#### Introduction to Kernel Methods

David Rosenberg

New York University

February 21, 2017

# Setup and Motivation

David Rosenberg (New York University) DS-GA 1003 February 21, 2017 2 / 22

# The Input Space $\mathfrak X$

- ullet Our general learning theory setup: no assumptions about  ${\mathcal X}$
- But  $\mathfrak{X} = \mathbf{R}^d$  for the specific methods we've developed:
  - Ridge regression
  - Lasso regression
  - Support Vector Machines
  - Perceptrons
- Our hypothesis space for these was all affine functions on  $R^d$ :

$$\mathcal{H} = \left\{ x \mapsto w^T x + b \mid w \in \mathbf{R}^d, b \in \mathbf{R} \right\}.$$

• What if we want to do prediction on inputs not natively in  $\mathbb{R}^d$ ?

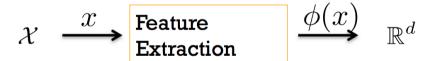
#### Feature Extraction

#### Definition

Mapping an input from X to a vector in  $\mathbb{R}^d$  is called **feature extraction** or **featurization**.

## Raw Input

# Feature Vector



• e.g. Quadratic feature map:  $\mathfrak{X} = \mathbf{R}^d$ 

$$\phi(x) = (x_1, \dots, x_d, x_1^2, \dots, x_d^2, \sqrt{2}x_1x_2, \dots, \sqrt{2}x_ix_j, \dots \sqrt{2}x_{d-1}x_d)^T.$$

# Linear Models with Explicit Feature Map

- Rather than take  $\mathfrak{X} = \mathbf{R}^d$ , let  $\mathfrak{X}$  be its own thing:
- ullet Input space:  ${\mathfrak X}$
- Introduce feature map  $\psi: \mathfrak{X} \to \mathbf{R}^d$
- The feature map maps into the feature space  $R^d$ .
- Hypothesis space of affine functions on feature space:

$$\mathcal{H} = \left\{ x \mapsto w^T \psi(x) + b \mid w \in \mathbf{R}^d, b \in \mathbf{R} \right\}.$$

## Linear Models Need Big Feature Spaces

- To get expressive hypothesis spaces using linear models,
  - need high-dimensional feature spaces
  - (What do we mean by expressive?)
- Very large feature spaces have two problems:
  - Overfitting
  - Memory and computational costs
- Overfitting we handle with regularization.
- Kernel methods can (sometimes) help with memory and computational costs.

### Some Methods Can Be "Kernelized"

#### **Definition**

A method is **kernelized** if inputs only appear inside inner products:  $\langle \psi(x), \psi(y) \rangle$  for  $x, y \in \mathcal{X}$ .

• The **kernel function** corresponding to  $\psi$  and inner product  $\langle \cdot, \cdot \rangle$  is

$$k(x,y) = \langle \psi(x), \psi(y) \rangle$$
.

- Why introduce this new notation k(x, y)?
- Turns out, we can often evaluate k(x, y) directly,
  - without explicitly computing  $\psi(x)$  and  $\psi(y)$ .
- For large feature spaces, can be much faster.

### Kernel Evaluation Can Be Fast

#### Example

Quadratic feature map

$$\phi(x) = (x_1, \dots, x_d, x_1^2, \dots, x_d^2, \sqrt{2}x_1x_2, \dots, \sqrt{2}x_ix_j, \dots, \sqrt{2}x_{d-1}x_d)^T$$

has dimension  $O(d^2)$ , but

$$k(w,x) = \langle \phi(w), \phi(x) \rangle = \langle w, x \rangle + \langle w, x \rangle^2$$

- Naively explicit computation of k(w,x):  $O(d^2)$
- Implicit computation of k(w,x): O(d)

## Kernels as Similarity Scores

- Can think of the kernel function as a **similarity score**.
- But this is not precise.
- There are many ways to design a similarity score.
  - A kernel function is special because it's an inner product.
  - Has many mathematical benefits.

### What's the Benefit of Kernelization?

- Computational (e.g. when feature space dimension d larger than sample size n).
- Access to infinite-dimensional feature spaces.
- 3 Allows thinking in terms of "similarity" rather than features.

10 / 22

Example: SVM

David Rosenberg (New York University) DS-GA 1003 February 21, 2017 11 / 22

### SVM Dual

Recall the SVM dual optimization problem

$$\sup_{\alpha} \sum_{i=1}^{n} \alpha_{i} - \frac{1}{2} \sum_{i,j=1}^{n} \alpha_{i} \alpha_{j} y_{i} y_{j} x_{j}^{T} x_{i}$$
s.t. 
$$\sum_{i=1}^{n} \alpha_{i} y_{i} = 0$$

$$\alpha_{i} \in \left[0, \frac{c}{n}\right] \ i = 1, \dots, n.$$

- Notice: x's only show up as inner products with other x's.
- Can replace  $x_i^T x_i$  by an arbitrary kernel  $k(x_i, x_i)$ .
- What kernel are we currently using?

# The Kernel Matrix (or the Gram Matrix)

#### Definition

For a set of  $\{x_1, \ldots, x_n\}$  and an inner product  $\langle \cdot, \cdot \rangle$  on the set, the **kernel matrix** or the **Gram matrix** is defined as

$$K = (\langle x_i, x_j \rangle)_{i,j} = \begin{pmatrix} \langle x_1, x_1 \rangle & \cdots & \langle x_1, x_n \rangle \\ \vdots & \ddots & \cdots \\ \langle x_n, x_1 \rangle & \cdots & \langle x_n, x_n \rangle \end{pmatrix}.$$

Then for the standard Euclidean inner product  $\langle x_i, x_i \rangle = x_i^T x_i$ , we have

$$K = XX^T$$

### SVM Dual with Kernel Matrix

$$\sup_{\alpha} \sum_{i=1}^{n} \alpha_{i} - \frac{1}{2} \sum_{i,j=1}^{n} \alpha_{i} \alpha_{j} y_{i} y_{j} K_{ji}$$
s.t. 
$$\sum_{i=1}^{n} \alpha_{i} y_{i} = 0$$

$$\alpha_{i} \in \left[0, \frac{c}{n}\right] i = 1, \dots, n.$$

- Once our algorithm works with kernel matrices, we can change kernel just by changing the matrix.
- Size of matrix:  $n \times n$ , where n is the number of data points.
- Recall with ridge regression, we worked with  $X^TX$ , which is  $d \times d$ , where d is feature space dimension.

## Some Kernels

David Rosenberg (New York University) DS-GA 1003 February 21, 2017 15 / 22

### Linear Kernel

- Input space:  $\mathfrak{X} = \mathbf{R}^d$
- Feature space:  $\mathcal{H} = \mathbf{R}^d$ , with standard inner product
- Feature map

$$\psi(x) = x$$
.

Kernel:

$$k(w,x) = w^T x$$

# Quadratic Kernel in R<sup>2</sup>

- Input space:  $\mathfrak{X} = \mathbb{R}^2$
- Feature space:  $\mathcal{H} = \mathbb{R}^5$
- Feature map:

$$\psi: (x_1, x_2) \mapsto (x_1, x_2, x_1^2, x_2^2, \sqrt{2}x_1x_2)$$

- Gives us ability to represent conic section boundaries.
- Define kernel as inner product in feature space:

$$k(w,x) = \langle \psi(w), \psi(x) \rangle$$

$$= w_1x_1 + w_2x_2 + w_1^2x_1^2 + w_2^2x_2^2 + 2w_1w_2x_1x_2$$

$$= w_1x_1 + w_2x_2 + (w_1x_1)^2 + (w_2x_2)^2 + 2(w_1x_1)(w_2x_2)$$

$$= \langle w, x \rangle + \langle w, x \rangle^2$$

Based on Guillaume Obozinski's Statistical Machine Learning course at Louvain, Feb 2014.

## Quadratic Kernel in R<sup>d</sup>

- Input space  $\mathfrak{X} = \mathbf{R}^d$
- Feature space:  $\mathcal{H} = \mathbf{R}^D$ , where  $D = d + \binom{d}{2} \approx d^2/2$ .
- Feature map:

$$\phi(x) = (x_1, \dots, x_d, x_1^2, \dots, x_d^2, \sqrt{2}x_1x_2, \dots, \sqrt{2}x_ix_j, \dots, \sqrt{2}x_{d-1}x_d)^T$$

Still have

$$k(w,x) = \langle \phi(w), \phi(x) \rangle$$
  
=  $\langle x, y \rangle + \langle x, y \rangle^2$ 

- Computation for inner product with explicit mapping:  $O(d^2)$
- Computation for implicit kernel calculation: O(d).

Based on Guillaume Obozinski's Statistical Machine Learning course at Louvain, Feb 2014.

# Polynomial Kernel in $\mathbf{R}^d$

- Input space  $\mathfrak{X} = \mathbf{R}^d$
- Kernel function:

$$k(w,x) = (1 + \langle w, x \rangle)^M$$

- ullet Corresponds to a feature map with all terms up to degree M.
- For any M, computing the kernel has same computational cost
- ullet Cost of explicit inner product computation grows rapidly in M.

# Radial Basis Function (RBF) / Gaussian Kernel

• Input space  $\mathfrak{X} = \mathbf{R}^d$ 

$$k(w,x) = \exp\left(-\frac{\|w-x\|^2}{2\sigma^2}\right),\,$$

where  $\sigma^2$  is known as the bandwidth parameter.

- Does it act like a similarity score?
- Why "radial"?
- Have we departed from our "inner product of feature vector" recipe?
  - Yes and no: corresponds to an infinite dimensional feature vector
- Probably the most common nonlinear kernel.

Kernel Trick: Overview

### Recap

- Given a kernelized ML algorithm.
- ② Can swap out the inner product for a new kernel function.
- New kernel may correspond to a high dimensional feature space.
- Once kernel matrix is computed, computational cost depends on number of data points, rather than the dimension of feature space.

Swapping out a linear kernel for a new kernel is called the kernel trick.