CS388: Natural Language Processing

Lecture 2: Binary Classification



Some slides adapted from Vivek Srikumar, University of Utah



credit: Machine Learning Memes on Facebook



Course enrollment

- Course website: slides, readings, office hours, syllabus
- Mini 1 out, due Tuesday
- Greg's office hours on Thursday are rescheduled to 9am-10am



Linear classification fundamentals

Three discriminative models: logistic regression, perceptron, SVM Different motivations but very similar update rules / inference!

- Optimization
- Sentiment analysis

This Lecture

Classification



- Datapoint x with label $y \in \{0, 1\}$
- For Embed datapoint in a feature space $f(x) \in \mathbb{R}^n$ but in this lecture f(x) and x are interchangeable
- Linear decision rule: $w^{\top}f(x) + b > 0$ $w^{\top}f(x) > 0$
- Can delete bias if we augment feature space: f(x) = [0.5, 1.6, 0.3][0.5, 1.6, 0.3, **1**]

Classification









Linear functions are powerful!





this movie was great! would watch again

that film was <mark>awful,</mark> I'll never watch again

- absence of certain words (great, awful)
- Steps to classification:
 - Turn examples like this into feature vectors
 - Pick a model / learning algorithm
 - Train weights on data to get our classifier

Classification: Sentiment Analysis



Surface cues can basically tell you what's going on here: presence or



this movie was great! would watch again

- Convert this example to a vector using bag-of-words features
- [contains the] [contains a] [contains was] [contains movie] [contains film] ... position 0 position 1 position 2 position 3 position 4
- f(x) = [0] \mathbf{O}
 - Very large vector space (size of vocabulary), sparse features (how many?)
 - Requires indexing the features (mapping them to axes)
 - More sophisticated feature mappings possible (tf-idf), as well as lots of other features: n-grams, character n-grams, parts of speech, lemmas, ...

Positive

1 U





. . .



- Data point $x = (x_1, ..., x_n)$, label $y \in \{0, 1\}$
- Generative models: probabilistic models of P(x,y)
 - Compute P(y|x), predict $\operatorname{argmax}_{y} P(y|x)$ to classify

$$P(y|x) = \frac{P(y)P(x|y)}{P(x)} \propto P(x)$$

P(x) P(x)

- Examples: Naive Bayes (see textbook), Hidden Markov Models
- biscriminative models model P(y|x) directly, compute $\operatorname{argmax}_{u} P(y|x)$
 - Examples: logistic regression
 - Cannot draw samples of x, but typically better classifiers

Generative vs. Discriminative Modeling

- (y)P(x|y)
- ional to"

Logistic Regression



$$P(y = +|x) = \text{logistic}(w^{\top}x)$$
$$P(y = +|x) = \frac{\exp(\sum_{i=1}^{n}x)}{1 + \exp(\sum_{i=1}^{n}x)}$$

To learn weights: maximize discriminative log likelihood of data (log P(y|x)) $\mathcal{L}(\{x_j, y_j\}_{j=1,...,n}) = \sum \log P(y_j | x_j)$ corpus-level LL one (positive) example LL $\mathcal{L}(x_j, y_j = +) = \log P(y_j = +|x_j)$ $= \sum_{i=1}^{n} w_i x_{ji} - \log\left(1 + \exp\left(\sum_{i=1}^{n} w_i x_{ji}\right)\right)$ sum over features

Logistic Regression







 $\mathcal{L}(x_{i}, y_{i} = +) = \log P(y_{i} = + |x_{i}|) =$



Logistic Regression

$$= \sum_{i=1}^{n} w_i x_{ji} - \log \left(1 + \exp \left(\sum_{i=1}^{n} w_i x_{ji} \right) \right)$$

$$= \sum_{i=1}^{n} w_i x_{ji} \left(\sum_{i=1}^{n} w_i x_{ji} \right)$$

$$= \sum_{i=1}^{n} \frac{\partial}{\partial w_i} \left(1 + \exp \left(\sum_{i=1}^{n} w_i x_{ji} \right) \right)$$

$$= \sum_{i=1}^{n} w_i x_{ji} \exp \left(\sum_{i=1}^{n} w_i x_{ji} \right)$$

$$= \sum_{i=1}^{n} w_i x_{ji} \exp \left(\sum_{i=1}^{n} w_i x_{ji} \right)$$

$$= \sum_{i=1}^{n} w_i x_{ji} \exp \left(\sum_{i=1}^{n} w_i x_{ji} \right)$$

$$= \sum_{i=1}^{n} w_i x_{ji} \exp \left(\sum_{i=1}^{n} w_i x_{ji} \right)$$

$$= \sum_{i=1}^{n} w_i x_{ji} \exp \left(\sum_{i=1}^{n} w_i x_{ji} \right)$$

$$= \sum_{i=1}^{n} w_i x_{ji} \exp \left(\sum_{i=1}^{n} w_i x_{ji} \right)$$

$$= \sum_{i=1}^{n} w_i x_{ji} \exp \left(\sum_{i=1}^{n} w_i x_{ji} \right)$$

$$= \sum_{i=1}^{n} w_i x_{ji} \exp \left(\sum_{i=1}^{n} w_i x_{ji} \right)$$

 $= x_{ji} - x_{ji} \frac{\exp\left(\sum_{i=1}^{n} w_i x_{ji}\right)}{1 + \exp\left(\sum_{i=1}^{n} w_i x_{ji}\right)} = x_{ji} \left(1 - P(y_j = +|x_j)\right)$





- Gradient of w_i on positive example
 - If P(+) is close to 1, make very little update Otherwise make w_i look more like x_{ii} , which will increase P(+)
- Gradient of w_i on negative examp
 - If P(+) is close to 0, make very little update Otherwise make w_i look less like x_{ii} , which will decrease P(+)
- Let $y_i = 1$ for positive instances, $y_i = 0$ for negative instances.
- Can combine these gradients as a

Logistic Regression

$$\mathbf{e} = x_{ji}(1 - P(y_j = +|x_j))$$

$$\mathsf{ole} = x_{ji}(-P(y_j = +|x_j))$$

$$x_j(y_j - P(y_j = 1|x_j))$$



- $f(x_1) = [1]$ (1) this movie was great! would watch again 1] + (2) I expected a great movie and left happy $f(x_2) = [1]$ + 1] (3) great potential but ended up being a flop $f(x_3) = [1]$ 0 [contains great] [contains movie] position 0 position 1 $w = [0, 0] \longrightarrow P(y = 1 | x_1) = \exp(0)/(1 + \exp(0)) = 0.5 \longrightarrow g = [0.5, 0.5]$ $\rightarrow q = [0.25, 0.25]$ $\rightarrow q = [-0.67, 0]$ $w = [0.75, 0.75] \rightarrow P(y = 1 | x_3) = \text{logistic}(0.75) \approx 0.67$ $P(y = +|x) = \text{logistic}(w \mid x)$ $w = [0.08, 0.75] \dots$ $x_j(y_j - P(y_j = 1|x_j))$

Example







 \boldsymbol{m}

Regularization

Regularizing an objective can mean many things, including an L2norm penalty to the weights:

$$\sum_{j=1}^{m} \mathcal{L}(x_j, y_j) - \lambda ||w||_2^2$$

- Keeping weights small can prevent overfitting
- For most of the NLP models we build, explicit regularization isn't necessary
 - Early stopping

 - Large numbers of sparse features are hard to overfit in a really bad way For neural networks: dropout and gradient clipping







Logistic Regression: Summary

Model

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

Inference

 $\operatorname{argmax}_{y} P(y|x)$

 $P(y = 1|x) \ge 0.5 \Leftrightarrow w^{\top}x \ge 0$

 $\frac{w_i x_i}{w_i x_i}$

Learning: gradient ascent on the (regularized) discriminative log-likelihood



Perceptron/SVM



Simple error-driven learning approach similar to logistic regression

Decision rule: $w^{\top}x > 0$

If incorrect: if positive, $w \leftarrow w + x$ if negative, $w \leftarrow w - x$

• Guaranteed to eventually separate the data if the data are separable

Perceptron





Support Vector Machines

Many separating hyperplanes — is there a best one?





Support Vector Machines

Many separating hyperplanes — is there a best one?





Constraint formulation: find w via following quadratic program:

Min	imize	$\ w\ _{2}^{2}$
s.t.	$\forall j$	$w^{\top} x_j \ge 1$ if $y_j = 1$
		$w^{\top} x_j \leq -1 \text{ if } y_j = 0$

As a single constraint:

 $\forall j \ (2y_j - 1)(w^{\top}x_j) \geq 1$

Support Vector Machines

minimizing norm with fixed margin <=> maximizing margin

Generally no solution (data is generally non-separable) — need slack!



N-Slack SVMs

Minimize
$$\lambda \|w\|_2^2 + \sum_{j=1}^m \xi_j$$

s.t. $\forall j \ (2y_j - 1)(w^\top x_j) \ge 1 - \xi_j \qquad \forall j \ \xi_j \ge 0$

The ξ_j are a "fudge factor" to make all constraints satisfied

Take the gradient of the objective: ∂ $\frac{\partial}{\partial w_i} \xi_j = 0 \text{ if } \xi_j = 0$

Looks like the perceptron! But updates more frequently

$$\xi_j = (2y_j - 1)x_{ji}$$
 if $\xi_j > 0$
= x_{ji} if $y_j = 1, -x_{ji}$ if $y_j = 1$

 \cap



*gradients are for maximizing things, which is why they are flipped





Logistic regression (unregularize

$$x(y - P(y = 1|x)) = x(y - \log x)$$

Perceptron (2y-1)x if classified incorrect 0 else

SVM (2y-1)x if not classified corre () else

Comparing Gradient Updates (Reference)

ed) $\operatorname{ogistic}(w^{ op}x))$	y = 1 for pos, 0 for neg
۱ у	
ectly with margin of 1	

Optimization



Four elements of a structured machine learning method:

Model: probabilistic, max-margin, deep neural network



- Training: gradient descent?

Structured Prediction

Inference: just maxes and simple expectations so far, but will get harder



- Stochastic gradient *ascent*
 - Very simple to code up
 - "First-order" technique: only relies on having gradient
- Newton's method
 - Second-order technique
 - Optimizes quadratic instantly
- Quasi-Newton methods: L-BFGS, etc. approximate inverse Hessian

$$w \leftarrow w + \alpha g, \quad g = \frac{\partial}{\partial w} \mathcal{L}$$

Can avg gradient over a few examples and apply update once (minibatch) Setting step size is hard (decrease when held-out performance worsens?)

$$w \leftarrow w + \left(\frac{\partial^2}{\partial w^2}\mathcal{L}\right)^{-1}g$$

Inverse Hessian: *n* x *n* mat, expensive!



Optimized for problems with sparse features

that get updated frequently

$$w_i \leftarrow w_i + \alpha \frac{1}{\sqrt{\epsilon + \sum_{\tau=1}^t g_{\tau,i}^2}} g_{t_i}$$

- Other techniques for optimizing deep models more later!

AdaGrad

Per-parameter learning rate: smaller updates are made to parameters

(smoothed) sum of squared gradients from all updates

Generally more robust than SGD, requires less tuning of learning rate

Duchi et al. (2011)



Implementation

Supposing k active features on an instance, gradient is only nonzero on k dimensions

$$w \leftarrow w + \alpha g, \quad g = \frac{\partial}{\partial w} \mathcal{L}$$

- k < 100, total num features = 1M+ on many problems</p>
- Be smart about applying updates!
- In PyTorch: applying sparse gradients only works for certain optimizers and sparse updates are very slow. The code we give you is much faster



this movie was great! would watch again

the movie was gross and overwrought, but I liked it

this movie was not really very enjoyable

- Bag-of-words doesn't seem sufficient (discourse structure, negation)
- There are some ways around this: extract bigram feature for "not X" for all X following the *not*





Bo Pang, Lillian Lee, Shivakumar Vaithyanathan (2002)





	Features	# of	frequency or	NB	ME	SVM
		features	presence?			
(1)	unigrams	16165	freq.	78.7	N/A	72.8
(2)	unigrams	"	pres.	81.0	80.4	82.9
(3)	unigrams+bigrams	32330	pres.	80.6	80.8	82.7
(4)	bigrams	16165	pres.	77.3	77.4	77.1
(5)	unigrams+POS	16695	pres.	81.5	80.4	81.9
(6)	adjectives	2633	pres.	77.0	77.7	75.1
(7)	top 2633 unigrams	2633	pres.	80.3	81.0	81.4
(8)	unigrams+position	22430	pres.	81.0	80.1	81.6

Simple feature sets can do pretty well!

Bo Pang, Lillian Lee, Shivakumar Vaithyanathan (2002)





Method	RT-s	M
MNB-uni	77.9	8
MNB-bi	79.0	8
SVM-uni	76.2	8
SVM-bi	77.7	<u>8</u>
NBSVM-uni	78.1	8
NBSVM-bi	<u>79.4</u>	8
RAE	76.8	8
RAE-pretrain	77.7	8
Voting-w/Rev.	63.1	8
Rule	62.9	8
BoF-noDic.	75.7	8
BoF-w/Rev.	76.4	8
Tree-CRF	77.3	8
BoWSVM		
Kim (2014) CNNs	81.5	8

- PQA 35.3 36.3 36.1

- **36.7** 35.3
- **86.3**
- 85.7
- **6.4**
- 1.7
- 1.8
- 1.8
- 84.1
- 6.1
- 9.5

Naive Bayes is doing well!

Ng and Jordan (2002) — NB can be better for small data

Before neural nets had taken off results weren't that great

Wang and Manning (2012)







Model

- Stanford Sentiment Treebank (SST) binary classification
- Best systems now: large pretrained networks
- 90 -> 97 over the last 2 years

XLNet-Large (ensemble)	(Y
2019)	

MT-DNN-ensemble (Liu et

Snorkel MeTaL(ensemble) (al., 2018)

MT-DNN (Liu et al., 2019)

Bidirectional Encoder Representations from Trans (Devlin et al., 2018)

Neural Semantic Encoder (Munkhdalai and Yu, 2017)

BLSTM-2DCNN (Zhou et al.

https://github.com/sebastianruder/NLP-progress/blob/master/english/sentiment_analysis.md

	Accuracy	Paper / Source	С
ang et al.,	96.8	XLNet: Generalized Autoregressive Pretraining for Language Understanding	Offi
al., 2019)	96.5	Improving Multi-Task Deep Neural Networks via Knowledge Distillation for Natural Language Understanding	Offi
Ratner et	96.2	Training Complex Models with Multi-Task Weak Supervision	Offi
	95.6	Multi-Task Deep Neural Networks for Natural Language Understanding	Offi
sformers	94.9	BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding	Offi

	89.7	Neural Semantic Encoders	
., 2017)	89.5	Text Classification Improved by Integrating Bidirectional LSTM with Two-dimensional Max Pooling	

 $\bullet \bullet \bullet$



• Logistic regression: $P(y = 1|x) = \frac{\exp\left(\sum_{i=1}^{n} w_i x_i\right)}{(1 + \exp\left(\sum_{i=1}^{n} w_i x_i\right))}$ Decision rule: $P(y = 1|x) \ge 0.5 \Leftrightarrow w^{\top} x \ge 0$ Gradient (unregularized): x(y - P(y = 1|x))**SVM:** Decision rule: $w^{\top}x > 0$

Recap

(Sub)gradient (unregularized): 0 if correct with margin of 1, else x(2y-1)

Logistic regression, SVM, and perceptron are closely related

SVM and perceptron inference require taking maxes, logistic regression has a similar update but is "softer" due to its probabilistic nature

- All gradient updates: "make it look more like the right thing and less like the wrong thing"
- Next time: multiclass classification

