Introduction to Machine Learning

Nonlinear Support Vector Machines Feature Generation for Nonlinear Separation

× × 0 × × ×



Learning goals

- Understand how nonlinearity can be introduced via feature maps in SVMs
- Know the limitation of feature maps

NONLINEARITY VIA FEATURE MAPS

- How to extend a linear classifier, e.g. the SVM, to nonlinear separation between classes?
- We could project the data from 2D into a richer 3D feature space!





NONLINEARITY VIA FEATURE MAPS / 2

In order to "lift" the data points into a higher dimension, we have to find a suitable **feature map** $\phi : \mathcal{X} \to \Phi$. Let us consider another example where the classes lie on two concentric circles:



× 0 0 × × ×

NONLINEARITY VIA FEATURE MAPS / 3

We apply the feature map $\phi(x_1, x_2) = (x_1, x_2, x_1^2 + x_2^2)$ to map our points into a 3D space. Now our data can be separated by a hyperplane.





NONLINEARITY VIA FEATURE MAPS / 4

The hyperplane learned in $\Phi \subset \mathbb{R}^3$ yields a nonlinear decision boundary when projected back to $\mathcal{X} = \mathbb{R}^2$.





Let us have a look at a similar nonlinear feature map $\phi : \mathbb{R}^2 \to \mathbb{R}^5$, where we collect all monomial feature extractors up to degree 2 (pairwise interactions and quadratic effects):

$$\phi(x_1, x_2) = (x_1^2, x_2^2, x_1 x_2, x_1, x_2).$$

For *p* features vectors, there are k_1 different monomials where the degree is exactly *d*, and k_2 different monomials up to degree *d*.

$$k_1 = \begin{pmatrix} d+p-1\\ d \end{pmatrix}$$
 $k_2 = \begin{pmatrix} d+p\\ d \end{pmatrix} - 1$

Which is quite a lot, if *p* is large.

× 0 0 × × ×

Let us see how well we can classify the 28×28 -pixel images of the handwritten digits of the MNIST dataset (70K observations across 10 classes). We use SVM with a nonlinear feature map which projects the images to a space of all monomials up to the degree *d* and *C* = 1:



0 0 X X 0 X X

For this scenario, with increasing degree *d* the test mmce decreases.

NB: We handle the multiclass task with the "one-against-one" approach. We are somewhat lazy and only use 700 observations to train (rest for testing). We do not do any tuning - as we always should for the SVM!

However, even a 16×16 -pixel input image results in infeasible dimensions for our extracted features (monomials up to degree *d*).



× 0 0 × 0 × ×

In this case, training classifiers like a linear SVM via dataset transformations will incur serious **computational and memory problems**.

Are we at a "dead end"? Answer: No, this is why kernels exist! × × 0 × × ×