

# Applied Machine Learning

## Lecture 5-2: Logistic regression and SVM

**Selpi (selpi@chalmers.se)**

The slides are further development of Richard Johansson's slides

February 7, 2020

# Overview

A bit on perceptron

logistic regression

training a logistic regression classifier

detour: multiclass linear classifiers

support vector classification

optimizing the LR and SVM objectives

# The perceptron algorithm and the simpler version

## Perceptron algorithm

$\mathbf{w} = (0, \dots, 0)$

repeat  $N$  times

  for  $(\mathbf{x}_i, y_i)$  in the training set

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

    if a pos. is misclassified

$\mathbf{w} = \mathbf{w} + \mathbf{x}_i$

    else if a neg. is misclassified

$\mathbf{w} = \mathbf{w} - \mathbf{x}_i$

return  $\mathbf{w}$

# The perceptron algorithm and the simpler version

## Perceptron algorithm

```
w = (0, ..., 0)
repeat N times
  for (xi, yi) in the training set
    score = w · xi
    if a pos. is misclassified
      w = w + xi
    else if a neg. is misclassified
      w = w - xi
return w
```

## The simpler version

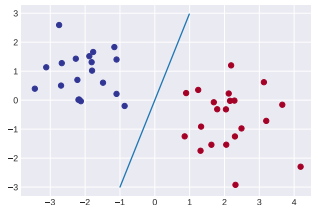
if the  $y_i$  are coded as +1  
or -1:

```
w = (0, ..., 0)
for (xi, yi) in the training set
  score = w · xi
  if  $y_i \cdot \text{score} \leq 0$ 
    w = w +  $y_i \cdot x_i$ 
return w
```

# how can we get the “certainty” of a linear classifier?

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

- ▶ **large positive score**: quite certain that  $\mathbf{x}$  belongs to the positive class
- ▶ **large negative score**: quite certain that  $\mathbf{x}$  belongs to the negative class
- ▶ **near zero**: we are unsure



# Overview

A bit on perceptron

logistic regression

training a logistic regression classifier

detour: multiclass linear classifiers

support vector classification

optimizing the LR and SVM objectives

# The logistic regression model

- ▶ **logistic regression** is a method to train a linear classifier that gives a probabilistic output
- ▶ how to get the probability? use a **logistic** or **sigmoid** function:

$$P(\text{positive output}|\mathbf{x}) = \frac{1}{1 + e^{-\text{score}}}$$

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

# The logistic regression model

- ▶ **logistic regression** is a method to train a linear classifier that gives a probabilistic output
- ▶ how to get the probability? use a **logistic** or **sigmoid** function:

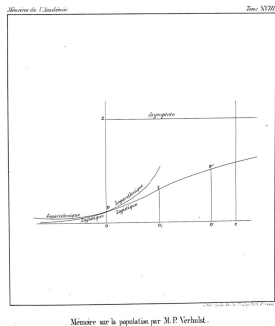
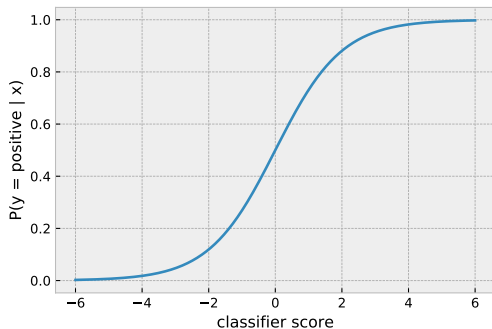
$$P(\text{positive output}|\mathbf{x}) = \frac{1}{1 + e^{-\text{score}}}$$

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

$$P(\text{negative output}|\mathbf{x}) = 1 - \frac{1}{1 + e^{-\text{score}}} = \frac{1}{1 + e^{\text{score}}}$$



# the logistic / sigmoid function [Verhulst, 1845]



## interpreting the linear classifier's score

- ▶ the output score  $\mathbf{w} \cdot \mathbf{x}$  can now be interpreted as the **log odds** in favor of the positive outcome
- ▶ **odds**: how much more likely is the positive outcome than the negative outcome?

$$\text{odds} = \frac{p}{1 - p}$$

*Knew that we ventured on such dangerous seas  
That if we wrought out life 'twas **ten to one***

Shakespeare, *Henry IV, Part II*, Act I, Scene 1 lines 183–4.

making it a bit more compact

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

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

## in scikit-learn

- ▶ LR is called `sklearn.linear_model.LogisticRegression`
- ▶ `predict_proba` gives the probability output

code example: using a logistic regression classifier

# Overview

A bit on perceptron

logistic regression

training a logistic regression classifier

detour: multiclass linear classifiers

support vector classification

optimizing the LR and SVM objectives

## recall: the maximum likelihood principle

- ▶ in a probabilistic model, we can train the model by selecting parameters that assign a **high probability to the data**
- ▶ in our case, the parameters are the weight vector  $\mathbf{w}$
- ▶ adjust  $\mathbf{w}$  so that each output label gets a high probability

## 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:

$$\mathcal{L}(\mathbf{w}) = P(y_1|\mathbf{x}_1) \cdot \dots \cdot P(y_m|\mathbf{x}_m)$$

- ▶ in our case, this means

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



rewriting a bit...

- ▶ we rewrite the previous formula

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

as

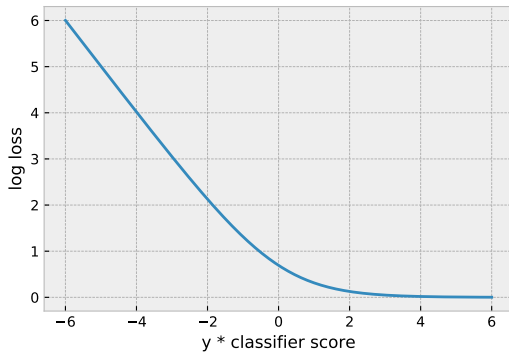
$$-\log \mathcal{L}(\mathbf{w}) = \text{Loss}(\mathbf{w}, \mathbf{x}_1, y_1) + \dots + \text{Loss}(\mathbf{w}, \mathbf{x}_m, y_m)$$

where

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

is called the **log loss** function

## plot of the log loss



# 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!

$$-\log \mathcal{L}(\mathbf{w}) = \text{Loss}(\mathbf{w}, \mathbf{x}_1, y_1) + \dots + \text{Loss}(\mathbf{w}, \mathbf{x}_m, y_m)$$

## regularization in logistic regression models

- ▶ just like we saw for linear regression models (Ridge and Lasso), we can add a regularizer that **keeps the weights small**
- ▶ most commonly, the  **$L_2$  regularizer**:

$$\|\mathbf{w}\|^2 = w_1 \cdot w_1 + \dots + w_n \cdot w_n = \mathbf{w} \cdot \mathbf{w}$$

- ▶ ...or an  **$L_1$  regularizer**:

$$\|\mathbf{w}\|_1 = |w_1| + \dots + |w_n|$$

which will do some feature selection

## combining the pieces

- ▶ we combine the loss and the regularizer:

$$\frac{1}{N} \cdot \sum_{i=1}^N \text{Loss}(\mathbf{w}, \mathbf{x}_i, y_i) + \frac{\lambda}{2} \cdot \|\mathbf{w}\|^2$$

- ▶ in this formula,  $\lambda$  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  $\lambda$  that is put before the loss

$$\frac{C}{N} \cdot \sum_{i=1}^N \text{Loss}(\mathbf{w}, \mathbf{x}_i, y_i) + \frac{1}{2} \|\mathbf{w}\|^2$$

check

$$\frac{C}{N} \cdot \sum_{i=1}^N \text{Loss}(\mathbf{w}, \mathbf{x}_i, y_i) + \frac{1}{2} \|\mathbf{w}\|^2$$

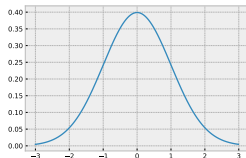
- ▶ how do we convert this into an algorithm?

# probabilistic justification of the regularizers

- ▶ by adding the regularizer, we are carrying out **maximum a posteriori** estimation instead of maximum likelihood

- ▶ in MAP estimation, we use a **prior** to push the parameters in a desired direction

$$\mathbf{w}_{\text{MAP}} = \arg \max_{\mathbf{w}} p(\text{data}|\mathbf{w}) \cdot p(\mathbf{w})$$

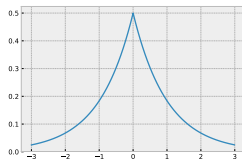


- ▶  $L_2$  regularizer = **Gaussian prior**

$$-\log p(\mathbf{w}) = \lambda \|\mathbf{w}\|^2 + \text{constant}$$

- ▶  $L_1$  regularizer = **Laplace prior**

$$-\log p(\mathbf{w}) = \lambda |\mathbf{w}| + \text{constant}$$



# Overview

A bit on perceptron

logistic regression

training a logistic regression classifier

**detour: multiclass linear classifiers**

support vector classification

optimizing the LR and SVM objectives



## two-class (binary) linear classifiers

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

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

- ▶ this is a **binary** (2-class) classifier:
  - ▶ return the first class if the score  $> 0$
  - ▶ ...otherwise the second class
- ▶ how can we deal with non-binary (**multi-class**) problems when using linear classifiers?

## approaches to multiclass classification

- ▶ **idea 1**: break down the complex problem into simpler problems, train a classifier for each separately
- ▶ **idea 2**: modify the learning algorithm so that it can handle the multiclass case directly

## idea 1: reduction from multiclass to binary

- ▶ **one-versus-rest** (“long jump”):
  - ▶ for each class  $c$ , make a binary classifier to distinguish  $c$  from all other classes
  - ▶ so if there are  $n$  classes, there are  $n$  classifiers
  - ▶ at test time, we select the class giving the highest score
- ▶ **one-versus-one** (“football league”):
  - ▶ for each pair of classes  $c_1$  and  $c_2$ , make a classifier to distinguish  $c_1$  from  $c_2$
  - ▶ if there are  $n$  classes, there are  $\frac{n \cdot (n-1)}{2}$  classifiers
  - ▶ at test time, we select the class that has most “wins”

## example

- ▶ assume we're training a classifier of fruits and we have the classes apple, orange, mango
- ▶ in one-vs-rest, we train the following three classifiers:
  - ▶ apple vs orange+mango
  - ▶ orange vs apple+mango
  - ▶ mango vs apple+orange
- ▶ in one-vs-one, we train the following three:
  - ▶ apple vs orange
  - ▶ apple vs mango
  - ▶ orange vs mango

## example (continued)

- ▶ we train classifiers to distinguish between apple, orange, and mango, using one-vs-rest
  - ▶ so we get  $\mathbf{w}_{\text{apple}}$ ,  $\mathbf{w}_{\text{orange}}$ ,  $\mathbf{w}_{\text{mango}}$
- ▶ for some instance  $\mathbf{x}$ , the respective scores are

$$[-1, 2.2, 1.5]$$

so our guess is orange

## in scikit-learn

- ▶ scikit-learn includes implementations of both of the methods we have discussed:
  - ▶ `OneVsRestClassifier`
  - ▶ `OneVsOneClassifier`
- ▶ however, the built-in algorithms (e.g. Perceptron, LogisticRegression) will do this automatically for you
  - ▶ they use one-versus-rest

## idea 2: multiclass learning algorithms

- ▶ is it good to separate the multiclass task into smaller tasks that are trained independently?
  - ▶ maybe training should be similar to testing?
- ▶ let's make a model where one-vs-rest is used while training
  - ▶ we'll see how this can be done for logistic regression

## binary LR: reminders

- ▶ the **logistic** or **sigmoid** function:

$$P(\text{positive output}|\mathbf{x}) = \frac{1}{1 + e^{-\text{score}}}$$

```
def sigmoid(score):  
    return 1 / (1 + np.exp(-score))
```

- ▶ when training, we minimize the log-loss

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



## multiclass LR using the softmax

- ▶ the **softmax** function is used in multiclass LR instead of the logistic:

$$P(y_i|\mathbf{x}) = \frac{e^{\text{score}_i}}{\sum_k e^{\text{score}_k}}$$

```
def softmax(scores):  
    expscores = np.exp(scores)  
    return expscores / sum(expscores)
```

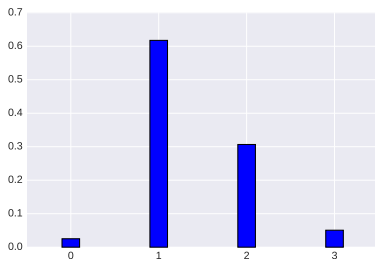
[exercise: make softmax numerically stable]

## softmax example

```
def softmax(scores):  
    expscores = np.exp(scores)  
    return expscores / sum(expscores)
```

```
scores = [-1, 2.2, 1.5, -0.3]  
print(softmax(scores))
```

```
array([ 0.02517067,  0.61750026,  0.30664156,  0.05068751])
```



## cross-entropy loss

- ▶ when training, the softmax probabilities lead to the **cross-entropy** loss instead of the log loss

$$\text{Loss}_{CE}(\mathbf{w}, \mathbf{x}_i, y_i) = -\log P(y_i | \mathbf{x}_i) = -\log \frac{e^{\text{score}_i}}{\sum_k e^{\text{score}_k}}$$

- ▶ just like the log-loss:
  - ▶ high probability for the correct label  $y_i \Rightarrow$  low loss
  - ▶ low probability for  $y_i \Rightarrow$  high loss

## multiclass LR in scikit-learn

- ▶ `LogisticRegression(multi_class='multinomial')`
- ▶ (otherwise, separate classifiers are trained independently)

# Overview

A bit on perceptron

logistic regression

training a logistic regression classifier

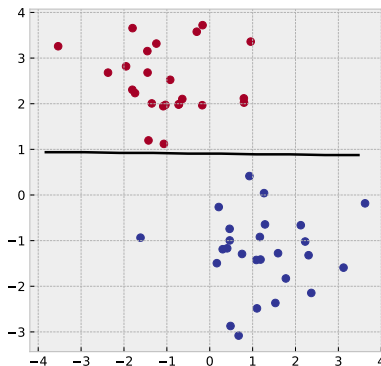
detour: multiclass linear classifiers

**support vector classification**

optimizing the LR and SVM objectives

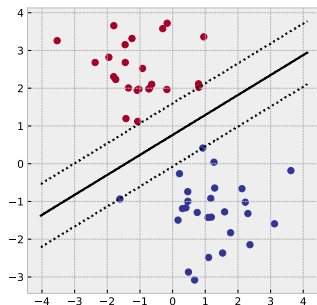
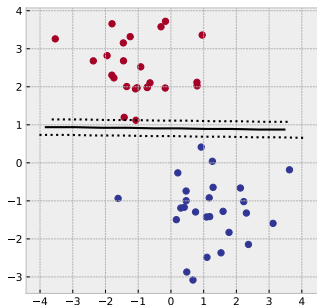
## geometric view

- ▶ using lines to separate 2-dimensional data
- ▶ using planes to separate 3-dimensional data



# margin of separation

- ▶ the **margin**  $\gamma$  denotes how well  $\mathbf{w}$  separates the classes:



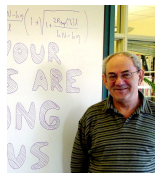
- ▶  $\gamma$  is the shortest distance from the separator to the nearest training instance

large margins are good

- a result from statistical learning theory:

$$\text{true error} \leq \text{training error} + \text{BigUglyFormula}\left(\frac{1}{\gamma^2}\right)$$

- larger margin  $\rightarrow$  better generalization




В. Н. ВАПНИК, А. Я. ЧЕРВОНЕНКИС

## ТЕОРИЯ РАСПОЗНАВАНИЯ ОБРАЗОВ

СТАТИСТИЧЕСКИЕ ПРОБЛЕМЫ ОБУЧЕНИЯ



# Structural risk minimization theorem

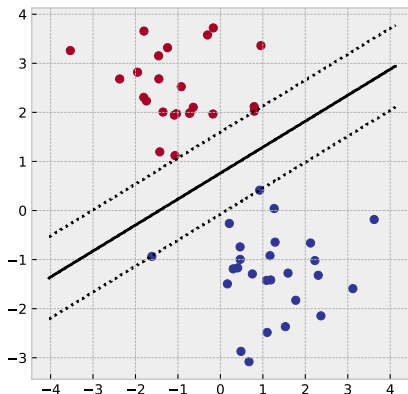

$$\text{error}_{\text{true}}(h) \leq \text{error}_{\text{train}}^{\gamma}(h) + C \sqrt{\frac{\frac{R^2}{\gamma^2} \ln m + \ln \frac{1}{\delta}}{m}}$$

$$\text{error}_{\text{train}}^{\gamma}(h) = \text{num. points with margin} < \gamma$$

[[source](#)]

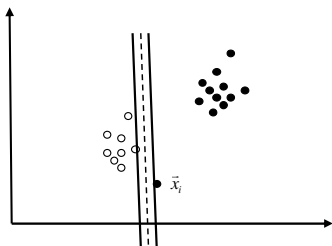
## support vector machines

- **support vector machines** (SVM) or support vector classifiers (SVC) are linear classifiers constructed by selecting the  $\mathbf{w}$  that maximizes the margin

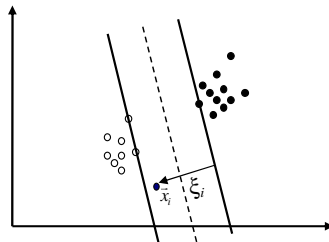


# soft-margin SVMs

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



A) Hard Margin SVM



B) Soft Margin SVM

## stating the SVM as 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 it can be shown (see Daumé's book) that the soft-margin SVM can be stated as minimizing

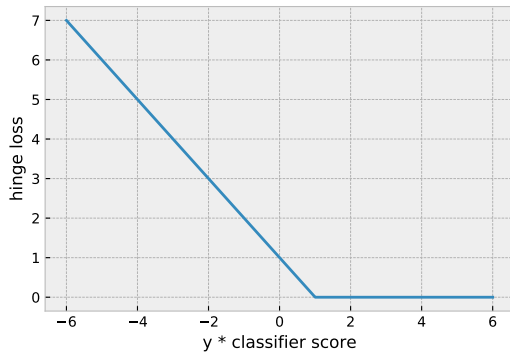
$$\frac{C}{N} \cdot \sum_{i=1}^N \text{Loss}(\mathbf{w}, \mathbf{x}_i, y_i) + \frac{1}{2} \|\mathbf{w}\|^2$$

where

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

is called the **hinge loss**

## plot of the hinge loss



in scikit-learn

- ▶ linear SVM is called `sklearn.svm.LinearSVC`

# Overview

A bit on perceptron

logistic regression

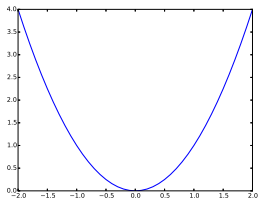
training a logistic regression classifier

detour: multiclass linear classifiers

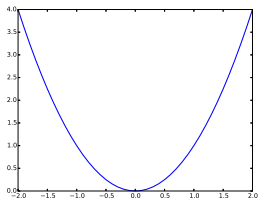
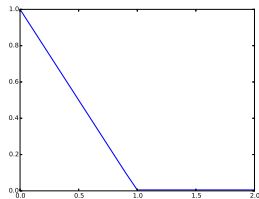
support vector classification

optimizing the LR and SVM objectives

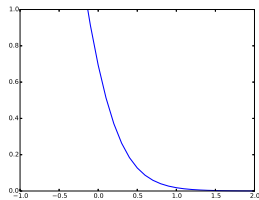
# SVM and LR have convex objective functions



+



+





# optimizing SVM and LR

- ▶ since the objective functions of SVM and LR are convex, we can find  $\mathbf{w}$  by stochastic gradient descent
- ▶ pseudocode:
  - ▶ set  $\mathbf{w}$  to some initial value, e.g. all zero
  - ▶ iterate a fixed number of times:
    - ▶ select a single training instance  $\mathbf{x}$
    - ▶ select a “suitable” step length  $\eta$
    - ▶ compute the gradient of the hinge loss or log loss
    - ▶ subtract step length  $\cdot$  gradient from  $\mathbf{w}$
- ▶ note the similarity to the perceptron!

## missing pieces

- ▶ setting the **learning rate**  $\eta$ 
  - ▶ in principle, one can try to select a “small enough” value of  $\eta$
  - ▶ in practice, it’s better to decrease  $\eta$  gradually
  - ▶ we can use the **Pegasos** algorithm to set  $\eta$  as follows:

$$\eta = \frac{C}{t} = \frac{1}{\lambda \cdot t}$$

where

- ▶  $t$  is the current step (1, 2, ...)
  - ▶  $C$  or  $\lambda$  is the loss/regularization tradeoff
- ▶ **gradients** for SVM and LR loss functions (hinge and log loss)

## Next lecture

- ▶ Gradient Boosting
- ▶ Evaluation methods