

inputs  
↓ ↓  
 $k(x, z)$  measures similarity between  $x$  and  $z$

kernel function

high similarity  $\Rightarrow$  large  $k(x, z)$

Logistic regression is already computing prediction given  $x$  based on

$$\sum_{i=1}^n a_i k(x, x^{(i)}) \begin{cases} \nearrow & \text{if } > 0, \text{ predict } +1 \\ \searrow & \text{if } < 0, \text{ predict } -1 \end{cases}$$

if we define  $k(x, x^{(i)}) = w^T z$

→ for each training example, compute its similarity to  $x$  and multiply by a learned weight  $a_i$

original logistic regression

1) Training:  $w^{(t)} \leftarrow w^{(t-1)} + \eta \underbrace{\sum_{i=1}^n \sigma(-y^{(i)} w^{(t-1)T} x^{(i)}) y^{(i)} x^{(i)}}_{\text{scalar}}$

2) Testing: 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:  $a_i^{(t)} \leftarrow a_i^{(t-1)} + \eta \cdot \sigma(-y^{(i)} w^{(t-1)T} x^{(i)}) \cdot y^{(i)}$   
 $= \sum_{j=1}^n a_j x^{(j)T} x^{(i)}$

$= a_i^{(t-1)} + \eta \cdot \sigma(-y^{(i)} \sum_{j=1}^n a_j \underbrace{k(x^{(j)}, x^{(i)})}_{\text{only place we use } x\text{'s}}) \cdot y^{(i)}$

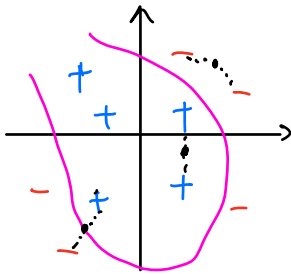
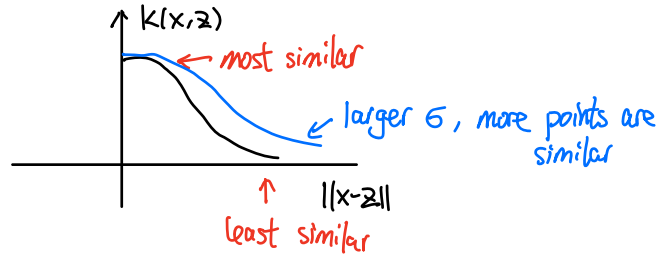
2) Testing: compute  $\sum_{j=1}^n a_j k(x^{(j)}, x)$

We can run the algorithm with any choice of the kernel function.

### 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 vs Features

$y$	$x_1$	$x_2$		$y$	$x_1$	$x_2$	1	$x_1^2$	$x_2^2$	$x_1 x_2$
+1	2	3	transform each row $\rightarrow$	+1	2	3	1	4	9	6
-1	0	1		-1	0	1	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  $\gg x$  long.

"kernel trick"

Use a kernelized algorithm so

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

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

Polynomial kernel for degree 2 (Quadratic kernel)

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

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

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

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

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

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

What about RBF?

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

for some  $\phi(x)$  that is infinite dimension

kernel logistic regression + RBF (doable)



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

Runtime: polynomial kernel degree  $p$

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

kernel:  $T$  iterations, each  $O(n^2 \cdot d)$   
↳ computing  $k(x, z)$

kernel pay  $O(n^2)$ , but enables using more features at no additional costs