

Logistic regression in use

Topics we'll cover

- ① A text classification problem
- ② Bag-of-words representation for text
- ③ Solution by logistic regression
- ④ Margin versus test error
- ⑤ Interpreting the model

Sentiment data

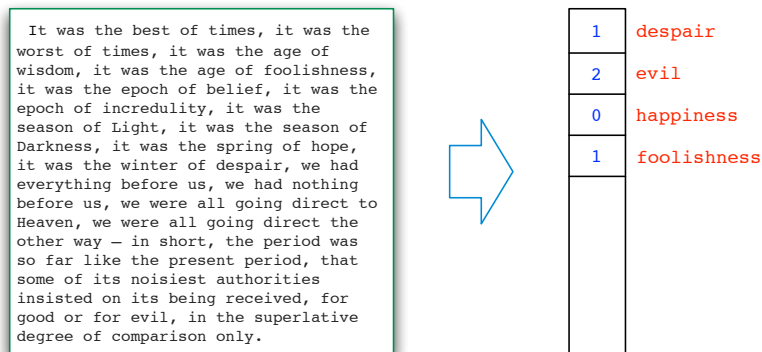
Data set: sentences from reviews on Amazon, Yelp, IMDB.
Each labeled as positive or negative.

- Needless to say, I wasted my money.
- He was very impressed when going from the original battery to the extended battery.
- I have to jiggle the plug to get it to line up right to get decent volume.
- Will order from them again!

2500 training sentences, 500 test sentences

Handling text data

Bag-of-words: vectorial representation of text sentences (or documents).



- Fix V = some vocabulary.
- Treat each sentence (or document) as a vector of length $|V|$:

$$x = (x_1, x_2, \dots, x_{|V|}),$$

where $x_i = \#$ of times the i th word appears in the sentence.

A logistic regression approach

Code positive as $+1$ and negative as -1 .

$$\Pr_{w,b}(y \mid x) = \frac{1}{1 + e^{-y(w \cdot x + b)}}$$

Given training data $(x^{(1)}, y^{(1)}), \dots, (x^{(n)}, y^{(n)}) \in \mathbb{R}^d \times \{-1, 1\}$, find w, b minimizing

$$L(w, b) = \sum_{i=1}^n \ln(1 + e^{-y^{(i)}(w \cdot x^{(i)} + b)})$$

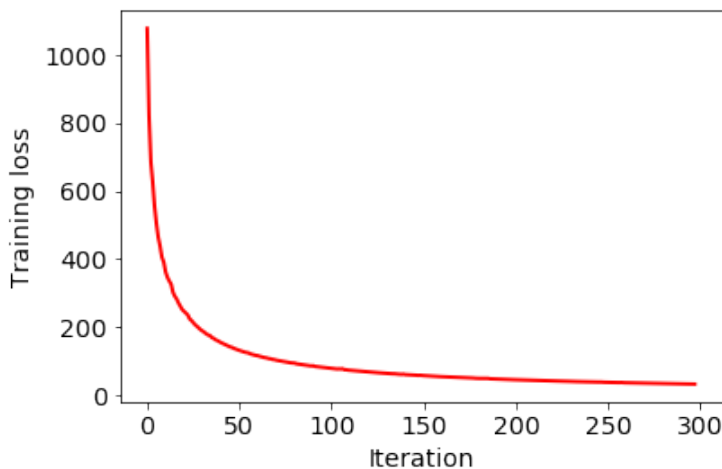
Convex problem with many solution methods, e.g.

- gradient descent, stochastic gradient descent
- Newton-Raphson, quasi-Newton

All converge to the optimal solution.

Local search in progress

Look at how loss function $L(w, b)$ changes over iterations of stochastic gradient descent.



Final model: **test error** 0.21.

Some of the mistakes

Not much dialogue, not much music, the whole film was shot as elaborately and aesthetically like a sculpture. 1

This film highlights the fundamental flaws of the legal process, that it's not about discovering guilt or innocence, but rather, is about who presents better in court. 1

You need two hands to operate the screen. This software interface is decade old and cannot compete with new software designs. -1

The last 15 minutes of movie are also not bad as well. 1

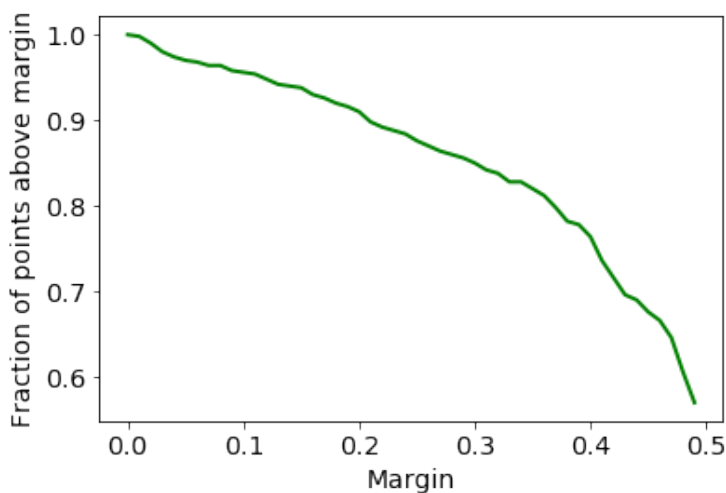
If you plan to use this in a car forget about it. -1

If you look for authentic Thai food, go else where. -1

Waste your money on this game. 1

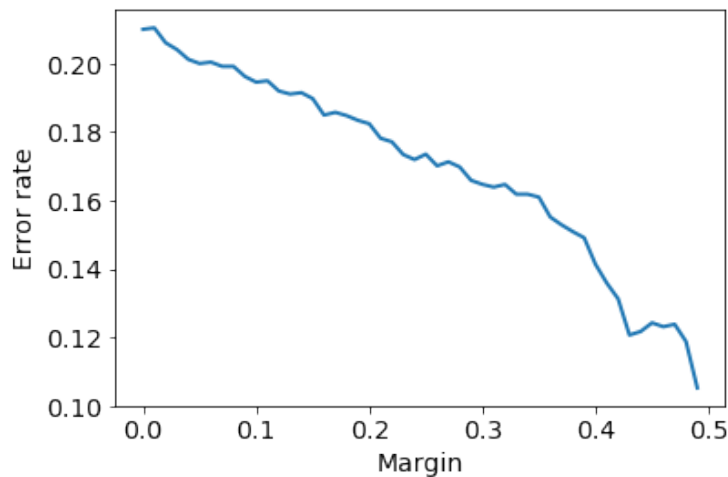
Margin and test error

$$\text{Margin on test pt } x = \left| \Pr_{w,b}(y = 1|x) - \frac{1}{2} \right|$$



Margin and test error

$$\text{Margin on test pt } x = \left| \Pr_{w,b}(y = 1|x) - \frac{1}{2} \right|$$



Interpreting the model

Words with the most positive coefficients

'sturdy', 'able', 'happy', 'disappoint', 'perfectly', 'remarkable', 'animation',
'recommendation', 'best', 'funny', 'restaurant', 'job', 'overly', 'cute', 'good', 'rocks',
'believable', 'brilliant', 'prompt', 'interesting', 'skimp', 'definitely', 'comfortable',
'amazing', 'tasty', 'wonderful', 'excellent', 'pleased', 'beautiful', 'fantastic',
'delicious', 'watch', 'soundtrack', 'predictable', 'nice', 'awesome', 'perfect', 'works',
'loved', 'enjoyed', 'love', 'great', 'happier', 'properly', 'liked', 'fun', 'screamy',
'masculine'

Words with the most negative coefficients

'disappointment', 'sucked', 'poor', 'aren', 'not', 'doesn', 'worst', 'average',
'garbage', 'bit', 'looking', 'avoid', 'roasted', 'broke', 'starter', 'disappointing', 'dont',
'waste', 'figure', 'why', 'sucks', 'slow', 'none', 'directing', 'stupid', 'lazy',
'unrecommended', 'unreliable', 'missing', 'awful', 'mad', 'hours', 'dirty', 'didn',
'probably', 'lame', 'sorry', 'horrible', 'fails', 'unfortunately', 'barking', 'bad', 'return',
'issues', 'rating', 'started', 'then', 'nothing', 'fair', 'pay'