

(1)

Spring 2013.

Note Title

## Kernel Trick

11/19/2006

Note that the final classifier depends on  $\underline{x}$  only by dot products.

EG.  $\underline{x} \cdot \underline{x}_\mu$  in final classifier.

$\underline{x}_\mu \cdot \underline{x}_\nu$  in the dual energy.

This motivates the Kernel Trick

Compute features  $\underline{\varphi}(\underline{x})$  and reformulate the problem in feature space.

→ i.e. seek a classifier of form

$$\text{Sign}(\underline{c} \cdot \underline{\varphi}(\underline{x}) + b).$$

Replace  $\underline{x}$  by  $\underline{\varphi}(\underline{x})$  everywhere in the primal & dual formulation.

Then the classifier only depends on the dot product of the  $\underline{\varphi}(\underline{x})$ 's.

On the kernel  $k(\underline{x}, \underline{x}') = \underline{\varphi}(\underline{x}) \cdot \underline{\varphi}(\underline{x}')$ .

(2)

Spring 2013.

Why does this help?

First, features can make it possible to classify data by hyperplanes.

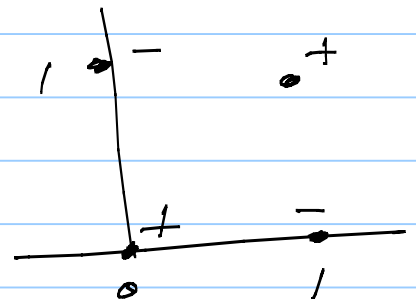
Example Logical X-OR,  $\underline{x} = (x_1, x_2)$   $x_j \in \{\pm 1\}$   
 $\omega_j \in \{\pm 1\}$

The X-OR (exclusive or) requires a decision rule  $\alpha(\underline{x})$  s.t.

$$\alpha(1, 1) = \alpha(-1, -1) = 1$$

$$\alpha(1, -1) = \alpha(-1, 1) = -1$$

Impossible to find a linear classifier to do this.



But define features  $\underline{\varphi}(x_1, x_2) = (x_1, x_2, x_1 x_2)$ .

Now the classifier  $\text{sign} \left\{ (0, 0, 1) \cdot \underline{\varphi}(x_1, x_2) \right\}$  can separate the data.

Moral: increasing the dimensionality of the data, by features, makes it possible to find separating hyperplanes

Spring 2013

(3)

Second, we do not need to specify the features  $\underline{\varphi}(x)$  explicitly, we only need to specify the kernel

$$K(\underline{x}, \underline{x}') = \underline{\varphi}(\underline{x}) \cdot \underline{\varphi}(\underline{x}')$$

Remember:

the dual problem reduces to

$$\begin{aligned} \text{maximizing } L_d(\underline{\alpha}) &= \sum_{\mu} \alpha_{\mu} - \frac{1}{2} \sum_{\mu, \nu} \alpha_{\mu} \alpha_{\nu} \omega_{\mu} \omega_{\nu} \underline{\varphi}(\underline{x}_{\mu}) \cdot \underline{\varphi}(\underline{x}_{\nu}) \\ &= \sum_{\mu} \alpha_{\mu} - \frac{1}{2} \sum_{\mu, \nu} \alpha_{\mu} \alpha_{\nu} \omega_{\mu} \omega_{\nu} K(\underline{x}_{\mu}, \underline{x}_{\nu}) \end{aligned}$$

$$\text{The solution } \underline{\hat{a}} = \sum_{\mu} \hat{\alpha}_{\mu} \omega_{\mu} \underline{\varphi}(\underline{x}_{\mu})$$

$$\begin{aligned} \underline{\hat{a}} \cdot \underline{\varphi}(\underline{x}) &= \sum_{\mu} \hat{\alpha}_{\mu} \omega_{\mu} \underline{\varphi}(\underline{x}) \cdot \underline{\varphi}(\underline{x}_{\mu}) \\ &= \sum_{\mu} \hat{\alpha}_{\mu} \omega_{\mu} K(\underline{x}, \underline{x}_{\mu}). \end{aligned}$$

(can solve for  $\hat{\alpha}$  as before).

(4.)

Spring 2013

What Kernels to Use?

There are many choices of kernels. The difficulty is knowing which one to use. As always, cross-validation is useful for checking whether a kernel can generalize

$$K(\underline{x}, \underline{x}') = (1 + \underline{x} \cdot \underline{x}')^d$$

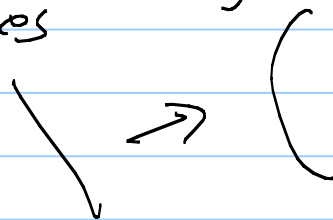
$$K(\underline{x}, \underline{x}') = e^{-\frac{1}{2\sigma^2} \|\underline{x} - \underline{x}'\|^2}$$

$$K(\underline{x}, \underline{x}') = \tanh\{C_1 \underline{x} \cdot \underline{x}' + C_2\}$$

Choice of best kernel is problem dependent.

Some kernels  $\rightarrow$  eg.  $(1 + \underline{x} \cdot \underline{x}')^d$  naturally generalize the idea of hyperplanes

Others  $\rightarrow$  eg.  $e^{-\frac{1}{2\sigma^2} \|\underline{x} - \underline{x}'\|^2}$  are similar to nearest neighbours.



(5)

Spring 2013,

When do kernels correspond to features?

i.e. if we specify  $K(\underline{x}, \underline{x}')$ , is it equal to  $\underline{\phi}(\underline{x}) \cdot \underline{\phi}(\underline{x}')$  for some features  $\underline{\phi}(\underline{x})$ ?

Theoretical results can be obtained, e.g. Mercer's Theorem.

Compute eigenfunctions of  $K(\underline{x}, \underline{x}')$

$$\int K(\underline{x}, \underline{x}') \psi(\underline{x}') d\underline{x}' = \lambda \psi(\underline{x})$$

with  $\int (\psi(\underline{x}))^2 d\underline{x} = 1$ .

Provided  $K(\underline{x}, \underline{x}')$  is positive definite, then the features are  $\underline{\phi}^M(\underline{x}) = \sqrt{\lambda_\mu} \psi_\mu(\underline{x})$ .

Similar to linear algebra expansion of a symmetric matrix in terms of eigenvalues

$$A_{ij} = \sum_{\mu} \lambda_{\mu} e_{i}^{\mu} e_{j}^{\mu}, \quad \text{where } \sum_j A_{ij} e_j^{\mu} = \lambda_{\mu} e_i^{\mu}.$$

If  $A_{ij}$  is the kernel.

$$A_{ij} = \sum_{\mu} (\lambda_{\mu}^{1/2} e_i^{\mu}) (\lambda_{\mu}^{1/2} e_j^{\mu}) = \sum_{\mu} \phi_i^{\mu} \phi_j^{\mu}$$

16)

Spring 2013,

## Kernel PCA

The kernel trick can be applied to an quadratic problem - e.g. PCA.

$$\underline{C} = \frac{1}{m} \sum_{k=1}^m (\underline{x}_k - \bar{\underline{x}}) (\underline{x}_k - \bar{\underline{x}})^T$$

w.l.o.g.  $\bar{\underline{x}} = \frac{1}{m} \sum_{k=1}^m \underline{x}_k = 0.$

Go to feature space  
 $\underline{x} \rightarrow \underline{\phi}(\underline{x})$

$$\rightarrow \underline{C} = \frac{1}{m} \sum_{k=1}^m \underline{\phi}(\underline{x}_k) \underline{\phi}(\underline{x}_k)^T$$

All non-zero eigenvalues  $\underline{e}$  of  $\underline{C}$  are of form

$$\underline{e} = \sum_{j=1}^m \alpha_j \underline{\phi}(\underline{x}_j), \text{ for some } \{\alpha_j\},$$

Substituting:  $\underline{C} \underline{e} = \lambda \underline{e}$

$$\rightarrow \frac{1}{m} \sum_{k=1}^m \underline{\phi}(\underline{x}_k) (\underline{\phi}(\underline{x}_k) \cdot \underline{e}) = \lambda \underline{e}$$

$$\rightarrow \frac{1}{m} \sum_{k=1}^m \underline{\phi}(\underline{x}_k) \sum_{j=1}^m \alpha_j (\underline{\phi}(\underline{x}_k) \cdot \underline{\phi}(\underline{x}_j)) = \lambda \sum_{j=1}^m \alpha_j \underline{\phi}(\underline{x}_j)$$

Equating coeffs of  $\underline{\phi}(\underline{x}_j)$  gives new eigenvalue equations

$$\frac{1}{m} \sum_j K(\underline{x}_k, \underline{x}_j) \alpha_j = \lambda \alpha_k \quad \left\| \begin{array}{l} \text{Index } \lambda^M \\ \alpha_k^M \end{array} \right.$$

(7)

Spring 2013,

$$\forall \mu \sum_j K(x_i, x_j) \alpha_j^\mu = \lambda^\mu \alpha_i^\mu$$

$$\mu = 1, \dots, m.$$

Solving this, gives us the eigenvectors

$$\underline{e}^\mu = \sum_{j=1}^m \alpha_j^\mu \underline{\phi}(x_j), \text{ eigenvalue } \lambda^\mu. \\ (\text{depends on } \phi)$$

But the projections of the data are

$$\underline{e}^\mu \cdot \underline{\phi}(x) = \sum_{j=1}^m \alpha_j^\mu K(x_j, x)$$

which is independent of  $\phi$ .  
(depends only on  $K$ ).

Hence:

the projection of the data onto the eigenvectors requires only knowing the kernel  $K(x_i, x_j)$  (i.e. not knowing  $\phi$ )

Knowledge of the kernel is used twice:  
(1) to compute the  $\{\alpha_j^\mu\}$ ,  
(2) to compute the projections  $\underline{e}^\mu \cdot \underline{\phi}(x)$ . //