

Kernels I

Kernels allow us to make complex, non-linear classifiers using Support Vector Machines.

Given x , compute new feature depending on proximity to landmarks l^1, l^2, l^3 .

To do this, we find the "similarity" of x and some landmark l^i

$$f_i = \text{similarity}(x, l^i) = \exp\left(-\frac{\|x - l^i\|^2}{2\sigma^2}\right)$$

This "similarity" function is called a **Gaussian Kernel**. It is a specific example of a kernel.

The similarity function can also be written as follows:

$$f_i = \text{similarity}(x, l^i) = \exp\left(-\frac{\sum_{j=1}^n (x_j - l_j^i)^2}{2\sigma^2}\right)$$

There are a couple properties of the similarity function:

$$\begin{aligned} \text{if } x \approx l^i, \text{ then } f_i &= \exp\left(-\frac{\approx 0^2}{2\sigma^2}\right) \approx 1 \\ \text{if } x \text{ is far from } l^i, \text{ then } f_i &= \exp\left(-\frac{(\text{large number})^2}{2\sigma^2}\right) \approx 0 \end{aligned}$$

In other words, if x and the landmark are close, then the similarity will be close to 1, and if x and the landmark are far away from each other, the similarity will be close to 0.

Each landmark gives us the features in our hypothesis:

$$\begin{aligned} l^1 &\rightarrow f_1 \\ l^2 &\rightarrow f_2 \\ l^3 &\rightarrow f_3 \\ &\dots \\ h_{\theta}(x) &= \theta_1 f_1 + \theta_2 f_2 + \dots \end{aligned}$$

σ^2 is a parameter of the Gaussian Kernel, and it can be modified to increase or decrease the **drop-off** of our feature f_i . Combined with looking at the values inside Θ , we can choose these landmarks to get the general shape of the decision boundary.

Kernels II

One way to get the landmarks is to put them in the **exact same** locations as all the training examples. This gives us m landmarks, with one landmark per training example.

Given example x :

$f_1 = \text{similarity}(x, l^1)$, $f_2 = \text{similarity}(x, l^2)$, $f_3 = \text{similarity}(x, l^3)$, and so on.

This gives us a "feature vector," f_i of all our features for example x_i . We may also set $f_0 = 1$ to correspond with θ_0 . Thus given training example x_i :

$$x^i \rightarrow [f_1^i = \text{similarity}(x^i, l^1), f_2 = \text{similarity}(x^i, l^2), f_3 = \text{similarity}(x^i, l^3)]$$

Now to get the parameters Θ we can use the SVM minimization algorithm but with f_i substituted in for x_i :

$$J(\theta) = C \sum_{i=1}^m [y^i \text{cost}_1(\theta^T f^i) + (1 - y^i)(\text{cost}_0(\theta^T f^i))] + \frac{1}{2} \sum_{j=1}^n \theta_j^2$$

Using kernels to generate $f(i)$ is not exclusive to SVMs and may also be applied to logistic regression. However, because of computational optimizations on SVMs, kernels combined with SVMs is much faster than with other algorithms, so kernels are almost always found combined only with SVMs.