Machine Learning for NLP Lecture 3: Optimization and machine learning



Richard Johansson

September 8, 2016

・ロト ・ 日 ・ ・ 日 ・ ・ 日 ・ ・ つ へ ()

linear classifiers

a linear classifier is a classifier that is defined in terms of a scoring function like this

 $score = w \cdot x$

- explanation of the parts:
 - x is a vector with features of what we want to classify (e.g. made with a DictVectorizer)
 - ▶ w is a vector representing which features the classifier thinks are important

ション ふゆ く 山 マ チャット しょうくしゃ

- is the dot product between the two vectors
- ▶ for now, we'll assume that there are two classes: binary classification
 - return the first class if the score > 0
 - otherwise the second class
- the essential idea: features are scored independently

geometric view

 geometrically, a linear classifier can be interpreted as separating the vector space into two regions with a line (plane, hyperplane)



▲□▶ ▲圖▶ ★ 国▶ ★ 国▶ - 国 - の Q @



optimization and machine learning

- we will now consider models that are less ad-hoc than the perceptron
- idea: define an objective function based on the fundamental tradeoff in machine learning:
 - how well we handle the training set (loss)
 - simplicity of the model (regularization)
- ...and then the training consists of applying optimization techniques to find the best w

ション ふゆ く 山 マ チャット しょうくしゃ

- we will consider two models:
 - logistic regression, based on probability
 - support vector classifier, based on geometry



in scikit-learn

LR is called sklearn.linear_model.LogisticRegression

▲□▶ ▲圖▶ ▲臣▶ ★臣▶ ―臣 …の�?

SVM is called sklearn.svm.LinearSVC



overview

logistic regression

support vector machines

basic optimization

practical information



◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 の�?

linear classifiers with probabilities?

linear classifiers select the outputs based on a scoring function:

score =
$$\mathbf{w} \cdot \mathbf{x}$$

how to convert the scores into probabilities?

idea: use a logistic or sigmoid function:

$$P(ext{positive output}|x) = rac{1}{1+e^{- ext{score}}}$$

where $e^{-\text{score}} = \text{math.exp}(-\text{score})$

this is formally a probability: always between 0 and 1, sum of probablities of possible outcomes = 1



the logistic function





conversely

$$P(\mathsf{negative output}|m{x}) = rac{1}{1+e^{\mathsf{score}}}$$



making it a bit more compact

 \blacktriangleright if we code the positive class as +1 and the negative class as -1, then we can write the probability a bit more neatly:

$$P(y|\mathbf{x}) = rac{1}{1 + e^{-y \cdot \mathsf{score}}}$$

・ロト ・ 日 ・ ・ 日 ・ ・ 日 ・ ・ つ へ ()



recall: the maximum likelihood principle

- select the model that gives a high probability to the data
- ▶ in our case, the "model" is the weight vector *w*
- ▶ adjust *w* so that each output label gets a high probability

ション ふゆ アメリア メリア しょうくの



the likelihood function

- formally, the "probability of the data" is defined by the likelihood function
- this is the product of the probabilities of all *m* individual training instances:

$$L(\boldsymbol{w}) = P(y_1|\boldsymbol{x}_1) \cdot \ldots \cdot P(y_T|\boldsymbol{x}_m)$$

in our case, this means

$$L(\boldsymbol{w}) = \frac{1}{1 + e^{-y_1 \cdot (\boldsymbol{w} \cdot \boldsymbol{x}_1)}} \cdot \ldots \cdot \frac{1}{1 + e^{-y_m \cdot (\boldsymbol{w} \cdot \boldsymbol{x}_m)}}$$

・ロト ・ 日 ・ ・ 日 ・ ・ 日 ・ ・ つ へ ()



rewriting a bit...

we rewrite the previous formula

$$L(\boldsymbol{w}) = \frac{1}{1 + e^{-y_1 \cdot (\boldsymbol{w} \cdot \boldsymbol{x}_1)}} \cdot \ldots \cdot \frac{1}{1 + e^{-y_m \cdot (\boldsymbol{w} \cdot \boldsymbol{x}_m)}}$$

as

$$-\log L(\boldsymbol{w}) = \operatorname{Loss}(\boldsymbol{w}, \boldsymbol{x}_1, y_1) + \ldots + \operatorname{Loss}(\boldsymbol{w}, \boldsymbol{x}_m, y_m)$$

where

$$Loss(\boldsymbol{w}, \boldsymbol{x}, \boldsymbol{y}) = \log(1 + \exp(-\boldsymbol{y} \cdot (\boldsymbol{w} \cdot \boldsymbol{x})))$$

▲□▶ ▲圖▶ ▲臣▶ ★臣▶ ―臣 …の�?

is called the log loss function



plot of the log loss



◆□▶ ◆圖▶ ◆臣▶ ◆臣▶ 臣 - のへで



recall: the fundamental tradeoff in machine learning



- goodness of fit: the learned classifier should be able to correctly classify the examples in the training data
- regularization: the classifier should be simple
- but so far in our LR description, we've just taken care of the first part!



what does it mean to "keep the classifier simple"?

- concretely, how can we add regularization to the LR model?
- the most common approach is to add a term that keeps the weights small
- formally, we say that the the squared length (norm) of the weight vector should be small:

$$|\boldsymbol{w}|^2 = w_1 \cdot w_1 + \ldots + w_n \cdot w_n = \boldsymbol{w} \cdot \boldsymbol{w}$$

・ロト ・ 日 ・ ・ 日 ・ ・ 日 ・ ・ つ へ ()



we combine the loss and the regularizer:

$$\sum \text{Loss}(\boldsymbol{w}, \boldsymbol{x}_i, y_i) + \lambda \cdot |\boldsymbol{w}|^2$$

- ▶ in this formula, \u03c6 is a "tweaking" parameter that controls the tradeoff between loss and regularization
- note: in some formulations (including scikit-learn), there is a parameter C instead of the λ that is put before the loss

ション ふゆ アメリア メリア しょうくの



this still doesn't look implementable...

we have an objective function that we want to minimize:

$$\sum \text{Loss}(\boldsymbol{w}, \boldsymbol{x}_i, y_i) + \lambda \cdot |\boldsymbol{w}|^2$$

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ●

but we still don't know how to convert this into code!



overview

logistic regression

support vector machines

basic optimization

practical information



◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 の�?

geometric view

 geometrically, a linear classifier can be interpreted as separating the vector space into two regions with a line (plane, hyperplane)



▲□▶ ▲圖▶ ★ 国▶ ★ 国▶ - 国 - の Q @



margin of separation

• the margin γ denotes how well \pmb{w} separates the classes:



(日) (四) (日) (日)

æ



► a result from statistical learning theory:





▲ロト ▲周ト ▲ヨト ▲ヨト 三日 - のへで

expected test error = training error + BigUglyFormula
$$(rac{1}{\gamma^2})$$

 \blacktriangleright larger margin \rightarrow better generalization



support vector machines

 support vector machines (SVMs) or support vector classifiers (SVC) are linear classifiers constructed by selecting the w that maximizes the margin



 note: the solution depends only on the borderline examples: the support vectors



soft-margin SVMs

- ▶ in some cases the dataset is inseparable, or nearly inseparable
- soft-margin SVM: allow some examples to be disregarded when maximizing the margin



・ロト ・ 日 ・ ・ 日 ・ ・ 日 ・

ж

stating the SVM as using an objective function

- the hard-margin and soft-margin SVM can be stated mathematically in a number of ways
- we'll skip the details, but with a bit of work we can show that the soft-margin SVM can be stated as minimizing

$$\sum \mathsf{Loss}(oldsymbol{w},oldsymbol{x}_i,y_i) + \lambda \cdot |oldsymbol{w}|^2$$

where

$$Loss(\boldsymbol{w}, \boldsymbol{x}, \boldsymbol{y}) = max(0, 1 - \boldsymbol{y} \cdot (\boldsymbol{w} \cdot \boldsymbol{x}))$$

・ロト ・ 日 ・ ・ 日 ・ ・ 日 ・ ・ つ へ ()

is called the hinge loss



plot of the hinge loss



◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 三臣 - のへで



overview

logistic regression

support vector machines

basic optimization

practical information



▲□▶ ▲圖▶ ▲臣▶ ★臣▶ 三臣 - のへで

optimization

- what is optimization?
- unconstrained optimization: find the x that gives us the minimal (or maximal) value of some function f:

$\min_{x} f(x)$

constrained optimization: find the x that gives us the minimal (or maximal) value of f, where x satisfies some extra conditions:

 $\min_{x} f(x)$
such that x > 0

ション ふゆ アメリア メリア しょうくの

today unconstrained optimization only

one-variable example







two-variable example





remember your highschool calculus...

- in your early school days, you might have seen the derivative of a function
- intuition: the derivative measures the slope



 if a "nice" function has a maximum or minimum, then the derivative will be zero there

OTHENBUR

the gradient

- the multidimensional equivalent of the derivative is called the gradient
- ▶ if f is a function of n variables, then the gradient is an n-dimensional vector, often written $\nabla f(x)$
- intuition: the gradient points in the uphill direction



again: the gradient is zero if we have an optimum



computing the gradient





Download page

POWERED BY THE WOLFRAM LANGUAGE



gradient descent

▶ as we saw, the gradient points in the uphill direction:



this intuition leads to a simple idea for finding the minimum:

- take a small step in the direction opposite to the gradient
- repeat until the gradient is close enough to zero
- this is called gradient descent

gradient descent, pseudocode

- the same thing again, in pseudocode:
 - 1. set x to some initial value, and select a suitable step size c
 - 2. compute the gradient $\nabla f(x)$
 - 3. if $\nabla f(x)$ is small enough, we are done
 - 4. otherwise, subtract $c \cdot \nabla f(x)$ from x and go back to step 2

ション ふゆ アメリア メリア しょうくの

conversely, to find the maximum we can do gradient ascent: then we instead add c · ∇f(x) to x



in Python

```
def gradient_ascent(x_init, y_init,
                threshold = 0.001,
                 steplength = 0.01:
x = x_{init}
y = y_{init}
done = False
while not done:
    gxy = gradient_of_my_function(x, y)
    x += steplength * gxy[0]
    y += steplength * gxy[1]
    if numpy.linalg.norm(gxy) < threshold:
        done = True
return (x, y)
```


let's optimize this function:

```
def f(x, y):
return math.exp(-(x-2)**2 - (y+1)**2)
```

```
its gradient is
```

```
def gradient_of_f(x, y):
return (-2*(x-2)*f(x, y), -2*(y+1)*f(x, y))
```

▲ロト ▲周ト ▲ヨト ▲ヨト ヨー のへで








































































will we always reach the top?

▶ yes, if

- there is actually a top
- the step is short enough
- the surface isn't too jumpy
- smarter versions of gradient ascent/descent try to adapt the step length so that we don't go too slow in the beginning, or bounce around the top at the end

・ロト ・ 日 ・ モート ・ 田 ・ うへで



```
gradient ascent example (2)
```

```
let's optimize another function:
```

```
def f(x, y):
    return math.exp( -(x-2)**2 - 0.5*(y+1)**2) \
        + 0.7 * math.exp( -0.7*(x+1)**2 - 0.8*(y-1)**2)
```

◆□▶ ◆□▶ ★□▶ ★□▶ □ のQ@



















local and global maxima/minima

some functions have local maxima or minima



 these functions are harder to optimize because the local (but not global) optima also have a gradient of 0



convex and concave functions

- ► a function is **convex** if it always curves downwards
- equivalently, if we draw a line between two points of the surface, the surface is always below the line



- the point of this: if we find a local optimum (gradient is 0) of a convex function, this is guaranteed to be the minimum
- conversely, a function is concave if it always curves upwards

stochastic gradient descent

in some cases it is cumbersome to compute the gradient

- because it depends on all the data in the training set
- stochastic gradient descent: simplify the computation by computing the gradient using just a small part
 - typically, a single training example
- pseudocode:
 - 1. set w to some initial value, and select a suitable step size c
 - 2. select a single training instance x
 - 3. compute the gradient $\nabla f(w)$ using x only
 - 4. if we are "done", stop
 - 5. otherwise, subtract $c \cdot \nabla f(w)$ from w and go back to step 2

ション ふゆ くり くり くし くし

(stopping criterion shouldn't be based on just a single instance)



SVM and LR have convex objective functions



UNIVERSITY OF GOTHENBURG

optimizing SVM and LR

- since the objective functions of SVM and LR are convex, we can find w by stochastic gradient descent
- pseudocode:
 - set w to some initial value, e.g. all zero
 - iterate a fixed number of times:
 - select a single training instance x
 - \blacktriangleright select a "suitable" step length η
 - compute the gradient of the hinge loss or log loss

◆□▶ ◆□▶ ◆□▶ ◆□▶ ● ● ●

- subtract step length · gradient from w
- note the similarity to the perceptron!



overview

logistic regression

support vector machines

basic optimization

practical information





some comments about assignment 2

- implement SVM and LR and test them in a document classifier
- we'll use the Pegasos algorithm see assignment page
- Pegasos works in an iterative fashion similar to the perceptronso if you start from my perceptron code this will be a breeze
- optional tasks to speed up the implementation using sparse vectors

◆□▶ ◆□▶ ◆□▶ ◆□▶ ● ● ●



some clarifications about the paper

- the important part of the paper is the pseudocode in Figure 1
- Pegasos adapts the step length η over time: long steps in the beginning, smaller in the end
- $\langle \boldsymbol{w}, \boldsymbol{x}
 angle$ is the dot product $\boldsymbol{w} \cdot \boldsymbol{x}$
- S is the training set, |S| is the size of S
- ► T is the number of steps in the algorithm.
 - this is a bit different from our perceptron, where we specified the number of times to process the whole training set.
- the optional line is there for theoretical reasons and can be ignored
- a subgradient is a gradient for "abrupt" functions such as the hinge loss



practical information

- solve the assignment individually
- two lab sessions next week
- deadline one week later: September 24

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ●



seminar tomorrow

 Peter will present a classical paper about document polarity classification

▲□▶ ▲圖▶ ▲臣▶ ★臣▶ ―臣 …の�?

I'll present some additional material



- on Friday so no seminar next week
- neural networks:
 - non-linear classifiers
 - they are also implemented using objective functions and optimization

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ●

