#### Kernels

#### Perceptron in action



$\mathbf{w} = 0$	
while some (x,y)	is
misclassified:	
w = w + yx	

## When does this fail?

When data is not *linearly separable*.



[2] Noise





## Systematic inseparability



Actual decision boundary is *quadratic*.

Quick fix: in addition to regular features

 $\mathbf{x} = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_d),$ add in extra features  $\mathbf{x}_1^2, \mathbf{x}_2^2, \dots, \mathbf{x}_d^2,$ 

 $\mathbf{x}_1 \mathbf{x}_2, \mathbf{x}_1 \mathbf{x}_3, \dots, \mathbf{x}_{d-1} \mathbf{x}_d$ 

The new, enhanced data vectors are of the form  $\Phi(\mathbf{x}) =$  $(1, \sqrt{2x_1}, ..., \sqrt{2x_d}, \mathbf{x_1}^2, ..., \mathbf{x_d}^2, \sqrt{2x_1x_2}, \sqrt{2x_1x_3}, ..., \sqrt{2x_{d-1}x_d})$ 

## Adding new features



Boundary is something like  $\mathbf{x}_1 = \mathbf{x}_2^2 + 5$ . This is *quadratic* in  $\mathbf{x} = (\mathbf{1}, \mathbf{x}_1, \mathbf{x}_2)$ but *linear* in  $\Phi(\mathbf{x}) = (\mathbf{1}, \sqrt{2\mathbf{x}_1}, \sqrt{2\mathbf{x}_2}, \mathbf{x}_1^2, \mathbf{x}_2^2, \sqrt{2\mathbf{x}_1\mathbf{x}_2})$ 

By embedding the data in a higher-dimensional feature space, we can keep using a linear classifier!

## Perceptron revisited

Learning in the higher-dimensional feature space: map each  $\mathbf{x}$  onto  $\Phi$  ( $\mathbf{x}$ ) and then run the regular perceptron; that is:

```
w = 0
while some y(w \cdot \Phi(x)) \leq 0:
w = w + y \Phi(x)
```

Everything works as before; final w is a weighted sum of various  $\Phi(x)$ .

**Problem:** number of features has now increased dramatically. For instance, if  $\mathbf{x} \in \mathbb{R}^{1,000}$  then  $\Phi(\mathbf{x}) \in \mathbb{R}^{1,000,000}$ !

## The kernel trick

[Aizenman, Braverman, Rozonoer, 1964] No need to explicitly write out either **w** or  $\Phi(\mathbf{x})$ !

[1] Keep w in sparse form If  $\mathbf{w} = (\mathbf{say}) \mathbf{a}_1 \Phi(\mathbf{x}^{(1)}) + \mathbf{a}_2 \Phi(\mathbf{x}^{(20)}) + \mathbf{a}_3 \Phi(\mathbf{x}^{(31)})$ , store it as a list, [(1,a<sub>1</sub>), (20,a<sub>2</sub>), (31,a<sub>3</sub>)].

#### [2] When do we ever access $\Phi(\mathbf{x})$ ?

Only to compute a dot product  $\mathbf{w} \cdot \Phi(\mathbf{x})$  (during training or future evaluation). In above example  $\mathbf{w} \cdot \Phi(\mathbf{x})$  is  $\mathbf{a}_1 \Phi(\mathbf{x}^{(1)}) \cdot \Phi(\mathbf{x}) + \mathbf{a}_2 \Phi(\mathbf{x}^{(20)}) \cdot \Phi(\mathbf{x}) + \mathbf{a}_3 \Phi(\mathbf{x}^{(31)}) \cdot \Phi(\mathbf{x})$ 

In general,  $\mathbf{w} \cdot \Phi(\mathbf{x})$  is a (weighted) sum of dot products  $\Phi(\mathbf{x}^{(i)}) \cdot \Phi(\mathbf{x})$ .

Can we compute such dot products without writing out the  $\Phi(\mathbf{x}) ' \mathbf{s}$ ?

## Kernel trick, cont'd

In 2-d:  

$$\Phi(\mathbf{x}) \cdot \Phi(\mathbf{z}) = (1, \sqrt{2x_1}, \sqrt{2x_2}, \mathbf{x_1}^2, \mathbf{x_2}^2, \sqrt{2x_1x_2}) \cdot (1, \sqrt{2z_1}, \sqrt{2z_2}, \mathbf{z_1}^2, \mathbf{z_2}^2, \sqrt{2z_1z_2}) = 1 + 2\mathbf{x}_1\mathbf{z}_1 + 2\mathbf{x}_2\mathbf{z}_2 + \mathbf{x}_1^2\mathbf{z}_1^2 + \mathbf{x}_2^2\mathbf{z}_2^2 + 2\mathbf{x}_1\mathbf{x}_2\mathbf{z}_1\mathbf{z}_2$$

$$= (1 + \mathbf{x}_1\mathbf{z}_1 + \mathbf{x}_2\mathbf{z}_2)^2$$

$$= (1 + \mathbf{x} \cdot \mathbf{z})^2$$

In d dimensions:

$$\Phi(\mathbf{x}) \cdot \Phi(\mathbf{z}) = (1, \sqrt{2x_1, ..., \sqrt{2x_d}}, \mathbf{x_1}^2, ..., \mathbf{x_d}^2, \sqrt{2x_1x_2}, \sqrt{2x_1x_3, ..., \sqrt{2x_{d-1}x_d}}) \cdot (1, \sqrt{2z_1, ..., \sqrt{2z_d}}, \mathbf{z_1}^2, ..., \mathbf{z_d}^2, \sqrt{2z_1z_2}, \sqrt{2z_1z_3, ..., \sqrt{2z_{d-1}z_d}}) = (1 + \mathbf{x_1}\mathbf{z_1} + \mathbf{x_2}\mathbf{z_2} + ... + \mathbf{x_d}\mathbf{z_d})^2 = (1 + \mathbf{x} \cdot \mathbf{z})^2$$

Computing dot products in the 1,000,000-dimensional feature space takes time proportional to just 1000, the original dimension!

Never need to write out  $\Phi(\mathbf{x})$  !

## Kernel trick

Why does it work?

- 1. The only time we ever use the data is to compute dot products  $\mathbf{w} \cdot \Phi(\mathbf{x})$ .
- 2. And w itself is a linear combination of  $\Phi(x)$ 's. If  $w = a_1 \Phi(x^{(1)}) + a_{22} \Phi(x^{(22)}) + a_{37} \Phi(x^{(37)})$ store it as  $[(1, a_1), (22, a_{22}), (37, a_{37})]$
- 3. Dot products  $\Phi(\mathbf{x}) \cdot \Phi(\mathbf{z})$  can be computed very efficiently.

## Quadratic kernel



## Quadratic kernel



## Polynomial decision surfaces

To get a decision surface which is a polynomial of order d:



Let  $\Phi(\mathbf{x})$  consist of all terms of order  $\leq \mathbf{d}$ , e.g.  $\mathbf{x}_1 \mathbf{x}_2 \mathbf{x}_3^{d-3}$ If  $\mathbf{x}$  has p coords, then  $\Phi(\mathbf{x})$  has about p<sup>d</sup> coords.

Same trick works:  $\Phi(\mathbf{x}) \cdot \Phi(\mathbf{z}) = (\mathbf{1} + \mathbf{x} \cdot \mathbf{z})^d$ !

The data is accessed only through the *kernel function*   $\mathbf{k}(\mathbf{x}, \mathbf{z}) = \Phi(\mathbf{x}) \cdot \Phi(\mathbf{z})$ . This is a measure of similarity.

# String kernels

Sequence data: eg. text documents speech signals protein sequences Each data point is a *sequence*.

```
Input space
X = {A,C,G,T}*
X = {English words}*
Different data points have different lengths.
```

## String kernels

#### For each substring **s**, define *feature* $\Phi_s(\mathbf{x}) = \# \text{ of times substring s}$ appears in $\mathbf{x}$

and

 $\Phi(\mathbf{x}) = (\Phi_s(\mathbf{x}): \text{ all strings } \mathbf{s})$ (vector with a coordinate for each string  $\mathbf{s}$  in the alphabet) Infinite-dimensional embedding!

Eg.

 $\Phi_{ar}$ (aardvark) = 2  $\Phi_{th}$ (aardvark) = 0

Linear classifier with such features is potentially very powerful.

## String kernels

We can compute dot products quickly!

To compute k(x,z) = Φ(x) · Φ(z): for each substring s of x count how often s appears in z

With dynamic programming, this just takes time proportional to product of the lengths of **x**, **z**.

*Postscript*: kernels over other interesting input spaces: parse trees, unordered collections of objects,...

## **Kernel function**

Now, shift attention:

away from embedding Φ (x)
[which we never explicitly construct anyway],
towards the *similarity measure* k (x, z)
[the thing we actually use].

Let's rewrite everything in terms of **k**. Suppose the final linear separator is

$$w = a_1 y_1 \Phi(x^{(1)}) + ... + a_m y_m \Phi(x^{(m)}).$$

Then the final classifier can also be written  $F(x) = sgn\{a_1y_1k(x^{(1)}, x) + ... + a_my_mk(x^{(m)}, x)\},\$ a weighted vote.

## **Kernel function**

As one varies  $\Phi$ , what kinds of similarity measures **k** are possible?

Any **k** satisfying a technical condition ("positive definiteness") corresponds to some embedding  $\Phi$  (x).

So: don't worry about  $\Phi$  and just pick a similarity measure **k** which suits the data at hand.

Popular choice: Gaussian kernel.  $\mathbf{k}(\mathbf{x}, \mathbf{x}') = \exp(-||\mathbf{x} - \mathbf{x}'||^2/s^2)$ The resulting classifiers  $\mathbf{F}(\mathbf{x}) = \operatorname{sgn}\{a_1y_1k(\mathbf{x}^{(1)}, \mathbf{x}) + \dots + a_my_mk(\mathbf{x}^{(m)}, \mathbf{x})\}$ are closely related to nearest neighbors.

## Gaussian kernel



## Gaussian kernel



## Generalization

With a powerful kernel: can fit any data set. Danger: *overfitting* – the resulting classifier might not perform well on future instances.

e.g. given 100 data points in **R**<sup>101</sup>. Can always find a linear classifier to fit them. But prediction on a future point is arbitrary.



All this becomes a lot more manageable if there is a substantial *margin* between the two classes.

Performance of the perceptron depends only on the margin m of the classifier, not the dimension! [Number of updates =  $1/m^2$ ]

So even infinite dimensions are fine!