Classifying with Regression

Reminders

- Midterm 1: this Friday
- My next help hours: Thursday 3-4pm
- CS Colloquium from 4-5pm: Elena Glassman
- Social Choice conference Oct. 14th-15th
- Over break: try to install PyTorch in your CS 333

Python environment



Midterm 1

- In-class programming midterm
- Bring your own laptop to work on
- I will have starter code and documentation for you to download at the beginning, then ask you to turn off wifi.
- At the end, you will submit your code on Gradescope.
- May bring a 4x6in note card

Word Embedding Exploration

Word Embedding Interactive Demo

https://www.cs.cmu.edu/~dst/ WordEmbeddingDemo/index.html

of Wad Embelding Techniques Summery Cons sparse bad long vectors storage collocation, frequency counts: interpretability easy to calc. complicated to learn word 2 vec: dense no interpretability SNOVY contextual embddings; Controls different even mure Word mesning

Logistic Regression

Generative and Discriminative Classifiers

Naive Bayes is a **generative** classifier.

Logistic regression is a discriminative classifier.

Generative and Discriminative Classifiers

Suppose we're distinguishing cat from dog images





me imagenet

Generative Classifier:

- Build a model of what's in a cat image
 - Knows about whiskers, ears, eyes
 - Assigns a probability to any image:
 - how cat-y is this image?





Also build a model for dog images

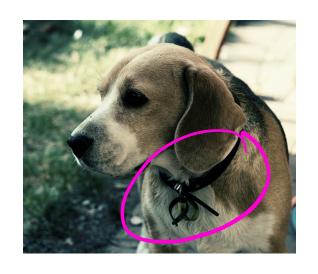
Now given a new image:

Run both models and see which one fits better

Discriminative Classifier

Just try to distinguish dogs from cats





Dogs have collars! Let's ignore everything else

Not easy to answer questions like: given that I observe a bone, how likely is it to be a dog?

But: discriminative classifiers generally perform better.

Finding the correct class c from a document d with Generative vs Discriminative Classifiers

Naive Bayes:
$$C_{MAP} = argmax P(dlc)P(c)$$

Logistic Regression:
$$C_{MAP} = argmax \qquad P(c \mid d)$$

$$c \in C$$

Logistic Regression Classifiers

Text Classification: definition

Input:

- a document x
- a fixed set of classes $C = \{c_1, c_2, ..., c_l\}$

Output: a predicted class $\hat{y} \in C$

Binary Classification in Logistic Regression

Given a series of input/output pairs:

$$(x^{(i)}, y^{(i)})$$

For each observation x⁽ⁱ⁾

- We represent $x^{(i)}$ by a **feature vector** $[x_1, x_2, ..., x_n]$
- We compute an output: a predicted class $\hat{y}^{(i)} \in \{0,1\}$

Features in logistic regression

For feature x_i , weight w_i tells is how important is x_i

- $x_i = "review contains 'awesome'": <math>w_i = +10$
- $x_j = "review contains 'abysmal'": <math>w_j = -10$
- x_k = "review contains 'mediocre'": $w_k = -2$

Logistic Regression for one observation x

Input observation: vector
$$x = [x_1, x_2, \dots x_n]$$

Weights: one per feature:
$$W = L_{w_1}, w_2 \cdots w_n J$$
Sometimes denoted $\Theta : \theta = L_{\theta_1}, \theta_2 \dots J$

Output: a predicted class 3 & £0,13

Moltinomial:
$$g \in \{0,1,\ldots,n\}$$

How to do classification

For each feature x_i, weight w_i tells us importance of x_i
• (Plus we'll have a bias b)

We'll sum up all the weighted features and the bias

$$Z = \left(\begin{array}{c} \tilde{Z} & W_{1} X_{1} \\ \vdots & 1 \end{array} \right) + b$$

$$Z = W \cdot X + b$$

$$I + Z \text{ is high, we predict } \mathcal{J} = 1$$

$$I + Z \text{ is law, we predict } \mathcal{J} = 0$$

But we want a probabilistic classifier

We need to formalize "sum is high".

We'd like a principled classifier that gives us a probability, just like Naive Bayes did

We want a model that can tell us:

$$p(y=1|x;W)$$

$$p(y=0|x;W)$$

$$p(y=0|x;W)$$

$$p(y=1|x;w)$$

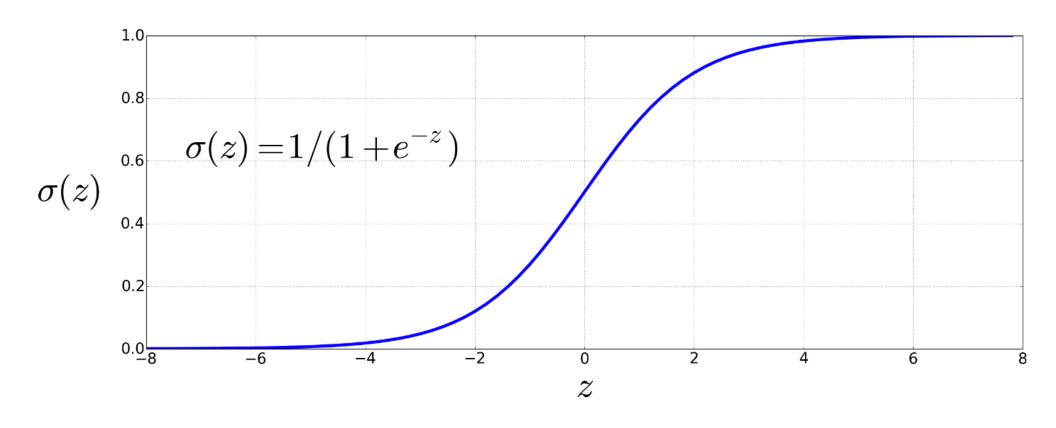
The problem: z isn't a probability, it's just a number!

$$z = w \cdot x + b$$

Solution: use a function of z that goes from 0 to 1

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

The very useful sigmoid function



Idea of logistic regression

Compute w·x+b, then pass it through the sigmoid function: $\sigma(w\cdot x+b)$.

Treat the result as a probability.

Making probabilities with sigmoids

$$P(y=1) = \sigma(w \cdot x + b)$$

$$= \frac{1}{1 + \exp(-(w \cdot x + b))}$$

$$P(y=0) = 1 - \sigma(w \cdot x + b)$$

$$= 1 - \frac{1}{1 + \exp(-(w \cdot x + b))}$$

$$= \exp(-(w \cdot x + b))$$

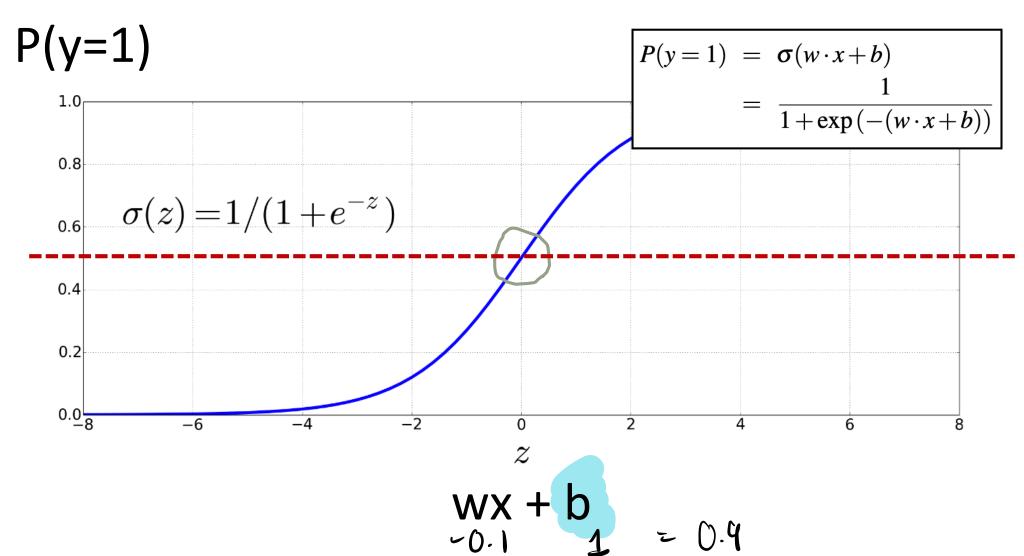
$$= \exp(-(w \cdot x + b))$$

Turning a probability into a classifier

$$decision(x) = \begin{cases} 1 & \text{if } P(y=1|x) > 0.5 \\ 0 & \text{otherwise} \end{cases}$$

0.5 here is called the **decision boundary**

The probabilistic classifier



The bias term b can be used to shift the intercept!

Turning a probability into a classifier

decision(x) =
$$\begin{cases} 1 & \text{if } P(y=1|x) > 0.5 \\ 0 & \text{otherwise} \end{cases}$$
 if $\mathbf{W} \cdot \mathbf{x} + \mathbf{b} > 0$

The two phases of logistic regression

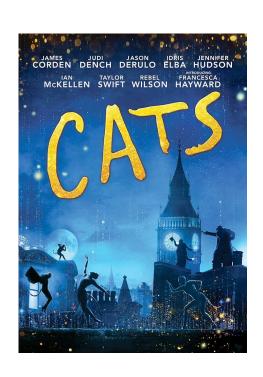
Training: we learn weights *w* and *b* using **stochastic gradient descent** and **cross-entropy loss**.

Test: Given a test example x we compute p(y|x) using learned weights w and b, and return whichever label (y = 1 or y = 0) is higher probability

Logistic Regression Example: Text Classification

Sentiment example: does y=1 or y=0?

CATS was a marvelous disaster, with witty charm and emotion throughout, cheeky charisma and crying no doubt... I personally went in expecting the worst movie I had ever seen - and it was far more awful and disappointing that I expected.

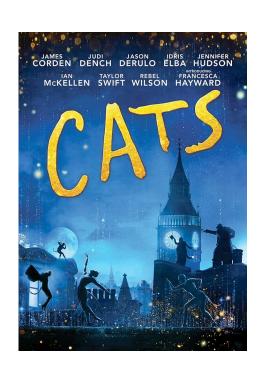




Sentiment example: does y=1 or y=0?



CATS was a marvelous disaster, with witty charm and emotion throughout, cheeky charisma and crying no doubt... I personally went in expecting the worst movie I had ever seen - and it was far more awful and disappointing that I expected.





Features

<u>Var</u>	Definition	Value ∈ doc
\mathbf{x}_1	count(positive lexicon) ∈ doc	6
X 2	count(negative lexicon) ∈ do	c 4
X 3	1 if "no"∈ doc else 0	1
X 4	1 if "!"∈ doc else 0	0
X 5	log(word count of doc)	1n(42) = 3.73
	X = [6 4 1 0	3.73]

Weights

Var	Definition	Value ∈ doc	W
\mathbf{x}_1	count(positive lexicon) ∈ doc	6	2.5
X ₂	count(negative lexicon) ∈ doc	4	-5
X 3	1 if "no"∈ doc else 0	1	- I· Z
X4	1 if "!"∈ doc else 0	0	2
X 5	log(word count of doc)	3.73	0.7

Bias
$$b = 0.$$

Classifying sentiment for input x

Classification in (binary) logistic regression: summary

Given:

- a set of classes: (+ sentiment,- sentiment)
- $^{\circ}$ a vector **x** of features [x_1 , x_2 , ..., x_n]
 - x1= count("awesome")
 - x2 = log(number of words in review)
- A vector **w** of weights [w₁, w₂, ..., w_n]
 - ° w_i for each feature f_i

$$P(y=1) = \sigma(w \cdot x + b)$$

$$= \frac{1}{1 + \exp(-(w \cdot x + b))}$$

Feature Representations

Example task: period disambiguation

This ends in a period.
The house at 465 Main St. is new.

$$x_1 = \begin{cases} 1 & \text{if "}Case(w_i) = \text{Lower"} \\ 0 & \text{otherwise} \end{cases}$$
 $x_2 = \begin{cases} 1 & \text{if "}w_i \in \text{AcronymDict"} \\ 0 & \text{otherwise} \end{cases}$
 $x_3 = \begin{cases} 1 & \text{if "}w_i = \text{St. \& }Case(w_{i-1}) = \text{Cap"} \\ 0 & \text{otherwise} \end{cases}$

Multi-class Regression

Multinomial Logistic Regression

Often we need more than 2 classes

- Positive/negative/neutral
- Parts of speech (noun, verb, adjective, adverb, preposition, etc.)
- Classify emergency SMSs into different actionable classes

If >2 classes we use multinomial logistic regression

- = Softmax regression
- = Multinomial logit
- = Maximum entropy modeling or MaxEnt

So "logistic regression" means binary (2 classes)

Multinomial Logistic Regression

The probability of everything must still sum to 1

```
P(positive|doc) + P(negative|doc) + P(neutral|doc) = 1
```

Need a generalization of the sigmoid called **softmax**

- Takes a vector $z = [z_1, z_2, ..., z_k]$ of k arbitrary values
- Outputs a probability distribution

Turns a vector $z = [z_1, z_2, ..., z_k]$ of k arbitrary values into probabilities

$$softmax(z_i) = \frac{e \times \rho(z_i)}{\sum_{j=1}^{K} e \times \rho(z_j)} | \leq i \leq k$$

softmax (z) =
$$\left[\frac{\exp(z_i)}{\sum_{i=1}^{k} \exp(z_i)} - \cdots \right] \frac{\exp(z_k)}{\sum_{i=1}^{k} \exp(z_i)}$$

Turns a vector $z = [z_1, z_2, ..., z_k]$ of k arbitrary values into probabilities

$$\operatorname{softmax}(z_i) = \frac{\exp(z_i)}{\sum_{j=1}^k \exp(z_j)} \quad 1 \le i \le k$$

The denominator $\sum_{i=1}^{K} \exp(\mathbf{z}_i)$ is used to normalize all the values into probabilities.

softmax(z) =
$$\left[\frac{\exp(z_1)}{\sum_{i=1}^k \exp(z_i)}, \frac{\exp(z_2)}{\sum_{i=1}^k \exp(z_i)}, ..., \frac{\exp(z_k)}{\sum_{i=1}^k \exp(z_i)}\right]$$

Turns a vector $z = [z_1, z_2, ..., z_k]$ of k arbitrary values into probabilities :

$$z = [0.6, 1.1, -1.5, 1.2, 3.2, -1.1]$$

softmax(z) =

Turns a vector $z = [z_1, z_2, ..., z_k]$ of k arbitrary values into probabilities :

$$z = [0.6, 1.1, -1.5, 1.2, 3.2, -1.1]$$

softmax(z) =
$$\left[\frac{\exp(z_1)}{\sum_{i=1}^k \exp(z_i)}, \frac{\exp(z_2)}{\sum_{i=1}^k \exp(z_i)}, ..., \frac{\exp(z_k)}{\sum_{i=1}^k \exp(z_i)}\right]$$

[0.055, 0.090, 0.006, 0.099, 0.74, 0.010]

Softmax in multinomial logistic regression

$$p(y = c|x) = \frac{\exp(w_c \cdot x + b_c)}{\sum_{j=1}^{K} \exp(w_j \cdot x + b_j)}$$

Input is still the dot product between weight vector *w* and input vector *x*, but now we need separate weight vectors for each of the *K* classes.

Features in binary versus multinomial logistic regression

Binary: positive weight

$$x_5 = \begin{cases} 1 & \text{if "!"} \in \text{doc} \\ 0 & \text{otherwise} \end{cases} \quad W_5 = 3.0$$

Multinominal: separate weights for each class:

Feature	Definition	$w_{5,+}$	$w_{5,-}$	$w_{5,0}$
$f_5(x)$	$\begin{cases} 1 & \text{if "!"} \in \text{doc} \\ 0 & \text{otherwise} \end{cases}$	3.5	3.1	-5.3