# Kernel Methods
Greg Mori - CMPT 419/726

Bishop PRML Ch. 6

## Non-linear Mappings

- In the lectures on linear models for regression and classification, we looked at models with $w^T\phi(x)$
- The feature space $\phi(x)$ could be high-dimensional
- This was good because if data aren't separable in original input space ($x$), they may be in feature space $\phi(x)$
- We'd like to avoid computing high-dimensional $\phi(x)$
- We'd like to work with $x$ which doesn't have a natural vector-space representation
  - e.g. graphs, sets, strings

## Kernel Trick

- In previous lectures on linear models, we would explicitly compute $\phi(x_i)$ for each datapoint
  - Run algorithm in feature space
- For some feature spaces, can compute dot product $\phi(x_i)^T\phi(x_j)$ efficiently
- Efficient method is computation of a kernel function $k(x_i, x_j) = \phi(x_i)^T\phi(x_j)$
- The kernel trick is to rewrite an algorithm to only have $x$ enter in the form of dot products
- The menu:
  - Kernel trick examples
  - Kernel functions

## A Kernel Trick

- Let's look at the nearest-neighbour classification algorithm
- For input point $x_i$, find point $x_j$ with smallest distance:

$$
\begin{aligned}
||x_i - x_j||^2 &= (x_i - x_j)^T(x_i - x_j) \\
&= x_i^T x_i - 2x_i^T x_j + x_j^T x_j
\end{aligned}
$$

- If we used a non-linear feature space $\phi(\cdot)$:

$$
\begin{aligned}
||\phi(x_i) - \phi(x_j)||^2 &= \phi(x_i)^T\phi(x_i) - 2\phi(x_i)^T\phi(x_j) + \phi(x_j)^T\phi(x_j) \\
&= k(x_i, x_i) - 2k(x_i, x_j) + k(x_j, x_j)
\end{aligned}
$$

- So nearest-neighbour can be done in a high-dimensional feature space without actually moving to it

## A Kernel Function

- Consider the kernel function $k(\boldsymbol{x}, \boldsymbol{z}) = (1 + \boldsymbol{x}^T \boldsymbol{z})^2$
- With $\boldsymbol{x}, \boldsymbol{z} \in \mathbb{R}^2$,

$$
\begin{aligned}
k(\boldsymbol{x}, \boldsymbol{z}) &= (1 + x_1 z_1 + x_2 z_2)^2 \\
&= 1 + 2x_1 z_1 + 2x_2 z_2 + x_1^2 z_1^2 + 2x_1 z_1 x_2 z_2 + x_2^2 z_2^2 \\
&= (1, \sqrt{2}x_1, \sqrt{2}x_2, x_1^2, \sqrt{2}x_1 x_2, x_2^2)(1, \sqrt{2}z_1, \sqrt{2}z_2, z_1^2, \sqrt{2}z_1 z_2, z_2^2)^T \\
&= \phi(\boldsymbol{x})^T \phi(\boldsymbol{z})
\end{aligned}
$$

- So this particular kernel function does correspond to a dot product in a feature space (is valid)
- Computing $k(\boldsymbol{x}, \boldsymbol{z})$ is faster than explicitly computing $\phi(\boldsymbol{x})^T \phi(\boldsymbol{z})$
  - In higher dimensions, larger exponent, much faster

## Why Kernels?

- Why bother with kernels?
  - Often easier to specify how similar two things are (dot product) than to construct explicit feature space $\phi$.
  - There are high-dimensional (even infinite) spaces that have efficient-to-compute kernels
    - Separability
- So you want to use kernels
  - Need to know when kernel function is valid, so we can apply the kernel trick

## Valid Kernels

- Given some arbitrary function $k(\boldsymbol{x}_i, \boldsymbol{x}_j)$, how do we know if it corresponds to a dot product in some space?
- Valid kernels: if $k(\cdot, \cdot)$ satisfies:
  - Symmetric; $k(\boldsymbol{x}_i, \boldsymbol{x}_j) = k(\boldsymbol{x}_j, \boldsymbol{x}_i)$
  - Positive definite; for any $\boldsymbol{x}_1, \ldots, \boldsymbol{x}_N$, the Gram matrix $\boldsymbol{K}$ must be positive semi-definite:

$$
\boldsymbol{K} = \begin{pmatrix} k(\boldsymbol{x}_1, \boldsymbol{x}_1) & k(\boldsymbol{x}_1, \boldsymbol{x}_2) & \ldots & k(\boldsymbol{x}_1, \boldsymbol{x}_N) \\ \vdots & \vdots & \ddots & \vdots \\ k(\boldsymbol{x}_N, \boldsymbol{x}_1) & k(\boldsymbol{x}_N, \boldsymbol{x}_2) & \ldots & k(\boldsymbol{x}_N, \boldsymbol{x}_N) \end{pmatrix}
$$

  - Positive semi-definite means $\boldsymbol{x}^T \boldsymbol{K} \boldsymbol{x} \geq 0$ for all $\boldsymbol{x}$

  then $k(\cdot, \cdot)$ corresponds to a dot product in some space $\phi$

  - a.k.a. Mercer kernel, admissible kernel, reproducing kernel

## Examples of Kernels

- Some kernels:
  - Linear kernel $k(\boldsymbol{x}_1, \boldsymbol{x}_2) = \boldsymbol{x}_1^T \boldsymbol{x}_2$
    - $\phi(\boldsymbol{x}) = \boldsymbol{x}$
  - Polynomial kernel $k(\boldsymbol{x}_1, \boldsymbol{x}_2) = (1 + \boldsymbol{x}_1^T \boldsymbol{x}_2)^d$
    - Contains all polynomial terms up to degree $d$
  - Gaussian kernel $k(\boldsymbol{x}_1, \boldsymbol{x}_2) = \exp(-||\boldsymbol{x}_1 - \boldsymbol{x}_2||^2 / 2\sigma^2)$
    - Infinite dimension feature space

## Constructing Kernels

- Can build new valid kernels from existing valid ones:
  - $k(x_1, x_2) = c k_1(x_1, x_2), c > 0$
  - $k(x_1, x_2) = k_1(x_1, x_2) + k_2(x_1, x_2)$
  - $k(x_1, x_2) = k_1(x_1, x_2) k_2(x_1, x_2)$
  - $k(x_1, x_2) = \exp(k_1(x_1, x_2))$
- Table on p. 296 gives many such rules

## More Kernels

- Stationary kernels are only a function of the difference between arguments: $k(x_1, x_2) = k(x_1 - x_2)$
  - Translation invariant in input space:
    $k(x_1, x_2) = k(x_1 + c, x_2 + c)$
- Homogeneous kernels, a. k. a. radial basis functions only a function of magnitude of difference: $k(x_1, x_2) = k(||x_1 - x_2||)$
- Set subsets $k(A_1, A_2) = 2^{|A_1 \cap A_2|}$, where $|A|$ denotes number of elements in $A$
- Domain-specific: think hard about your problem, figure out what it means to be similar, define as $k(\cdot, \cdot)$, prove positive definite (Feynman algorithm)

## Perceptron Classifier - Kernelized

- Recall the perceptron $y(x) = f(w^T \phi(x))$
- The update rule for the perceptron is

$$w^{(\tau+1)} = w^{(\tau)} + \underbrace{\eta \phi(x_n) t_n}_{if \ incorrect}$$

- Hence,

$$w^{(\tau+1)} = w^{(0)} + \alpha_1 \phi(x_1) + \alpha_2 \phi(x_2) + \ldots \alpha_N \phi(x_N)$$

- The classifier is then

$$f(w^T \phi(x)) = f(w^{(0),T}\phi(x) + \alpha_1 \phi(x_1)^T \phi(x) + \alpha_2 \phi(x_2)^T \phi(x) + \ldots)$$

- Kernelized! (init $w^{(0)} = 0$)
- Similar trick can be done for the update rule

## Regression - Kernelized

- Regularized least squares regression can also be kernelized
- Kernelized solution is

$$y(x) = k(x)^T (K + \lambda I_N)^{-1} t \quad \text{vs.} \quad \phi(x)(\Phi^T \Phi + \lambda I_M)^{-1} \Phi^T t$$

for original version

- $N$ is number of datapoints (size of Gram matrix $K$)
- $M$ is number of basis functions (size of matrix $\Phi^T \Phi$)
- Bad if $N > M$, but good otherwise

## Conclusion

- Readings: Ch. 6.1-6.2 (pp. 291-297)
- Many algorithms can be re-written with only dot products of features
  - We've seen NN, perceptron, regression; also PCA, SVMs (later)
- Non-linear features, or domain-specific similarity measurements are useful
- Dot products of non-linear features, or similarity measurements, can be written as kernel functions
  - Validity by positive semi-definiteness of kernel function
- Can have algorithm work in non-linear feature space without actually mapping inputs to feature space
  - Advantageous when feature space is high-dimensional