L3: Text Classification

Spring 2022



COS 484

Assignment 1 will be out later today — due Tuesday, Feb 15, 9:30am (in 2 weeks)



Sentiment analysis

Why classify?

- Authorship attribution
- Language detection
- News categorization
- and many more!

Text classification



- Inputs:
 - A document *d*
 - A set of classes $C = \{c_1, c_2, c_3, \dots, c_m\}$
- Output:
 - Predicted class *c* for document *d*



Rule-based classification

• Combinations of features on words in document, meta-data

IF there exists word w in document d such that w in [good, great, extra-ordinary, ...], THEN output Positive IF email address ends in [ithelpdesk.com, makemoney.com, spinthewheel.com, ...] THEN output SPAM

- + Can be very accurate
- Rules may be hard to define (and some even unknown to us!)
- Expensive
- Not easily generalizable

VADER-Sentiment-Analysis

VADER (Valence Aware Dictionary and sEntiment Reasoner) is a lexicon and rule-based sentiment analysis tool that is specifically attuned to sentiments expressed in social media. It is fully open-sourced under the [MIT



Supervised Learning: Let's use statistics!

Let the machine figure out the best patterns using data

Inputs:

- Set of *m* classes $C = \{c_1, c_2, ..., c_m\}$

Output:

• Trained classifier, $F: d \rightarrow c$

• Set of *n* 'labeled' documents: $\{(d_1, c_1), (d_2, c_2), \dots, (d_n, c_n)\}$

Key questions: a) What is the form of F? How do we learn F? b)

Types of supervised classifiers



Naive Bayes



Support vector machines



Logistic regression



k-nearest neighbors

Multinomial Naive Bayes

• Simple classification model making use of Bayes rule

• Bayes Rule:

P(c|d) =

$$d - document$$

$$C - class$$

$$P(c) P(d|c)$$

$$P(d)$$



Predicting a class



1 d-document c- class 1

 $T_{MAP} = a_{gmax} P(c|d)$ cEC

How to represent P(d | c)?

• Option 1: represent the entire sequence of words

•
$$P(w_1, w_2, \ldots, w_K | c)$$

- Option 2: Bag of words
 - Assume position of each word is irrelevant (both absolute and relative)

•
$$P(w_1, w_2, \dots, w_K | c) = P(w_1$$

• Probability of each word is *conditionally independent* of the other words given class **c**

(too many sequences!)

 $|c)P(w_2|c)\ldots P(w_k|c)$



Bag of words (BoW)

I love this movie! It's sweet, but with satirical humor. The dialogue is great and the adventure scenes are fun... It manages to be whimsical and romantic while laughing at the conventions of the fairy tale genre. I would recommend it to just about anyone. I've seen it several times, and I'm always happy to see it again whenever I have a friend who hasn't seen it yet!



it	6
1	5
the	4
to	3
and	3
seen	2
yet	1
would	1
whimsical	1
times	1
sweet	1
satirical	1
adventure	1
genre	1
fairy	1
humor	1
have	1
great	1

Predicting with Naive Bayes

• We now have:

P(d|c) P(c)

Predicting with Naive Bayes

• We now have:

$$P(d|c) P(c)$$

 $P(w_1, w_2, \dots, w_k | c) P(c)$

Predicting with Naive Bayes

• We now have:

$$P(d|c) P(c)$$

$$P(w_{1}, w_{2}, \dots, w_{k}|c) P(c)$$

$$P(c) \xrightarrow{k} P(w_{1}|c)$$

$$i=1$$

$$(using Bow assumption)$$





9





Generate the entire data set one document at a time

Estimatin

Maximum likelihood estimates:

 $\hat{P}(c_j) = \frac{(ount(class))}{\sum_{c}}$

g the model
asymmetry
$$P(c) \xrightarrow{k} P(w_i)$$

 $c \xrightarrow{i=1} p(w_i)$
 $c \xrightarrow{k} p(w$



- What if count('amazing', positive) = 0?
 - Implies P('amazing' | positive) = 0
- Given a review document, d = ".... most amazing movie ever ..."



- Data sparsity



• Laplace smoothing:

 $\hat{P}(\omega_i | c) = c_0$



• Effective in practice

Solution: Smoothing!

$$\frac{Sunt(w_{i},c) + \alpha}{\sum_{w} (ount(w,c)] + \alpha |v|} K Vocabulary}$$

Overall process

Input: Set of annotated documents $\{(d_i, c_i)\}_{i=1}^n$

A. Compute vocabulary V of all words

B. Calculate $\hat{P}(c_j) = \frac{\text{Count}(c_j)}{n}$

C. Calculate $\hat{P}(w_i | c_j) = -$

D. (Prediction) Given do

 $c_{MAP} = \arg r$

$$Count(w_i, c_j) + \alpha$$
$$\sum_{w \in V} \left[Count(w, c_j) + \alpha \right]$$

cument
$$d = (w_1, w_2, \dots, w_k)$$

$$\max_c \hat{P}(c) \prod_{i=1}^K \hat{P}(w_i | c)$$

prior

Naive Bayes as a language model

Which class assigns the higher probability to s?

Model pos		Model neg	
0.1	I	0.2	
0.1	love	0.001	love
0.01	this	0.01	this
0.05	fun	0.005	fun
0.1	film	0.1	film



Sentence s				
	love	this	fun	film
0.1 0.2	0.1 0.00	0.01 1 0.01	0.05 0.005	0.1 0.1
	A) pos	B) neg	C) both	equal

Naive Bayes as a language model

Which class assigns the higher probability to s?

Model pos		Model neg	
0.1	I	0.2	
0.1	love	0.001	love
0.01	this	0.01	this
0.05	fun	0.005	fun
0.1	film	0.1	film

Sentence s				
	love	this	fun	film
0.1 0.2	0.1 0.001	0.01 0.01	0.05 0.005	0.1 0.1

P(s|pos) > P(s|neg)

Rank	Category	Feature	Rank	Category	
1	Subject	Number of capitalized words	1	Subject	Min fo
2	Subject	Sum of all the character lengths of words	2	Subject	Min fo
3	Subject	Number of words containing letters and numbers	3	Subject	Min of cha
4	Subject	Max of ratio of digit characters to all characters of each word	4	Subject	Min fo
5	Header	Hour of day when email was sent	5	Subject	Max of th
		(a)			(b
		Spam URLs Feat	tures		
1	URL	The number of all URLs in an email	1	Header	Day of
2	URL	The number of unique URLs in an email	2	Payload	Ν
3	Payload	Number of words containing letters and numbers	3	Payload	Sum of all
4	Payload	Min of the compression ratio for the bz2 compressor	4	Header	Minute of
5	Payload	Number of words containing only letters	5	Header	Hour o

Top features for spam detection

Features

Feature

- of the compression ratio or the bz2 compressor
- of the compression ratio or the zlib compressor
- naracter diversity of each word
- of the compression ratio or the lzw compressor
- he character lengths of words
- week when email was sent
- Number of characters
- the character lengths of words
- of hour when email was sent
- of day when email was sent

- In general, Naive Bayes can use any set of features, not just words:
 - URLs, email addresses, Capitalization, ...
 - Domain knowledge crucial to performance



Evaluating a classifier

• Precision: % of selected classes that are correct



• Recall: % of correct items selected

 $\operatorname{Recall}(+) = \frac{TP}{TP + FN}$

$$Precision(-) = \frac{TN}{TN + FN}$$

$$\mathsf{Recall}(-) = \frac{TN}{TN + FP}$$

- Combined measure using precision and recall
- Harmonic mean of Precision and Recall

$$F_1 = \frac{2 \cdot \Pr}{\Pr}$$

• Or more generally,

$$F_{\beta} = \frac{(1+\beta^2) \cdot |\beta^2 \cdot \text{Pre}|}{\beta^2 \cdot \text{Pre}|}$$

F-Score

ecision · Recall

sision + Recall

Precision · **Recall** ecision + Recall

- Very fast, low storage requirements
- Robust to irrelevant features Irrelevant features cancel each other without affecting results
- Very good in domains with many equally important features Decision trees suffer from fragmentation in such cases — especially if little data
- Optimal if the independence assumptions hold If assumed independence is correct, this is the 'Bayes optimal' classifier
- A good dependable baseline for text classification However, other classifiers can give better accuracy

Advantages of Naive Bayes

Failings of Naive Bayes (1)

Independence assumptions are too strong

x1	x2	Class: x ₁ XOR x ₂
1	1	0
0	1	1
1	0	1
0	0	0

- XOR problem: Naive Bayes cannot learn a decision boundary
- Both variables are jointly required to predict class

Failings of Naive Bayes (2)

Class imbalance

- One or more classes have more instances than others in data
- Data skew causes NB to prefer one class over the other
- Potential solution: Complement Naive Bayes (Rennie et al., 2003)

$$\widehat{O}(w_i | \widetilde{C}_j) = \sum_{\substack{c \neq c_j \\ c \neq c_j}} (\operatorname{ourt}(w_i, c)) \rightarrow (\operatorname{ourt} \# \operatorname{times} w \operatorname{dd})$$

$$\operatorname{occurs} in \operatorname{classes}$$

$$\operatorname{other} \operatorname{than} c$$

$$\sum_{\substack{c \neq c_j \\ c \neq c_j \\ w}} (\operatorname{ourt}(w_i, c))$$

Logistic Regression



Logistic Regression

- Powerful supervised model
- Baseline approach for many NLP tasks
- Connections with neural networks
- Binary (two classes) or multinomial (>2 classes)



Discriminative Model

- Logistic Regression is a *discriminative* model
- Naive Bayes: generative model



Dog: a domesticated carnivorous mammal with a long snout, nonretractable claws, and a barking, howling, or whining voice.

- Inputs:
 - 1. Classification instance in a **feature representation**
 - 2. Classification function to compute \hat{y} using $P(\hat{y} | x)$
 - 3. Loss function (for learning)
 - 4. Optimization algorithm

Overview

• Train phase: Learn the parameters of the model to minimize loss function

• Test phase: Apply parameters to predict class given a new input x

- Input observation: $x^{(i)}$
- Feature vector: $[x_1, x_2, \dots, x_d] = \mathbf{x}$
- Feature j of ith input : $x_i^{(i)}$

I. Feature representation

I love this movie! It's sweet, but with satirical humor. The dialogue is great and the adventure scenes are fun... It manages to be whimsical and romantic while laughing at the conventions of the fairy tale genre. I would recommend it to just about anyone. I've seen it several times, and I'm always happy to see it again whenever I have a friend who hasn't seen it yet!

 $\chi^{(i)}$



Bag of words

2. Classification function

- Given: Input feature vector $\mathbf{x} = [x_1, x_2, \dots, x_d]$
- Output: $P(y = 1 | \mathbf{x})$ and $P(y = 0 | \mathbf{x})$ (binary classification)
- Require a function, $F : \mathbb{R}^d \to [0,1]$
- Sigmoid:



 e^{z}

 $1 + e^{-z}$ $1 + e^{z}$

Weights and Biases

- Which features are important and how much?
- Learn a vector of **weights** and a **bias**
- Weights: Vector of real numbers, $\mathbf{w} = [w_1, w_2, \dots, w_d]$
- Bias: Scalar intercept, b
- Given input features \mathbf{x}_{i} : $z = \mathbf{w} \cdot \mathbf{x} + b$
- Therefore, $f(\mathbf{w} \cdot \mathbf{x} + b) = \frac{c}{1 + e^{\mathbf{w} \cdot \mathbf{x} + b}}$

 $e^{\mathbf{w}\cdot\mathbf{x}+b}$

Putting

• Compute probabilities: P(

P(y=1)

P(y=0) = 1 -

• Decision boundary:

ting it together
ies:
$$P(y = 1 | \mathbf{x}) = \frac{1}{1 + e^{-z}}$$

 $= 1) = \sigma(\mathbf{w} \cdot \mathbf{x} + b) = \frac{1}{1 + e^{-(\mathbf{w} \cdot \mathbf{x} + b)}}$
 $= 1 - \sigma(\mathbf{w} \cdot \mathbf{x} + b)$
 $= 1 - \frac{1}{1 + e^{-(\mathbf{w} \cdot \mathbf{x} + b)}} = \frac{e^{-(\mathbf{w} \cdot \mathbf{x} + b)}}{1 + e^{-(\mathbf{w} \cdot \mathbf{x} + b)}}$
 $\hat{y} = \begin{cases} 1 & \text{if } P(y = 1 | \mathbf{x}) > 0.5 \\ 0 & \text{otherwise} \end{cases}$

Example: Sentiment classification



Var	Definition
x_1	count(positive lexicor
x_2	count(negative lexico
<i>x</i> ₃	$\begin{cases} 1 & \text{if "no"} \in \text{doc} \\ 0 & \text{otherwise} \end{cases}$
x_4	count(1st and 2nd pro
<i>x</i> ₅	$\begin{cases} 1 & \text{if "!"} \in \text{doc} \\ 0 & \text{otherwise} \end{cases}$
x_6	log(word count of do



Example: Sentiment classification

Var	Definition
x_1	count(positive lexicon)
x_2	count(negative lexicon
<i>x</i> ₃	$\begin{cases} 1 & \text{if "no"} \in \text{doc} \\ 0 & \text{otherwise} \end{cases}$
x_4	count(1st and 2nd prop
<i>x</i> ₅	$\begin{cases} 1 & \text{if "!"} \in \text{doc} \\ 0 & \text{otherwise} \end{cases}$
r.	log(word count of doc

 $\log(\text{word count of doc})$ x_6

• Assume weights $\mathbf{w} = [2.5, -5.0, -1.2, 0.5, 2.0, 0.7]$ and bias b = 0.1

$$p(+|x) = P(Y = 1|x) = \sigma(w \cdot x + b)$$

= $\sigma([2.5, -5.0, -1.2, 0.5, 2.0, 0.7] \cdot [3, 2, 1, 3, 0, 4.15] + 0.1)$
= $\sigma(.805)$

- = 0.69

$$p(-|x) = P(Y = 0|x) = 1 - \sigma(x)$$

= 0.31



 $(w \cdot x + b)$

Designing features

- Most important rule: Data is key!
- Linguistic intuition (e.g. part of speech tags, parse trees)
- Complex combinations

$$x_{1} = \begin{cases} 1 & \text{if } ``Case(w_{i}) = \text{Lower''} \\ 0 & \text{otherwise} \end{cases}$$

$$x_{2} = \begin{cases} 1 & \text{if } ``w_{i} \in \text{AcronymDict''} \\ 0 & \text{otherwise} \end{cases}$$

$$x_{3} = \begin{cases} 1 & \text{if } ``w_{i} = \text{St. } \& Case(w_{i-1}) = \text{Cap''} \\ 0 & \text{otherwise} \end{cases}$$

- Feature templates
 - Sparse representations, hash only seen features into index
 - Ex. Trigram(logistic regression classifier) = Feature #78
- Advanced: Representation
 learning (we will see this later!)

Logistic Regression: what's good and what's not

- More freedom in designing features
 - No strong independence assumptions like Naive Bayes
 - More robust to correlated features ("San Francisco" vs "Boston")
 —LR is likely to work better than NB
 - Can even have the same feature twice! (why?)
- May not work well on small datasets (compared to Naive Bayes)
- Interpreting learned weights can be challenging

- We have our **classification function** how to assign weights and bias?
- - true $y : L(\hat{y}, y)$
 - **Optimization algorithm** for updating weights

3. Learning

• Goal: prediction \hat{y} as close as possible to actual label y

• **Distance metric/Loss function** between predicted \hat{y} and

Loss function

- Assume $\hat{y} = \sigma(\mathbf{w} \cdot \mathbf{x} + b)$
- $L(\hat{y}, y) =$ Measure of difference between \hat{y} and y. But what form?
- Maximum likelihood estimation (conditional):
 - input x
 - Similar to language models!

• Choose w and b such that $\log P(y|x)$ is maximized for true labels y paired with

• where we chose parameters to maximize $\log P(w_t | w_{t-n}, \dots, w_{t-1})$ given a corpus

Cross Entropy loss for a single instance

- Assume a single data point (x, y) and two possible classes to choose from
- **Classifier probability:** P(y|x) =
- Log probability: $\log P(y|x) = \log[\hat{y}^y(1-\hat{y})^{1-y}]$ = $y \log \hat{y} + (1-y)\log(1-\hat{y})$ (maximize this)
- Loss: $-\log P(y|x) = -[y\log \hat{y} + (1-y)\log(1-\hat{y})]$ (minimize this)
 - $y = 1 \implies -\log \hat{y}$, and $y = 0 \implies -\log(1 \hat{y})$

$$\hat{y}^{y}(1-\hat{y})^{1-y}$$
 (compact notation

ר)

Cross Entropy loss

- For n data points $(x^{(i)}, y^{(i)})$,
- Classifier probability: $\prod_{i=1}^{n} P(y | x) = \prod_{i=1}^{n} \hat{y}^{y} (1 \hat{y})^{1-y}$

• Loss:
$$-\log \prod_{i=1}^{n} P(y | x) = -\sum_{i=1}^{n} \log x_{i-1}$$

$$L_{CE} = -\sum_{i=1}^{n} \left[y \log \hat{y} + \right]$$

P(y|x)

 $+(1-y)\log(1-\hat{y})$

Example: Computing CE Loss

Var	Definition
x_1	$count(positive lexicon) \in$
x_2	$count(negative lexicon) \in$
<i>x</i> ₃	$\begin{cases} 1 & \text{if "no"} \in \text{doc} \\ 0 & \text{otherwise} \end{cases}$
x_4	count(1st and 2nd pronou
<i>x</i> ₅	$\begin{cases} 1 & \text{if "!"} \in \text{doc} \\ 0 & \text{otherwise} \end{cases}$
x_6	log(word count of doc)

- If y = 1 (positive sentiment), $L_{CE} = -\log(0.69) = 0.37$
- If y = 0 (negative sentiment), $L_{CE} = -\log(0.31) = 1.17$



Properties of CE Loss

•
$$L_{CE} = -\sum_{i=1}^{n} [y^{(i)} \log \hat{y}^{(i)} + \sum_{i=1}^{n} [y^{(i)} \log \hat{y}^{(i)}]]$$

• What values can this loss take?

A) 0 to ∞ B) $-\infty$ to ∞ C) $-\infty$ to 0



- $(1 y^{(i)})\log(1 \hat{y}^{(i)})]$

D) 1 to ∞

Properties of CE Loss

•
$$L_{CE} = -\sum_{i=1}^{n} [y^{(i)} \log \hat{y}^{(i)} + \sum_{i=1}^{n} [y^{(i)} \log \hat{y}^{(i)}] + \dots$$

- Ranges from 0 (perfect predictions) to ∞
- Lower the value, better the classifier
- Cross-entropy between the true distribution P(y|x) in the data and predicted distribution $P(\hat{y}|x)$

 $(1 - y^{(i)})\log(1 - \hat{y}^{(i)})]$

4. Optimization

 $\theta = [w; b]$

$$\hat{\theta} = \arg\min_{\theta} \frac{1}{n} \sum_{i=1}^{n}$$

- Gradient descent:
 - Find direction of steepest slope
 - Move in the opposite direction

• We have our **classification function** and **loss function** - how do we find the best w and b?

 $L_{CE}(y^{(i)}, x^{(i)}; \theta)$



 $\theta^{t+1} = \theta^t - \theta^t -$

$$\eta \frac{d}{d\theta} f(x;\theta)$$

Gradient descent for LR

- Cross entropy loss for logistic regression is **convex** (i.e. has only one global minimum)
 - No local minima to get stuck in
- Deep neural networks are not so easy
 - Non-convex



Learning Rate

• Updates:
$$\theta^{t+1} = \theta^t - \eta \frac{d}{d\theta} f(x; \theta)$$

- Magnitude of movement along gradient
- Higher/faster learning rate = larger updates to parameters



Recap: Logistic regression

• Inputs:

- Classification instance in a **feature representation** 1.
- **Classification function** to compute \hat{y} using $P(\hat{y} | x)$ 2.
- 3. Loss function (for learning)
- 4. Optimization algorithm

• Train phase: Learn the parameters of the model to minimize loss function

• Test phase: Apply parameters to predict class given a new input x

Gradient descent with vector weights

- In LR: weight w is a vector
- Express slope as a partial derivative of loss w.r.t each weight:

$$\nabla_{\theta} L(f(x;\theta),y)) = \begin{cases} \frac{\partial}{\partial w_1} L(f(x;\theta),y) \\ \frac{\partial}{\partial w_2} L(f(x;\theta),y) \\ \vdots \\ \frac{\partial}{\partial w_n} L(f(x;\theta),y) \end{cases}$$

• Updates: $\theta^{(t+1)} = \theta^t - \eta \nabla L(f(x;\theta), y)$



Gradient for logistic regression

•
$$L_{CE} = -\sum_{i=1}^{n} \left[y^{(i)} \log \sigma(\mathbf{w} \cdot \mathbf{x}^{(i)} + b) + (1 - y^{(i)}) \log(1 - \sigma(\mathbf{w} \cdot \mathbf{x}^{(i)} + b)) \right]$$

• Gradient,
$$\frac{dL_{CE}(\mathbf{w}, b)}{dw_{j}} = \sum_{i=1}^{n} [\sigma(\mathbf{w} \cdot \mathbf{x}^{(i)} + b) - y^{(i)}]x_{j}^{(i)}$$

$$\frac{dL_{CE}(\mathbf{w}, b)}{db} = \sum_{i=1}^{n} [\sigma(\mathbf{w} \cdot \mathbf{x}^{(i)} + b) - y^{(i)}]$$
input
frature

Gradient,
$$\frac{dL_{CE}(\mathbf{w}, b)}{dw_{j}} = \sum_{i=1}^{n} \left[\sigma(\mathbf{w} \cdot \mathbf{x}^{(i)} + b) - y^{(i)} \right] x_{j}^{(i)}$$

$$\frac{dL_{CE}(\mathbf{w}, b)}{db} = \sum_{i=1}^{n} \left[\sigma(\mathbf{w} \cdot \mathbf{x}^{(i)} + b) - y^{(i)} \right]$$
input
fratuel

Stochastic Gradient Descent

•	Online optimization	function STC # where:
	Compute loss and minimize after each	# fis # xi # yi
	training example	$\theta \leftarrow 0$ repeat til dor For each tr
	Per Instance	1. Optio Comp Comp
	LOSS	$2. g \leftarrow \nabla$ $3. \theta \leftarrow \theta$ return θ

- DCHASTIC GRADIENT DESCENT(L(), f(), x, y) returns θ L is the loss function is a function parameterized by θ
- is the set of training inputs $x^{(1)}$, $x^{(2)}$,..., $x^{(n)}$ is the set of training outputs (labels) $y^{(1)}$, $y^{(2)}$, ..., $y^{(n)}$

ne # see caption raining tuple $(x^{(i)}, y^{(i)})$ (in random order) onal (for reporting): pute $\hat{y}^{(i)} = f(x^{(i)}; \theta)$ $\nabla_{\boldsymbol{\theta}} L(f(x^{(i)}; \boldsymbol{\theta}), y^{(i)})$ $\theta - \eta g$

- # How are we doing on this tuple? # What is our estimated output \hat{y} ?
- pute the loss $L(\hat{y}^{(i)}, y^{(i)})$ # How far off is $\hat{y}^{(i)}$ from the true output $y^{(i)}$? # How should we move θ to maximize loss? # Go the other way instead



Stochastic Gradient Descent

- Online optimization
- Compute loss and minimize after each training example



Regularization

- Training objective: $\hat{\theta} = \arg \max_{\theta}$
- This might fit the training set too well! (including noisy features)
- Poor generalization to the unseen test set Overfitting
- **Regularization** helps prevent overfitting

$$\hat{\theta} = \arg \max_{\theta} \left[\sum_{i=1}^{n} \log_{\theta} \right]$$

$$\sum_{i=1}^{n} \log P(y^{(i)} | x^{(i)})$$



L2 regularization

•
$$R(\theta) = ||\theta||^2 = \sum_{j=1}^d \theta_j^2$$

- Euclidean distance of weight vector θ from origin
- L2 regularized objective:

$$\hat{\theta} = \arg \max_{\theta} \left[\sum_{i=1}^{n} \log P(y^{(i)} | x^{(i)}) - \alpha \sum_{j=1}^{d} \theta_j^2 \right]$$

LI Regularization

•
$$R(\theta) = ||\theta||_1 = \sum_{j=1}^d |\theta_j|$$

- Manhattan distance of weight vector θ from origin
- L1 regularized objective:

$$\hat{\theta} = \arg \max_{\theta} \left[\sum_{i=1}^{n} \log P(y^{(i)} | x^{(i)}) - \alpha \sum_{j=1}^{d} |\theta_j| \right]$$

- L2 is easier to optimize simpler derivation
 - L1 is complex since the derivative of $|\theta|$ is not continuous at 0
- L2 leads to many small weights (due to θ^2 term)
 - L1 prefers sparse weight vectors with many weights set to 0 (i.e. far fewer features used)

L2 vs L1 regularization



Multinomial Logistic Regression

- What if we have more than 2 classes? (e.g. Part of speech tagging, named entity recognition)
- Need to model $P(y = c | x) \quad \forall c \in C$
- Generalize **sigmoid** function to **softmax**
 - $\operatorname{softmax}(z_i) = -$

$$\frac{e^{z_i}}{\sum_{j=1}^k e^{z_j}} \quad 1 \le i \le k$$
Normalization

- If z = [0, 1, 2, 3, 4], then
- For multinomial LR,

$$P(y = c \mid x) = \frac{e^{\mathbf{w}_c \cdot \mathbf{x} + b_c}}{\sum_{j=1}^k e^{\mathbf{w}_j \cdot \mathbf{x} + b_j}}$$

Softmax

• Similar to sigmoid, softmax squashes values towards 0 or 1

• softmax(z) = ([0.0117, 0.0317, 0.0861, 0.2341, 0.6364])

- Features need to include both input (x) and class (c)
- There were implicit in binary case



Features in multinomial LR

tion	Wt
f "!" \in doc	_1 5
otherwise	-4.5
f "!" \in doc	26
otherwise	2.0
f "!" \in doc	12
otherwise	1.5

Learning

• Generalize binary loss to multinomial CE loss:

$$L_{CE}(\hat{y}, y) = -\sum_{c=1}^{k} 1\{y = c\} \log P(y = c \mid x)$$
$$= -\sum_{c=1}^{k} 1\{y = c\} \log \frac{e^{w_c \cdot x + b_c}}{\sum_{j=1}^{k} e^{w_j \cdot x + b_j}}$$

Gradient:

$$\frac{dL_{CE}}{dw_c} = -\left(1\{y=c\} - P(y=c)\right)$$
$$= -\left(1\{y=c\} - \frac{e^{w_c \cdot x+1}}{\sum_{j=1}^k e^{w_j}}\right)$$

Binary CE Loss: $-\log P(y|x) = -[y \log \hat{y} + (1-y)\log(1-\hat{y})]$

x))x

 $-b_c$ $w_j \cdot x + b_j$ X

