

# 2/9/2023: Kernels (Part 2)

inputs to predictor

$K(x, z)$  measures similarity between  $x$  &  $z$   
- Similar  $x$  &  $z$  should have large  $K(x, z)$   
Kernel function

Logistic Regression is already computing prediction given input  $x$  based on

$$\sum_{i=1}^n a_i k(x, x^{(i)})$$

If  $> 0$  predict +1  
If  $< 0$  predict -1

If we define  $k(x, z) = x^T z$

For each training example, compute its similarity to  $x$  & multiply by a learned weight  $a_i$

## Original Logistic Regression

1) Training:  $w^{(t)} \leftarrow w^{(t-1)} + \eta \sum_{i=1}^n \delta(-y^{(i)}) w^{(t-1)T} x^{(i)} y^{(i)}$

For example:  $+ 0.7 x^{(1)} - 0.3 x^{(2)} + 0.2 x^{(3)}$   
( $n=3$ )

2) Testing: we compute  $w^T x$

## Kernel Logistic Regression (Equivalent mathematically)

Define  $a \in \mathbb{R}^n$ ,  $w = \sum_{i=1}^n a_i x^{(i)}$

1) Training:  $a_i^{(t)} \leftarrow a_i^{(t-1)} + \eta \cdot \delta(-y^{(i)}) w^{(t-1)T} x^{(i)} y^{(i)}$

For example:

$$\begin{aligned} a_1^{(t)} &\leftarrow a_1^{(t-1)} + 0.7 \\ a_2^{(t)} &\leftarrow a_2^{(t-1)} - 0.3 \\ a_3^{(t)} &\leftarrow a_3^{(t-1)} + 0.2 \end{aligned}$$

$$= \sum_{j=1}^n a_j x^{(j)T} x^{(i)}$$

$$= a_i^{(t-1)} + \eta \cdot \delta(-y^{(i)}) \left( \sum_{j=1}^n a_j^{(t-1)} k(x^{(j)}, x^{(i)}) \right) y^{(i)}$$

only place we see  $x$ 's.

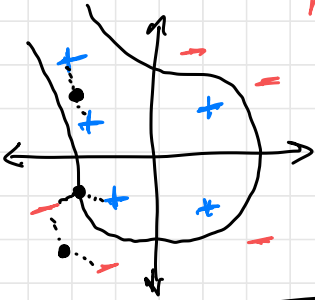
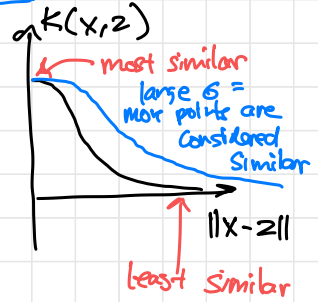
2) Similarly, testing: compute  $\sum_{j=1}^n a_j k(x^{(j)}, x)$   
 only place  $x$  is used

Big payoff: we can run kernelized algorithm with any choice of kernel functions

### Radial Basis Function (RBF)

$$K(x, z) = \exp\left(-\frac{\|x-z\|^2}{2\sigma^2}\right)$$

↑ hyperparameter



In practice, RBF is a very popular way to learn non-linear decision boundary.

### Kernels & Features

| $y$ | $x_1$ | $x_2$ | transform each row | $y$ | $\sqrt{2}x_1$ | $\sqrt{2}x_2$ | 1 | $x_1^2$ | $x_2^2$ | $\sqrt{2}x_1x_2$ |
|-----|-------|-------|--------------------|-----|---------------|---------------|---|---------|---------|------------------|
| +1  | 2     | 3     | →                  | +1  | $2\sqrt{2}$   | $3\sqrt{2}$   | 1 | 4       | 9       | $6\sqrt{2}$      |
| -1  | 0     | 1     |                    | -1  | 0             | $1\sqrt{2}$   | 1 | 0       | 1       | 0                |

Call this transformation  $\phi: \mathbb{R}^2 \rightarrow \mathbb{R}^6$

We can run normal logistic regression with  $\phi$ , but it would take 3x long

Idea: "Kernel trick"

Use a kernelized algorithm, so

$$K(x, z) = \phi(x)^T \phi(z)$$

In many cases, can compute this directly (i.e. without creating  $\phi(x)$  and  $\phi(z)$ )

Polynomial Kernel for degree 2 (Quadratic Kernel)

$$K(x, z) = (x^T z + 1)^2 = \phi(x)^T \phi(z)$$

when  $\phi(x) = \begin{bmatrix} 1 \\ \sqrt{2} x_1 \\ \sqrt{2} x_2 \\ x_1^2 \\ x_2^2 \\ \sqrt{2} x_1 x_2 \end{bmatrix}$

You can compute  $\phi(x)^T \phi(z)$  without "running"  $\phi$

In general, for any original dimension of  $x$ 's  
for any degree  $p$  hyperparameter

$$K(x, z) = (x^T z + 1)^p = \phi(x)^T \phi(z)$$

for some  $\phi$  that has all monomials of degree  $\leq p$

What about RBF?

Fact:

$$\exp\left(-\frac{\|x-z\|^2}{2\sigma^2}\right) = \phi(x)^T \phi(z)$$

for some  $\phi(x)$  that is infinite dimensional.

Kernel Logistic regression + RBF (doable)  
equivalent to

Logistic regression +  $d$ -dimensional features (not doable)

Practical Considerations: Polynomial kernel, degree  $p$ , original dimension is  $d$

Original:  $T$  iterations, each  $O(n \cdot d^p)$   
 $\uparrow$  size of  $\phi(x)$

Kernel:  $T$  iterations, each  $O(n^2 \cdot d)$  computing  $K(x, z)$

Kernel pays  $O(n^2)$ , but enables using more features at no additional cost