# Kernel Trick

He He

Slides based on Lecture 4d from David Rosenberg's course material.

CDS, NYU

March 2, 2021

# SVM with Explicit Feature Map

- Let $\psi : \mathcal{X} \to \mathsf{R}^d$ be a feature map.

- The SVM objective (with explicit feature map):

$$\min_{w \in \mathsf{R}^d} \frac{1}{2}\|w\|^2 + \frac{c}{n} \sum_{i=1}^{n} \max\left(0, 1 - y_i w^T \psi(x_i)\right).$$

- Computation is costly if $d$ is large (e.g. with high-degree monomials)

- Last time we mentioned an equivalent optimization problem from Lagrangian duality.

# SVM Dual Problem

- By Lagrangian duality, it is equivalent to solve the following dual problem:

$$\text{maximize} \quad \sum_{i=1}^{n} \alpha_i - \frac{1}{2} \sum_{i,j=1}^{n} \alpha_i \alpha_j y_i y_j \psi(x_j)^T \psi(x_i)$$

$$\text{s.t.} \quad \sum_{i=1}^{n} \alpha_i y_i = 0 \quad \text{and} \quad \alpha_i \in \left[0, \frac{c}{n}\right] \quad \forall i.$$

- If $\alpha^*$ is an optimal value, then

$$w^* = \sum_{i=1}^{n} \alpha_i^* y_i \psi(x_i) \quad \text{and} \quad \hat{f}(x) = \sum_{i=1}^{n} \alpha_i^* y_i \psi(x_i)^T \psi(x).$$

- Key observation: $\psi(x)$ only shows up in inner products with another $\psi(x')$ *for both training and inference*.

# Compute the Inner Products

Consider 2D data. Let's introduce degree-2 monomials using $\psi : R^2 \to R^3$.

$$(x_1, x_2) \mapsto (x_1^2, \sqrt{2} x_1 x_2, x_2^2).$$

The inner product is

$$\begin{aligned}
\psi(x)^T \psi(x') &= x_1^2 {x_1'}^2 + (\sqrt{2} x_1 x_2)(\sqrt{2} x_1' x_2') + x_2^2 {x_2'}^2 \\
&= (x_1 x_1')^2 + 2(x_1 x_1')(x_2 x_2') + (x_2 x_2')^2 \\
&= (x_1 x_1' + x_2 x_2')^2 \\
&= (x^T x')^2
\end{aligned}$$

We can calculate the inner product $\psi(x)^T \psi(x')$ without accessing the features $\psi(x)$!

## Compute the Inner Products

Now, consider monomials up to degree-2:

$$(x_1, x_2) \mapsto (1, \sqrt{2}x_1, \sqrt{2}x_2, x_1^2, \sqrt{2}x_1x_2, x_2^2).$$

The inner product can be computed by

$$\psi(x)^T \psi(x') = (1 + x^T x')^2 \quad \text{(check)}.$$

More generally, for features maps producing monomials up to degree-$p$, we have

$$\psi(x)^T \psi(x') = (1 + x^T x')^p.$$

(Note that the coefficients of each monomial in $\psi$ may not be 1)

**Kernel trick**: we do not need explicit features to calculate inner products.

- Using explicit features: $O(d^p)$
- Using implicit computation: $O(d)$

# Kernel Function

# The Kernel Function

- **Input space**: $\mathcal{X}$
- **Feature space**: $\mathcal{H}$ (a Hilbert space, e.g. $\mathsf{R}^d$)
- **Feature map**: $\psi : \mathcal{X} \to \mathcal{H}$
- The **kernel function** corresponding to $\psi$ is

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

  where $\langle \cdot, \cdot \rangle$ is the inner product associated with $\mathcal{H}$.

Why introduce this new notation $k(x, x')$?

- We can often evaluate $k(x, x')$ without explicitly computing $\psi(x)$ and $\psi(x')$.

When can we use the kernel trick?

# Some Methods Can Be "Kernelized"

## Definition

A method is **kernelized** if every feature vector $\psi(x)$ only appears inside an inner product with another feature vector $\psi(x')$. This applies to both the optimization problem and the prediction function.

The SVM Dual is a kernelization of the original SVM formulation.

Optimization:

$$\text{maximize} \quad \sum_{i=1}^{n} \alpha_i - \frac{1}{2} \sum_{i,j=1}^{n} \alpha_i \alpha_j y_i y_j \psi(x_j)^T \psi(x_i)$$

$$\text{s.t.} \quad \sum_{i=1}^{n} \alpha_i y_i = 0 \quad \text{and} \quad \alpha_i \in \left[0, \frac{c}{n}\right] \quad \forall i.$$

Prediction:

$$\hat{f}(x) = \sum_{i=1}^{n} \alpha_i^* y_i \psi(x_i)^T \psi(x).$$

# The Kernel Matrix

### Definition

The **kernel matrix** for a kernel $k$ on $x_1, \ldots, x_n \in \mathcal{X}$ is

$$\langle \phi(x_i), \phi(x_j) \rangle$$

$$K = \big(k(x_i, x_j)\big)_{i,j} = \begin{pmatrix} k(x_1, x_1) & \cdots & k(x_1, x_n) \\ \vdots & \ddots & \ldots \\ k(x_n, x_1) & \cdots & k(x_n, x_n) \end{pmatrix} \in \mathrm{R}^{n \times n}.$$

- In ML this is also called a **Gram matrix**, but traditionally (in linear algebra), Gram matrices are defined without reference to a kernel or feature map.

$$V_{ij} = \langle v_i, v_j \rangle$$

# The Kernel Matrix

- The kernel matrix summarizes all the information we need about the training inputs $x_1, \ldots, x_n$ to solve a kernelized optimization problem.

- In the kernelized SVM, we can replace $\psi(x_i)^T \psi(x_j)$ with $K_{ij}$:

$K_{n \times n}$

$$\text{maximize}_\alpha \quad \sum_{i=1}^{n} \alpha_i - \frac{1}{2} \sum_{i,j=1}^{n} \alpha_i \alpha_j y_i y_j K_{ij}$$

$\langle \psi(x_i), \psi(x_j) \rangle$

$$\text{s.t.} \quad \sum_{i=1}^{n} \alpha_i y_i = 0 \quad \text{and} \quad \alpha_i \in \left[0, \frac{c}{n}\right] \; i = 1, \ldots, n.$$

# Kernel Methods

Given a kernelized ML algorithm (i.e. all $\psi(x)$'s show up as $\langle \psi(x), \psi(x') \rangle$),

- Can swap out the inner product for a new kernel function.

- New kernel may correspond to a very high-dimensional feature space.

- Once the kernel matrix is computed, the computational cost depends on number of data points $n$, rather than the dimension of feature space $d$.

- Useful when $d >> n$.

- Computing the kernel matrix may still depend on $d$ and the essence of the **trick** is getting around this $O(d)$ dependence.

# Example Kernels

# Kernels as Similarity Scores



$$= \langle \psi(x), \psi(x') \rangle$$

$$\cos \theta = \frac{\langle x, x' \rangle}{\|x\| \, \|x'\|}$$

- Often useful to think of the $k(x, x')$ as a **similarity score** for $x$ and $x'$.

- We can design similarity functions without thinking about the explicit feature map, e.g. "string kernels", "graph kerners".

- How do we know that our kernel functions actually correspond to inner products in some feature space?

- Explicitly construct $\psi(x) : \mathcal{X} \to \mathrm{R}^d$ (e.g. monomials) and define $k(x, x') = \psi(x)^T \psi(x')$.

- Directly define the kernel function $k(x, x')$ ("similarity score"), and verify it corresponds to $\langle \psi(x), \psi(x') \rangle$ for some $\psi$.

There are many theorems to help us with the second approach.

# Linear Algebra Review: Positive Semidefinite Matrices

### Definition

A real, symmetric matrix $M \in \mathbb{R}^{n \times n}$ is **positive semidefinite (psd)** if for any $x \in \mathbb{R}^n$,

$$x^T M x \geqslant 0.$$

### Theorem

*The following conditions are each necessary and sufficient for a symmetric matrix $M$ to be positive semidefinite:*

- *$M$ can be factorized as $M = R^T R$, for some matrix $R$.*

- *All eigenvalues of $M$ are greater than or equal to $0$.*

# Positive Definite Kernel

## Definition

A symmetric function $k : \mathcal{X} \times \mathcal{X} \to \mathbb{R}$ is a **positive definite (pd)** kernel on $\mathcal{X}$ if for any finite set $\{x_1, \ldots, x_n\} \in \mathcal{X}$ ($n \in \mathbb{N}$), the kernel matrix on this set

$$K = \left( k(x_i, x_j) \right)_{i,j} = \begin{pmatrix} k(x_1, x_1) & \cdots & k(x_1, x_n) \\ \vdots & \ddots & \cdots \\ k(x_n, x_1) & \cdots & k(x_n, x_n) \end{pmatrix}$$

is a positive semidefinite matrix.

- Symmetric: $k(x, x') = k(x', x)$

- The kernel matrix needs to be positive semidefinite for any finite set of points.

- Equivalent definition: $\sum_{i=1}^{n} \sum_{j=1}^{n} \alpha_i \alpha_j k(x_i, x_j) \geqslant 0$ given $\alpha_i \in \mathbb{R} \ \forall i$.

# Mercer's Theorem

## Theorem

*A symmetric function $k(x, x')$ can be expressed as an inner product*

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

*for some $\psi$ if and only if $k(x, x')$ is **positive definite**.*

- Proving a kernel function is positive definite is typically not easy.

- But we can construct new kernels from valid kernels.

$$K = Q \wedge Q^\top$$
$$= \left( \wedge^{1/2} Q^\top \right)^\top \left( \wedge^{1/2} Q^\top \right)$$
$$\phi(x_i) = \left( \wedge^{1/2} Q^\top \right)[:, i]$$

# Generating New Kernels from Old

- Suppose $k, k_1, k_2 : \mathcal{X} \times \mathcal{X} \to \mathbb{R}$ are pd kernels. Then so are the following:

$$
\begin{aligned}
k_{\text{new}}(x, x') &= \alpha k(x, x') \quad \text{for } \alpha \geqslant 0 \quad \text{(non-negative scaling)} \\
k_{\text{new}}(x, x') &= k_1(x, x') + k_2(x, x') \quad \text{(sum)} \\
k_{\text{new}}(x, x') &= k_1(x, x') k_2(x, x') \quad \text{(product)} \\
k_{\text{new}}(x, x') &= k(\psi(x), \psi(x')) \text{ for any function } \psi(\cdot) \quad \text{(recursion)} \\
k_{\text{new}}(x, x') &= f(x) f(x') \text{ for any function } f(\cdot) \quad (f \text{ as 1D feature map})
\end{aligned}
$$

$\updownarrow$
$\phi$

- Lots more theorems to help you construct new kernels from old.

# Linear Kernel

- Input space: $\mathcal{X} = \mathsf{R}^d$

- Feature space: $\mathcal{H} = \mathsf{R}^d$, with standard inner product

- Feature map
$$\psi(x) = x$$

- Kernel:
$$k(x, x') = x^T x'$$

# Quadratic Kernel in $\mathsf{R}^d$

- Input space $\mathcal{X} = \mathsf{R}^d$

- Feature space: $\mathcal{H} = \mathsf{R}^D$, where $D = d + \binom{d}{2} \approx d^2/2$.

- Feature map:

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

- Then for $\forall x, x' \in \mathsf{R}^d$

$$
\begin{aligned}
k(x, x') &= \langle \psi(x), \psi(x') \rangle \\
&= \langle x, x' \rangle + \langle x, x' \rangle^2
\end{aligned}
$$

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

- Computation for implicit kernel calculation: $O(d)$.

# Polynomial Kernel in $R^d$

- Input space $\mathcal{X} = R^d$

- Kernel function:

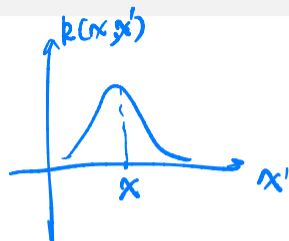$$k(x, x') = \left(1 + \langle x, x' \rangle\right)^M$$

- Corresponds to a feature map with all monomials up to degree $M$.

- For any $M$, computing the kernel has same computational cost

- Cost of explicit inner product computation grows rapidly in $M$.

# Radial Basis Function (RBF) / Gaussian Kernel

Input space $\mathcal{X} = \mathsf{R}^d$

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

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



- Probably the most common nonlinear kernel.

- 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

## Remaining Questions

Our current recipe:

- Recognize kernelized problem: $\psi(x)$ only occur in inner products $\psi(x)^T\psi(x')$

- Pick a kernel function ("similarity score")

- Compute the kernel matrix ($n$ by $n$ where $n$ is the dataset size)

- Optimize the model and make predictions by accessing the kernel matrix

Next: When can we apply kernelization?