

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

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 $\phi(\mathbf{x}_i)$ for each datapoint
 - Run algorithm in feature space
- For some feature spaces, can compute dot product $\phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)$ efficiently
- Efficient method is computation of a kernel function $k(\mathbf{x}_i, \mathbf{x}_j) = \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)$
- 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

Kernel Trick

- In previous lectures on linear models, we would explicitly compute $\phi(\mathbf{x}_i)$ for each datapoint
 - Run algorithm in feature space
- For some feature spaces, can compute dot product $\phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)$ efficiently
- Efficient method is computation of a kernel function
$$k(\mathbf{x}_i, \mathbf{x}_j) = \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)$$
- 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

Kernel Trick

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

Kernel Trick

- In previous lectures on linear models, we would explicitly compute $\phi(\mathbf{x}_i)$ for each datapoint
 - Run algorithm in feature space
- For some feature spaces, can compute dot product $\phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)$ efficiently
- Efficient method is computation of a kernel function
$$k(\mathbf{x}_i, \mathbf{x}_j) = \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)$$
- The **kernel trick** is to rewrite an algorithm to only have \mathbf{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 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 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 $k(\mathbf{x}, \mathbf{z}) = (1 + \mathbf{x}^T \mathbf{z})^2$
- With $\mathbf{x}, \mathbf{z} \in \mathbb{R}^2$,



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

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

A Kernel Function



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

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

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

A Kernel Function



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

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

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

A Kernel Function



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

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

- So this particular kernel function does correspond to a dot product in a feature space (is valid)
- Computing $k(\mathbf{x}, \mathbf{z})$ is faster than explicitly computing $\phi(\mathbf{x})^T \phi(\mathbf{z})$
 - 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 ϕ .
 - 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 $k(\mathbf{x}_i, \mathbf{x}_j)$, how do we know if it corresponds to a dot product in some space?
- Valid kernels: if $k(\cdot, \cdot)$ satisfies:
 - Symmetric; $k(\mathbf{x}_i, \mathbf{x}_j) = k(\mathbf{x}_j, \mathbf{x}_i)$
 - Positive definite; for any $\mathbf{x}_1, \dots, \mathbf{x}_N$, the Gram matrix K must be positive semi-definite:

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

- Positive semi-definite means $\mathbf{x}^T K \mathbf{x} \geq 0$ for all \mathbf{x}
- then $k(\cdot, \cdot)$ corresponds to a dot product in some space ϕ
- a.k.a. Mercer kernel, admissible kernel, reproducing kernel

Valid Kernels

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

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

- Positive semi-definite means $\mathbf{x}^T \mathbf{K} \mathbf{x} \geq 0$ for all \mathbf{x}
- then $k(\cdot, \cdot)$ corresponds to a dot product in some space ϕ
- a.k.a. Mercer kernel, admissible kernel, reproducing kernel

Examples of Kernels

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

Constructing Kernels

- Can build new valid kernels from existing valid ones:
 - $k(\mathbf{x}_1, \mathbf{x}_2) = ck_1(\mathbf{x}_1, \mathbf{x}_2)$, $c > 0$
 - $k(\mathbf{x}_1, \mathbf{x}_2) = k_1(\mathbf{x}_1, \mathbf{x}_2) + k_2(\mathbf{x}_1, \mathbf{x}_2)$
 - $k(\mathbf{x}_1, \mathbf{x}_2) = k_1(\mathbf{x}_1, \mathbf{x}_2)k_2(\mathbf{x}_1, \mathbf{x}_2)$
 - $k(\mathbf{x}_1, \mathbf{x}_2) = \exp(k_1(\mathbf{x}_1, \mathbf{x}_2))$
- Table on p. 296 gives many such rules

More Kernels

- **Stationary kernels** are only a function of the difference between arguments: $k(\mathbf{x}_1, \mathbf{x}_2) = k(\mathbf{x}_1 - \mathbf{x}_2)$
 - Translation invariant in input space:
 $k(\mathbf{x}_1, \mathbf{x}_2) = k(\mathbf{x}_1 + \mathbf{c}, \mathbf{x}_2 + \mathbf{c})$
- **Homogeneous kernels**, a. k. a. **radial basis functions** only a function of magnitude of difference: $k(\mathbf{x}_1, \mathbf{x}_2) = k(\|\mathbf{x}_1 - \mathbf{x}_2\|)$
- 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 $y(\mathbf{x}) = f(\mathbf{w}^T \phi(\mathbf{x}))$
- The update rule for the perceptron is

$$\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(\tau)} + \underbrace{\eta \phi(\mathbf{x}_n) t_n}_{\text{if incorrect}}$$

- Hence,

$$\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(0)} + \alpha_1 \phi(\mathbf{x}_1) + \alpha_2 \phi(\mathbf{x}_2) + \dots + \alpha_N \phi(\mathbf{x}_N)$$

- The classifier is then

$$f(\mathbf{w}^T \phi(\mathbf{x})) = f(\mathbf{w}^{(0),T} \phi(\mathbf{x}) + \alpha_1 \phi(\mathbf{x}_1)^T \phi(\mathbf{x}) + \alpha_2 \phi(\mathbf{x}_2)^T \phi(\mathbf{x}) + \dots)$$

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

Perceptron Classifier - Kernelized

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

$$\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(\tau)} + \underbrace{\eta \phi(\mathbf{x}_n) t_n}_{\text{if incorrect}}$$

- Hence,

$$\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(0)} + \alpha_1 \phi(\mathbf{x}_1) + \alpha_2 \phi(\mathbf{x}_2) + \dots + \alpha_N \phi(\mathbf{x}_N)$$

- The classifier is then

$$f(\mathbf{w}^T \phi(\mathbf{x})) = f(\mathbf{w}^{(0),T} \phi(\mathbf{x}) + \alpha_1 \phi(\mathbf{x}_1)^T \phi(\mathbf{x}) + \alpha_2 \phi(\mathbf{x}_2)^T \phi(\mathbf{x}) + \dots)$$

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

Perceptron Classifier - Kernelized

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

$$\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(\tau)} + \underbrace{\eta \phi(\mathbf{x}_n) t_n}_{\text{if incorrect}}$$

- Hence,

$$\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(0)} + \alpha_1 \phi(\mathbf{x}_1) + \alpha_2 \phi(\mathbf{x}_2) + \dots + \alpha_N \phi(\mathbf{x}_N)$$

- The classifier is then

$$f(\mathbf{w}^T \phi(\mathbf{x})) = f(\mathbf{w}^{(0),T} \phi(\mathbf{x}) + \alpha_1 \phi(\mathbf{x}_1)^T \phi(\mathbf{x}) + \alpha_2 \phi(\mathbf{x}_2)^T \phi(\mathbf{x}) + \dots)$$

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

Regression - Kernelized

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

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

for original version

- N is number of datapoints (size of Gram matrix \mathbf{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