

Q1.

## Kernel PCA.

Spring 2014

Note Title

5/11/2014

Now we show that the kernel trick can be applied to PCA.

In this lecture we set the means to zero  $\rightarrow \text{e.g. } \sum_{i=1}^N x_i = 0$   
this is easy to arrange (by subtraction).  $\sum_{i=1}^N \phi(x_i) = 0$

The output of PCA is the projection of the data  $\{x_i\}_{i=1}^N$  onto the subspace defined by  $\{e_1, \dots, e_M\}$   
i.e. the coefficients  $(x_i \cdot e_1, x_i \cdot e_2, \dots, x_i \cdot e_M)$ .  $i = 1 \text{ to } N$

This depends on dot products, which suggests we can use the kernel trick if we replace  $x$  by  $\phi(x)$ .

But what about the  $e$ 's? They are eigenvectors of the correlation function. How do they change if we replace  $x$  by  $\phi(x)$ ?

To understand this, we need another way to

compute the eigenvectors of  $C = \frac{1}{N} \sum_{i=1}^N x_i x_i^T$

Claim. The eigenvectors with non-zero eigenvalues can be expressed in form  $e = \sum_{i=1}^N \alpha_i x_i$ . The eigenvectors with zero eigenvalues are of form  $e$ , st.  $x_i \cdot e = 0, i = 1 \text{ to } N$ .

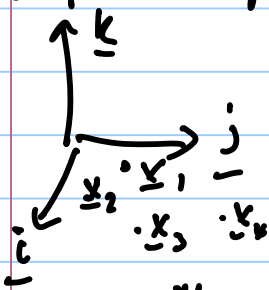
Proof. Suppose  $x_i \cdot e = 0, i = 1 \text{ to } N$

then  $C e = \frac{1}{N} \sum_{i=1}^N x_i (x_i \cdot e) = 0$ , hence  $e$  is an eigenvector with zero eigenvalue. This proves the last sentence.

The remaining eigenvectors must be orthogonal to the zero eigenvectors. Hence they are of form  $\sum_{i=1}^N \alpha_i e_i$ .

(2) Intuition: Suppose we have data in 3-dim space spanned by axes  $\underline{i}, \underline{j}, \underline{k}$ . Suppose all the data lies in the plane spanned by  $\underline{i}, \underline{j}$

11/19/2006



i.e.  $\underline{x}_1 = x_{1i} \underline{i} + y_{1j} \underline{j}$

$\underline{x}_2 = x_{2i} \underline{i} + y_{2j} \underline{j}$

$\underline{x}_n = x_{ni} \underline{i} + y_{nj} \underline{j}$

Then all vectors in the  $\underline{i}, \underline{j}$  plane can be

written as  $\underline{x} = \sum_{i=1}^N \alpha_i \underline{x}_i$ .

The eigenvectors of the correlation function lie in the  $\underline{i}, \underline{j}$  plane, except for the zero eigenvalue with eigenvector  $\underline{e} = \underline{k}$

Now replace  $\underline{x}$  by  $\underline{\phi}(\underline{x})$

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

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

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

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

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

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

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

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

(3)

$$\lambda_N \sum_j K(x_i, x_j) \alpha_j^\mu = \lambda^\mu \alpha_i^\mu \quad \mu = 1, \dots, p.$$

Solving this, gives us the eigenvectors

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

But the projections of the data are  $\underline{e}^\mu \cdot \underline{\phi}(x)$

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

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  $\underline{e}^\mu \cdot \underline{\phi}(x)$ . //