

---

# Kernels

Lyle Ungar

# What is a kernel?

---

- **$k(\mathbf{x}, \mathbf{y})$** 
  - Measures the *similarity* between a pair of points  $\mathbf{x}$  and  $\mathbf{y}$
- **Required properties**
  - Symmetric and positive semi-definite (PSD)
  - Often tested using a *Kernel Matrix*,
    - a matrix  $K$  with elements  $K_{ij} = k(\mathbf{x}_i, \mathbf{x}_j)$   
where the points are all pairs of rows of a matrix  $X$  of predictors
  - A *PSD matrix* has only non-negative singular values
- **Uses**
  - Anywhere you want to replace inner products  $\mathbf{x}^T \mathbf{y}$  with inner products of  $\phi(\mathbf{x})^T \phi(\mathbf{y})$

# Properties of a Kernel

---

Definition: A finitely positive semi-definite function  $k : x \times y \rightarrow R$  is a symmetric function of its arguments for which matrices formed by restriction on any finite subset of points is positive semi-definite.

$$\alpha^T K \alpha \geq 0 \quad \forall \alpha$$

Theorem: A function  $k : x \times y \rightarrow R$  can be written as  $k(x, y) = \langle \Phi(x), \Phi(y) \rangle$  where  $\Phi(x)$  is a feature map  $x \rightarrow \Phi(x) \in F$  iff  $k(x, y)$  satisfies the semi-definiteness property.

Relevance: We can now check if  $k(x, y)$  is a proper kernel using only properties of  $k(x, y)$  itself,

i.e. without the need to know the feature map!

# Where are kernels used?

---

- **Nearest neighbors**
  - Measure similarity in the kernel space
- **Linear regression**
  - Map points to new, transformed feature space
- **Logistic regression**
  - Map points to new, transformed feature space
- **SVMs and Perceptrons**
- **PCA**
  - $\text{SVD}[X^T X]$
- **CCA**
  - $\text{SVD}[(X^T X)^{-1/2} (X^T Y) (Y^T Y)^{-1/2}]$

What is the most common kernel method for linear regression?

What are we seeking to accomplish with kernels for classification?

What is the main benefit for PCA?

When are kernels most necessary for CCA?

# Kernels form a dual representation

---

- Start with an  $n \times p$  matrix  $X$  of predictors
- Generate an  $n \times n$  kernel matrix  $K$ 
  - with elements  $K_{ij} = k(x_i, x_j)$

When did we work in the dual?  
Why is it not bad to generate a  
potentially much larger feature  
space

# The “kernel trick” avoids computing $\phi(\mathbf{x})$

---

- $k(\mathbf{x}, \mathbf{y}) = \phi(\mathbf{x})^T \phi(\mathbf{y})$
- So we can compute  $k(\mathbf{x}, \mathbf{y})$  and never compute the expanded features  $\phi(\mathbf{x})$

# How are kernels selected?

---

- **Linear kernel**
  - $k(x,y) = x^T y$
- **Gaussian kernel**
  - $k(x,y) = \exp(-\|x - y\|^2/\sigma^2)$
- **Quadratic kernel**
  - $k(x,y) = (x^T y)^2$  or  $(x^T y + 1)^2$
- **Combinations and transformations of kernels**

## True or false

---

- Kernels in effect transform observations  $\mathbf{x}$  to a higher dimension space  $\phi(\mathbf{x})$
- Since kernels measure similarity,  
 $\mathbf{k}(\mathbf{x},\mathbf{y}) > \mathbf{k}(\mathbf{x},\mathbf{x})$  for  $\mathbf{x} \neq \mathbf{y}$
- If there exists a pair of points  $\mathbf{x}$  and  $\mathbf{y}$  such that  $\mathbf{k}(\mathbf{x},\mathbf{y}) < \mathbf{0}$ , then  $\mathbf{k}(\cdot)$  is not a kernel
- A quadratic kernel, when used in linear regression, gives results very similar to including quadratic interaction terms in the regression
- Any function  $\phi(\mathbf{x})$  can be used to generate a kernel using  $\mathbf{k}(\mathbf{x},\mathbf{y}) = \phi(\mathbf{x})^T \phi(\mathbf{y})$
- Any distance metric  $\mathbf{d}(\mathbf{x},\mathbf{y})$  can be used to generate a kernel using  $\mathbf{k}(\mathbf{x},\mathbf{y}) = \exp(-\mathbf{d}(\mathbf{x},\mathbf{y}))$