

Machine Learning for NLP

Lecture 6: Kernel-based classifiers



**UNIVERSITY OF
GOTHENBURG**

Richard Johansson

October 5, 2015

- ▶ **kernels** give us an interesting connection between **linear** and **example-based** classifiers
 - ▶ a **linear** classifier computes a score for each feature, and then sums the scores
 - ▶ an **example-based** classifier uses a similarity function to compare a new instance to the training examples
 - ▶ informally, kernels are **similarity functions**; formally, they are **dot products** in some transformed vector space
- ▶ we start from the linear classifiers and show that many of them have an alternative example-based form
- ▶ the selling point: get rid of feature engineering, use a similarity function instead

the primal and dual forms: reflection

$$\begin{array}{ll} \text{primal} & \text{dual} \\ \text{score}(\mathbf{x}) = \mathbf{w} \cdot \mathbf{x} & \text{score}(\mathbf{x}) = \sum_i \alpha_i \cdot (\mathbf{x}_i \cdot \mathbf{x}) \end{array}$$

- ▶ the dual form can be seen as some sort of example-based classifier
- ▶ why do we say that the primal and the dual are related?
 - ▶ can we go from the dual to the primal?

perceptron in the primal and dual forms

initialize to all zeros:

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

dual: $\alpha = (0, \dots, 0)$

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

if y_i is positive and $\text{score}(\mathbf{x}_i) \leq 0$

add \mathbf{x}_i to the classifier

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

dual: ?

else if y_i is negative and $\text{score}(\mathbf{x}_i) \geq 0$

subtract \mathbf{x}_i from the classifier

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

dual: ?

return the classifier

perceptron in the primal and dual forms

initialize to all zeros:

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

dual: $\alpha = (0, \dots, 0)$

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

if y_i is positive and $\text{score}(\mathbf{x}_i) \leq 0$

add \mathbf{x}_i to the classifier

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

dual: $\alpha_i = \alpha_i + 1$

else if y_i is negative and $\text{score}(\mathbf{x}_i) \geq 0$

subtract \mathbf{x}_i from the classifier

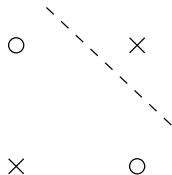
primal: $\mathbf{w} = \mathbf{w} - \mathbf{x}_i$

dual: $\alpha_i = \alpha_i - 1$

return the classifier

a simple example of linear inseparability: an “XOR” situation

<i>very good</i>	Positive
<i>very bad</i>	Negative
<i>not good</i>	Negative
<i>not bad</i>	Positive

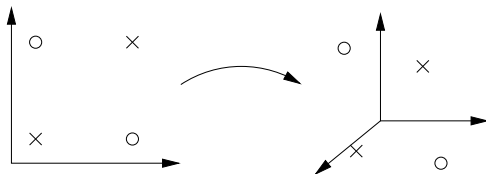


mapping into a larger vector space

- ▶ we may add “useful combinations” of features to make the dataset separable:

<i>very good</i>	very-good	Positive
<i>very bad</i>	very-bad	Negative
<i>not good</i>	not-good	Negative
<i>not bad</i>	not-bad	Positive

- ▶ from a geometrical viewpoint: we are creating a feature space with a higher dimensionality:



- ▶ lots of features → LOTS of combinations

example: XOR dataset converted into 3 dimensions

- ▶ let's apply the function

$$\phi([x_1, x_2]) = [x_1^2, x_2^2, \sqrt{2} \cdot x_1 \cdot x_2]$$

```
X = numpy.array([[1, 1, sqrt(2)*1*1],  
                [1, 0, sqrt(2)*1*0],  
                [0, 1, sqrt(2)*0*1],  
                [0, 0, sqrt(2)*0*0]])  
Y = ['no', 'yes', 'yes', 'no']
```

```
clf = LinearSVC()  
clf.fit(X, Y)
```

```
# in the 3-dimensional space, we get 100% accuracy  
print(accuracy_score(Y, clf.predict(X)))
```

overview

the primal and dual forms

mapping feature vectors into higher-dimensional spaces

kernels in classifiers

- ▶ let's combine the two ideas we've been discussing:
 - ▶ converting examples \mathbf{x} into higher-dimensional vectors $\phi(\mathbf{x})$
 - ▶ using the dual form of the classifiers
- ▶ then we get:

$$\text{score}(\mathbf{x}) = \sum_i \alpha_i \cdot (\phi(\mathbf{x}_i) \cdot \phi(\mathbf{x}))$$

- ▶ we mentioned previously that it seems like a problem that $\phi(\mathbf{x})$ is huge...

the “kernel trick”

- ▶ in the dual form, the feature vectors $\phi(\mathbf{x})_1, \dots, \phi(\mathbf{x})_n$ are used in **dot products only**:

$$\text{score}(\mathbf{x}) = \sum_i \alpha_i \cdot (\phi(\mathbf{x})_i \cdot \phi(\mathbf{x}))$$

- ▶ a **kernel** K is a function that corresponds to a dot product in some transformed vector space

$$\text{score}(\mathbf{x}) = \sum_i \alpha_i \cdot K(\mathbf{x}_i, \mathbf{x})$$

- ▶ the “kernel trick”: in some cases, we may compute the kernel **without actually computing ϕ**

example: quadratic kernel

- ▶ recall the previous example, where we had

$$\phi([x_1, x_2]) = [x_1^2, x_2^2, \sqrt{2} \cdot x_1 \cdot x_2]$$

- ▶ in this case, we can compute the high-dimensional dot product without actually making the high-dimensional vectors:

$$K(\mathbf{a}, \mathbf{b}) = \phi(\mathbf{a}) \cdot \phi(\mathbf{b}) = (\mathbf{a} \cdot \mathbf{b})^2$$

- ▶ this is called a **quadratic kernel**

quadratic kernel: derivation

$$\begin{aligned}\phi([a_1, a_2]) \cdot \phi([b_1, b_2]) &= [a_1^2, a_2^2, \sqrt{2}a_1a_2] \cdot [b_1^2, b_2^2, \sqrt{2}b_1b_2] = \\ &= a_1^2b_1^2 + a_2^2b_2^2 + 2a_1a_2b_1b_2 = \\ &= (a_1b_1 + a_2b_2)^2 = \\ &= ([a_1, a_2] \cdot [b_1, b_2])^2\end{aligned}$$

some common kernels

- ▶ many of the commonly used kernels are just some nonlinear function applied to the normal dot product
- ▶ **polynomial kernel**: $K(\mathbf{a}, \mathbf{b}) = (\mathbf{a} \cdot \mathbf{b} + c_0)^d$
 - ▶ ...where d is called the **degree** and c_0 the **offset**
 - ▶ this category includes the linear and quadratic kernels
 - ▶ in general, a polynomial kernel with degree d will implicitly form all combinations of d features
- ▶ **RBF kernel**: $K(\mathbf{a}, \mathbf{b}) = \exp(-\gamma \cdot |\mathbf{a} - \mathbf{b}|^2)$

decision boundaries: linear and quadratic SVM

- ▶ a kernel-based classifier is linear in the high-dimensional space, but non-linear in the original space
- ▶ example: linear SVM compared to SVM with quadratic kernel

