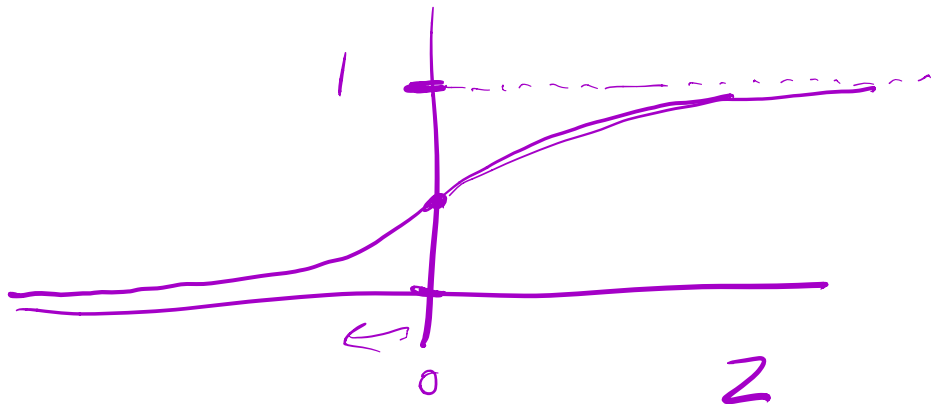
II
weighted sum of feature values

linear function

Why the logistic function?

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

$$z=0 \quad g = \frac{1}{1+e^0} = \frac{1}{2}$$



"sigmoid"

S-shape

outputs a prob

NB as Log-Linear Model

- What are the **features** in Naive Bayes?

word counts $\vec{x} = (n_{\text{dog}}, n_{\text{cat}}, \dots, n_{\text{zebra}})$

- What are the **weights** in Naive Bayes?

↪ word assoc w/ label

$$\beta_w^{\text{NB}} = \log \frac{p(w|y=1)}{p(w|y=0)}$$




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**

Training Time!

Learning Weights

- given: a set of **feature vectors** and **labels**

$\vec{x}_1, \vec{x}_2, \dots, \vec{x}_n$ y_1, y_2, \dots, y_n

- goal: learn the weights. $\Rightarrow \vec{\beta}$

Learning Weights

x_{00}

x_{10}

x_{00} x_{01} \dots x_{0m} y_0

x_{10} x_{11} \dots x_{1m} y_1

\vdots \vdots \ddots \vdots \vdots

x_{n0} x_{n1} \dots x_{nm} y_n

n examples; xs - features; ys - class

X

Y

Learning Weights

We know:

β weights

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

Choose β

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

Q: all $\vec{\beta} = \vec{0}$ $p(y=1/x; \beta) = \frac{1}{1+e^0}$ $z = \beta^T x = 0$

Learning Weights = 0.5

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

$$\beta^{MLE} = \arg \max_{\beta} \log P(y_0, \dots, y_n | \mathbf{x}_0, \dots, \mathbf{x}_n; \beta)$$

$$= \arg \max_{\beta} \sum_{i=1}^n \log P(y_i | \vec{x}_i; \beta)$$

↓
add std. label

(# "n...", # "pr", # "a p...")

$$B = (2, 3, -2) \Rightarrow 60\% \text{ acc}$$

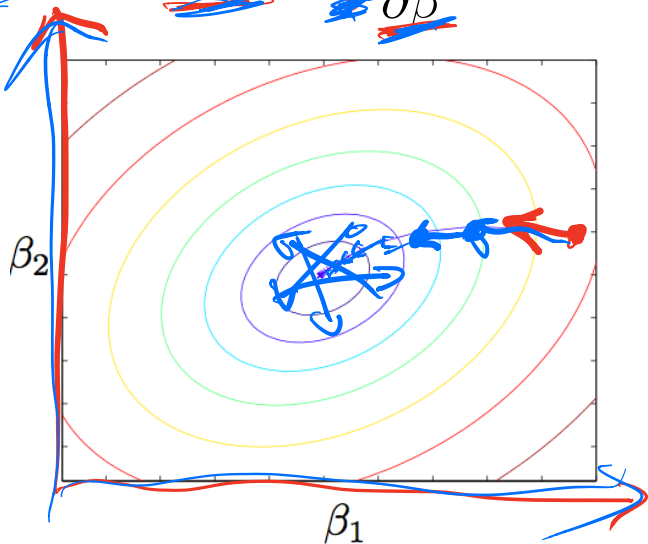
$$B = (1, 6, 3, .5) \Rightarrow 73\% \text{ acc}$$

$$B = (-1, -2, 6.1) \Rightarrow 78\% \text{ acc}$$

Gradient ascent learning

- Follow direction of *steepest ascent*. Iterate

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



weights \downarrow n \downarrow \downarrow fixed \downarrow changes during learning

$$\ell(\beta) = \sum_{i=1}^n \log P(y_i | \mathbf{x}_i; \beta)$$

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.

Pros & Cons

- LogReg doesn't assume independence *b/w feats*
 - better calibrated probabilities
~~~~~
⇒ allows for ~~negative~~ repetitive features
- NB is faster to train; less likely to overfit
~~~~~ ~~~~~~~~~

NB & Log Reg

- Both are linear models:

$$z = \sum_{i=0}^{|\mathbf{X}|} \beta_i x_i$$

- Training is different:
 - NB: weights trained independently
 - LogReg: weights trained jointly

LogReg: Important Details!

- Overfitting / regularization *"L2 penalty"*
- Visualizing decision boundary / bias term
- Multiclass LogReg *"Softmax fun":
multiclass sigmoid*

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

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 count thresholding: throw out features that occur in $< L$ documents (e.g. $L=5$). This is OK, and makes training faster, but not as good as....
- Regularized logistic regression: add a new term to penalize solutions with large weights. Controls the **bias/variance tradeoff**.

$$\beta^{\text{MLE}} = \arg \max_{\beta} [\log p(y_1..y_n | x_1..x_n, \beta)]$$
$$\beta^{\text{Regul}} = \arg \max_{\beta} \left[\log p(y_1..y_n | x_1..x_n, \beta) - \lambda \underbrace{\sum_j (\beta_j)^2}_{\text{"Quadratic penalty" or "L2 regularizer": Squared distance from origin}} \right]$$

“Regularizer constant”:
Strength of penalty

“Quadratic penalty”
or “L2 regularizer”:
Squared distance from origin

Visualizing a classifier in feature space

Feature vector $x = (1, \text{count "happy", count "hello", ...})$
Weights/parameters $\beta =$

"Bias term"
↓

50% prob where

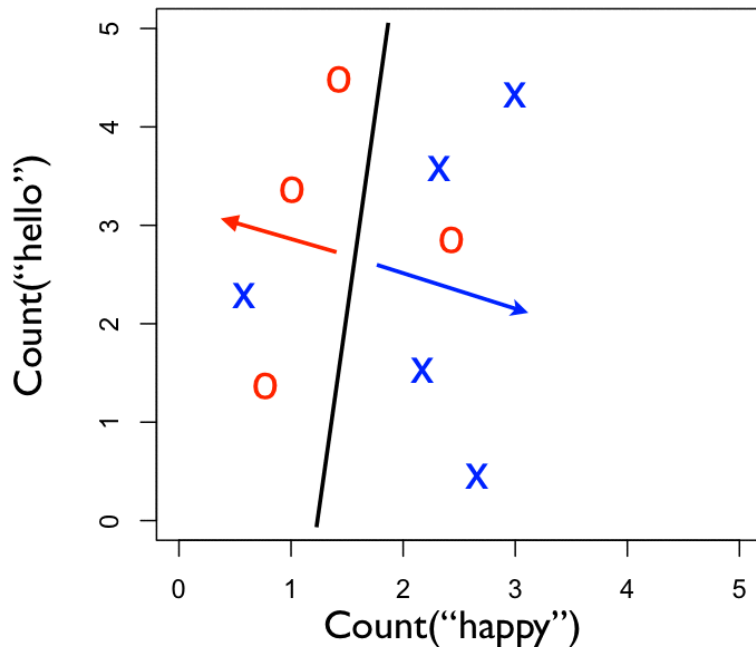
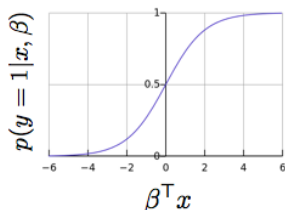
$$\beta^T x = 0$$

Predict $y=1$ when

$$\beta^T x > 0$$

Predict $y=0$ when

$$\beta^T x \leq 0$$



Binary vs. Multiclass

- 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: