

Logistic Regression for Text Classification

CS 485, Spring 2024
Applications of Natural Language Processing
https://people.cs.umass.edu/~brenocon/cs485_s24/

Brendan O'Connor
College of Information and Computer Sciences
University of Massachusetts Amherst

[With slides from Ari Kobren and SLP3]

BOW linear model for text classif.

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

$$\beta_{w,k} \in \mathbb{R}$$

$$\beta_{d_0, pos} = +3.0$$

- Prediction rule: choose class y with highest score

$$S_k = \sum_{w \in V} X_w \beta_{w,k}$$

$$\hat{y} = \underset{k=1..K}{\text{argmax}} S_k$$

$X_w = \#$ times " w "
occurs for X

"I love"

$$\left. \begin{array}{l} X_I = 1 \\ X_{love} = 0 \\ X_{char} = 0 \\ \vdots \end{array} \right\} |V|$$

Keyword count as linear model

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

$$B_{w,k} = \begin{cases} 1 & \text{if } w \text{ is in lexicon "k"} \\ 0 & \text{otherwise} \end{cases}$$

$$B_{\text{awesome}, \text{POS}} = 1$$

$$B_{\text{awesome}, \text{NEG}} = 0$$

- Prediction rule: choose class y with highest score

$$S_k = \sum_w B_{w,k} x_w \quad \text{= \# of times } w \text{ appears in doc } d \text{ in lexicon "k"}$$

$$S_{\text{POS}} = (\dots) \quad S_{\text{NEG}} = (\dots)$$

Naive Bayes as linear model

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

$$B_{w,k} = \log p(w/k)$$

$$X_w = \# \text{ instances of word "w" in doc}$$

- Prediction rule: choose class y with highest score

$$S_k = \sum_w x_w B_{w,k} = \sum_w x_w \log p(w/k)$$

$$\hat{y} = \underset{k}{\operatorname{argmax}} \log [p(y=k) p(x/y=k)]$$

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**
 - a linear classification model, trained to be good at *prediction*
 - allows for *features*
 - used within more complex models (neural networks)

Motivation: feature engineering

- For Naive Bayes, we used counts of each word in the vocabulary (BOW representation). But why not also use...
 - Number of words from "CS485 Crowdscore Positive Lexicon"
 - ...from "CS485 Crowdscore Negative Lexicon" ... or another...
 - Phrases?
 - Words/phrases with negation markers?
 - Number of "!" occurrences?
 - or...?

$y = x$
→ $y = NEG$



Var	Definition	Value in Fig. 5.2
x_1	count(positive lexicon words \in doc)	3
x_2	count(negative lexicon words \in doc)	2
x_3	$\begin{cases} 1 & \text{if "no" } \in \text{ doc} \\ 0 & \text{otherwise} \end{cases}$	1
x_4	count(1st and 2nd pronouns \in doc)	3
x_5	$\begin{cases} 1 & \text{if "!" } \in \text{ doc} \\ 0 & \text{otherwise} \end{cases}$	0
x_6	$\ln(\text{word count of doc})$	$\ln(66) = 4.19$

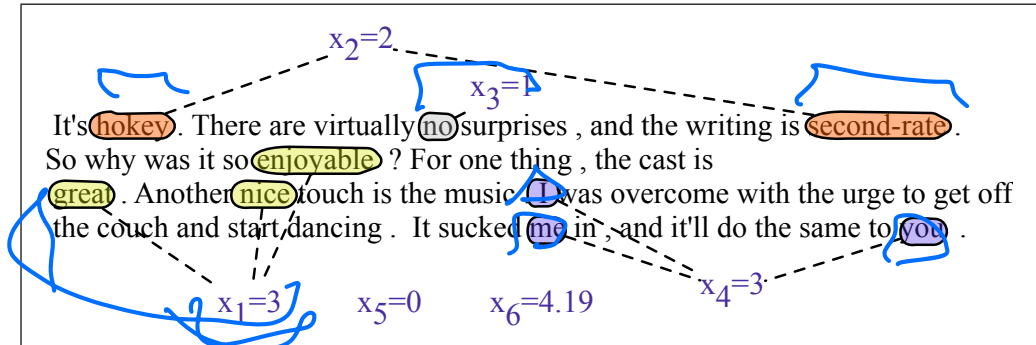


Figure 5.2 A sample mini test document showing the extracted features in the vector x .

- Logistic regression can accommodate **any arbitrary features**
- 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)

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

- given **weights** (betas)

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

Classification: LogReg (II)

- Compute the **dot product**

$$s = \vec{\beta} \cdot \vec{x} = \vec{\beta}^T \vec{x} = \langle \vec{\beta}, \vec{x} \rangle$$

"inner product"

$$= \sum_{j=1}^J \beta_j x_j$$

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

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

$$P(y=0|x) = 1 - P(y=1|x)$$

$y=0$
 $y=1$

LogReg Exercise

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

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

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

$$s = \beta^T \mathbf{x} = 1(-1) + 1(-1) + 1(4) = 2$$

$$\underline{P(y=1 | \mathbf{x})} = g(s) = \frac{1}{1 + e^{-s}} = \frac{1}{1 + e^{-2}} = .8807$$

Classification: Dot Product

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

Why the **logistic function**?

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

$$s = 0 \Rightarrow \frac{1}{1 + e^0} = \frac{1}{2}$$

$$P(y=1|x) \in [0, 1]$$

$$s \rightarrow \infty$$

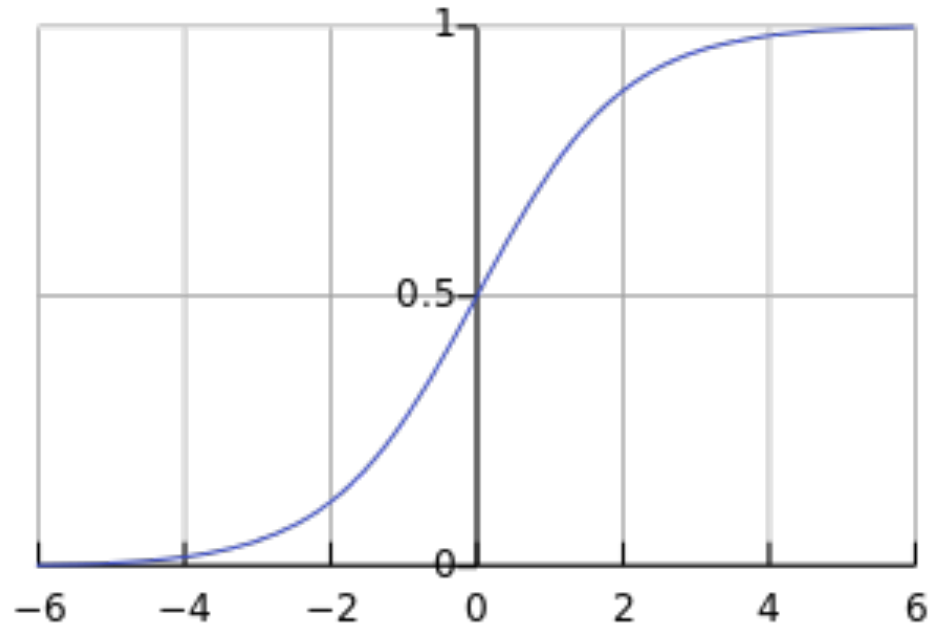
$$\frac{1}{1 + e^{-\infty}} \rightarrow \frac{1}{1 + 0} = 1$$

$$s \rightarrow -\infty$$

$$\frac{1}{1 + e^{\infty}} \rightarrow \frac{1}{\infty} = 0$$

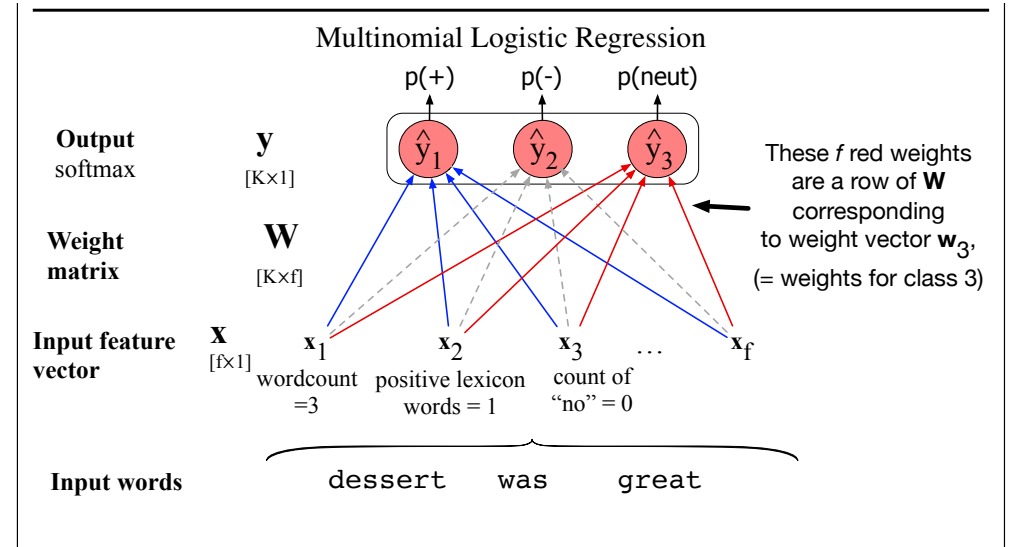
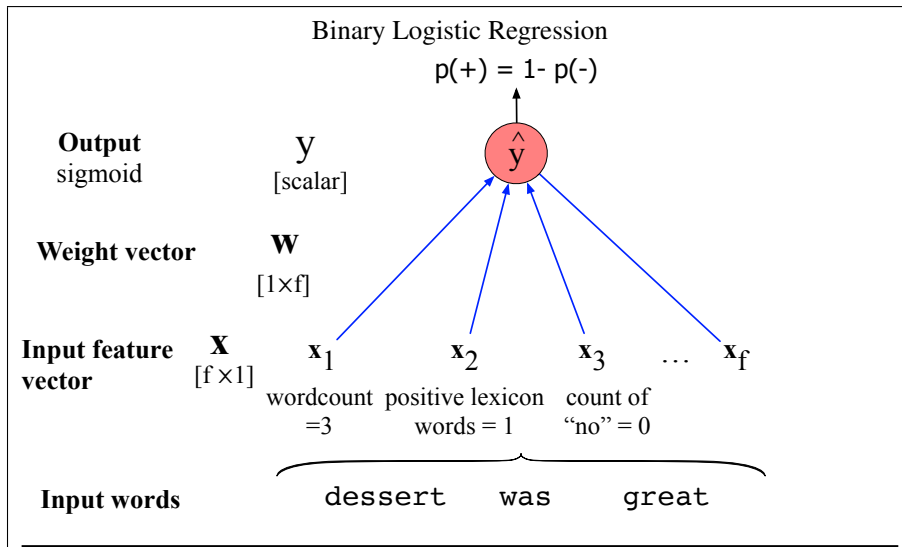
Logistic Function

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



Multiclass Logistic Regression

- Generalize to $K > 2$ classes
- Each class has its own weight vector (across all features; e.g. BOW counts)



Multiclass Logistic Regression

- Weight vector for each class

$$\vec{\beta}_k = [\beta_{k,1}, \dots, \beta_{k,f}] \quad \forall k \in 1..K$$

- Prediction: dot product for each class

$\forall k:$

$$s_k = (\vec{\beta}_k)^T X = \sum_{j=1}^f \beta_{k,j} X_j$$

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

$\forall k:$

$$\frac{e^{s_k}}{\sum_{l=1}^K e^{s_l}} = P(y=k | x)$$

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 Weights

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

Learning Weights

We know:

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

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

Learning Weights

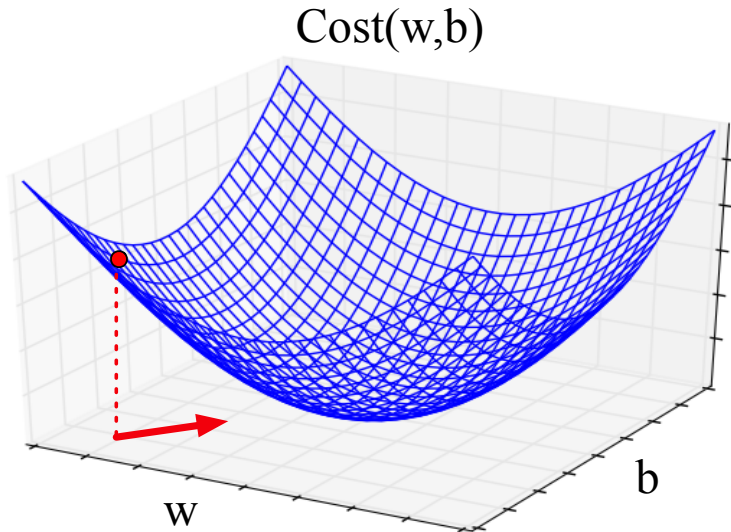
So let's try to maximize probability of the entire dataset - **maximum likelihood estimation**

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

Gradient ascent/descent learning

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

- Follow direction of *steepest ascent*. Iterate: $\beta^{(new)} = \beta^{(old)} + \eta \frac{\partial \ell}{\partial \beta}$



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

GD is a generic method for optimizing differentiable functions — widely used in machine learning!

Pros & Cons

- LogReg doesn't assume independence
 - better calibrated probabilities
- NB is faster to train; less likely to overfit

NB & Log Reg

- Both are linear models:

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

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

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.
- How did we regularize Naive Bayes and language modeling?

- For logistic regression: L2 regularization for training

Regularization tradeoffs

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

Visualizing a classifier in feature space

“Bias term”
↓
Feature vector $x = (1, \text{count “happy”, count “hello”, …})$
Weights/parameters $\beta =$

50% prob where

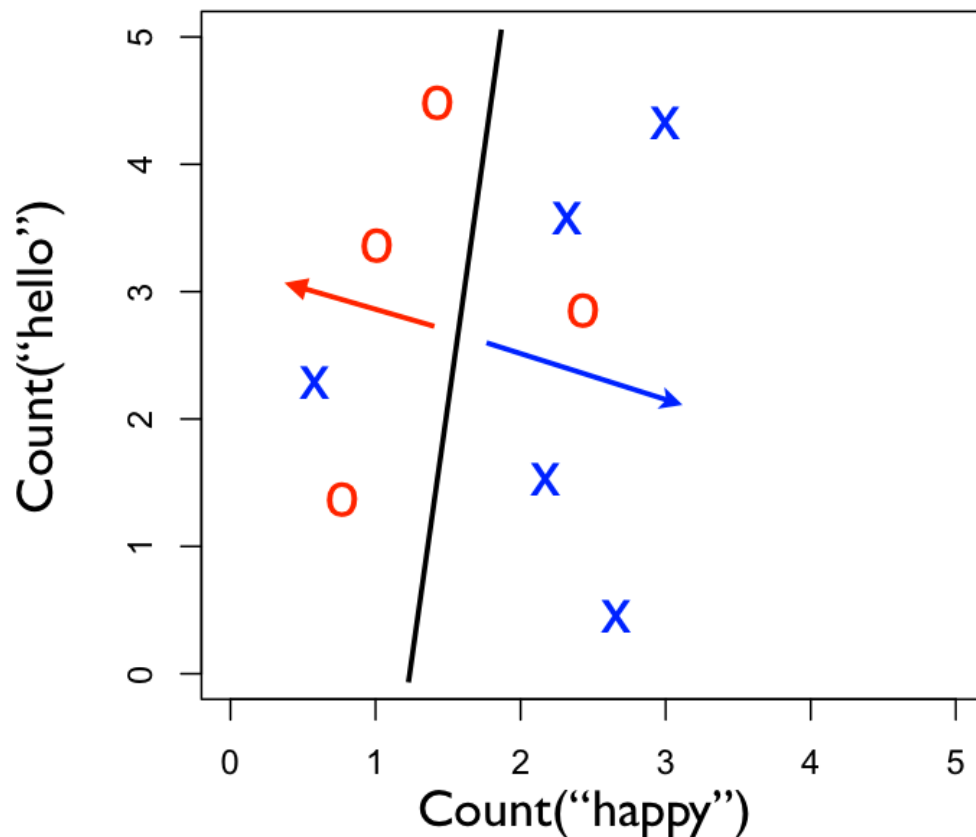
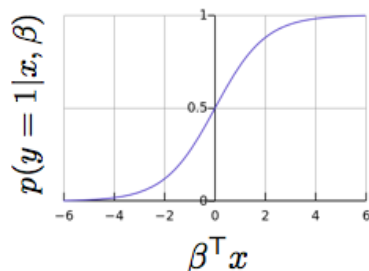
$$\beta^T x = 0$$

Predict $y=1$ when

$$\beta^T x > 0$$

Predict $y=0$ when

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



Logistic regression wrap-up

- Given you can extract features from your text, logistic regression is the best, easy-to-use, method
 - Logistic regression with BOW features is an excellent baseline method to try at first
 - Will be a foundation for more sophisticated models, later in course
- Always regularize your LR model
- We recommend using the implementation in scikit-learn
 - Useful: CountVectorizer to help make BOW count vectors
- Next: but where do the LABELS in supervised learning come from?