# Kernel methods

Narayana Santhanam

EE 645

Jan 22, 2023

## Brief recap

Linear to non-linear

Support Vector Classification
Ridge Regression
    $\ell_2$ regularization makes it kernelizable

Gaussian process regression
    conditional means of Gaussians $=$ ridge regression
    though ridge computes mean, this is Bayesian
    predictions gaussian (with known variance)

# Complexity

If $x_1, \ldots, x_n$ are the training points, kernel $k(\cdot, \cdot)$, need the kernel Gram matrix:

$$\begin{bmatrix} k(x_1, x_1) & k(x_1, x_2) & \ldots & k(x_1, x_n) \\ \vdots & \vdots & \vdots & \vdots \\ k(x_n, x_1) & k(x_n, x_2) & \ldots & k(x_n, x_n) \end{bmatrix}$$

The above matrix has $n^2$ entries (and we often need to invert matrices of this size. Complexity is quadratic or worse.

In certain cases, we can get to linear complexity in $n$ (training size)

# Deeper look into kernels

In certain cases, we can get to linear complexity in $n$ (training size)
approximate solutions, not exact

# Deeper look into kernels

In certain cases, we can get to linear complexity in $n$ (training size)
    approximate solutions, not exact
    NeurIPS Test of Time award for influential papers

## Deeper look into kernels

In certain cases, we can get to linear complexity in $n$ (training size)
   approximate solutions, not exact
   NeurIPS Test of Time award for influential papers

Makes large training sets feasible
Kernel methods have many of the "amazing" features neural nets
have

# Deeper look into kernels

In certain cases, we can get to linear complexity in $n$ (training size)
  approximate solutions, not exact
  NeurIPS Test of Time award for influential papers

Makes large training sets feasible
Kernel methods have many of the "amazing" features neural nets have
  Can often fit any random permutations of labels

# Deeper look into kernels

In certain cases, we can get to linear complexity in $n$ (training size)
   approximate solutions, not exact
   NeurIPS Test of Time award for influential papers

Makes large training sets feasible
Kernel methods have many of the "amazing" features neural nets have
   Can often fit any random permutations of labels
   ... yet do not misuse power and overfit!

# Function positive definiteness

Function $k(x, y)$ is said to be positive (semi-)definite if

# Function positive definiteness

Function $k(x, y)$ is said to be positive (semi-)definite if for all $n$ and all $x_1, \ldots, x_n$,

# Function positive definiteness

Function $k(x, y)$ is said to be positive (semi-)definite if for all $n$ and all $x_1, \ldots, x_n$, the Gram matrix is positive semi-definite.

## Function positive definiteness

Function $k(x, y)$ is said to be positive (semi-)definite if for all $n$ and all $x_1, \ldots, x_n$, the Gram matrix is positive semi-definite.

This means that for all vectors $w = \begin{bmatrix} w_1 \\ \vdots \\ w_n \end{bmatrix}$

$$\begin{bmatrix} w_1 & \ldots & w_n \end{bmatrix} \begin{bmatrix} k(x_1, x_1) & k(x_1, x_2) & \ldots & k(x_1, x_n) \\ \vdots & \vdots & \vdots & \vdots \\ k(x_n, x_1) & k(x_n, x_2) & \ldots & k(x_n, x_n) \end{bmatrix} \begin{bmatrix} w_1 \\ \vdots \\ w_n \end{bmatrix} \geq 0$$

Not enough that all entries of Gram matrix $\geq 0$

Any positive (semi-)definite $k$ is allowed to be a kernel

UNIVERSITY
of HAWAI'I
MĀNOA

# Making new kernels from old

Linear combinations with non-negative coeffs
    if $k_1$ and $k_2$ are two kernels, so is $\alpha k_1(\mathsf{x}, \mathsf{y}) + \beta k_2(\mathsf{x}, \mathsf{y})$

## Making new kernels from old

Linear combinations with non-negative coeffs
  if $k_1$ and $k_2$ are two kernels, so is $\alpha k_1(\mathsf{x},\mathsf{y}) + \beta k_2(\mathsf{x},\mathsf{y})$

Product of kernels
  if $k_1$ and $k_2$ are two kernels, so is $k_1(\mathsf{x},\mathsf{y})k_2(\mathsf{x},\mathsf{y})$

UNIVERSITY
*of* HAWAI'I
MĀNOA

# Making new kernels from old

Linear combinations with non-negative coeffs
if $k_1$ and $k_2$ are two kernels, so is $\alpha k_1(x, y) + \beta k_2(x, y)$

Product of kernels
if $k_1$ and $k_2$ are two kernels, so is $k_1(x, y) k_2(x, y)$

If $g(x)$ is any function $k(x, y) = g(x) g(y)$ is a kernel

University
of Hawaiʻi
MĀNOA

# Making new kernels from old

Linear combinations with non-negative coeffs
   if $k_1$ and $k_2$ are two kernels, so is $\alpha k_1(x, y) + \beta k_2(x, y)$

Product of kernels
   if $k_1$ and $k_2$ are two kernels, so is $k_1(x, y)k_2(x, y)$

If $g(x)$ is any function $k(x, y) = g(x)g(y)$ is a kernel

If $k(x, y)$ is any kernel, so are $\exp(k(x, y))$ and $k(f(x), f(y))$

# Examples

Radial basis function (for scale parameter $s > 0$)

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

## Examples

Radial basis function (for scale parameter $s > 0$)

$$k(\mathrm{x}, \mathrm{y}) = \exp\left(-\frac{||\mathrm{x} - \mathrm{y}||^2}{2s}\right)$$

This function is positive semi-definite because

$$\exp\left(-\frac{||\mathrm{x} - \mathrm{y}||^2}{2s}\right) = \exp\left(-\frac{||\mathrm{x}||^2}{2s}\right)\exp\left(-\frac{||\mathrm{y}||^2}{2s}\right)\exp\left(\frac{\mathrm{x}^T\mathrm{y}}{s}\right)$$

Exponential/Laplace kernel

$$k(x, y) = \exp(-||x - y||/\lambda)$$

# Examples

Exponential/Laplace kernel

$$k(x, y) = \exp(-||x - y||/\lambda)$$

Positive semi-definiteness of this function not trivial
but follows easily from Bochner's theorem

## Examples

Exponential/Laplace kernel

$$k(x, y) = \exp(-||x - y||/\lambda)$$

Positive semi-definiteness of this function not trivial

but follows easily from Bochner's theorem

... as for the whole class of Matern kernels and a host of others

# Bochner's Theorem

Need this for two reasons

    finding kernels
    faster computation

# Bochner's Theorem

Consider kernels $k(x, y)$ where dependency only via $||x - y||$
  rbf, Matern
  not examples: polynomial

### Bochner

$k(x - y)$, $x, y \in \mathbb{R}^d$ is positive semi-definite iff it is the
($d$−dimensional) Fourier transform of a finite positive measure on
$\mathbb{R}^d$ (think pdf).

## Fourier transform

Let $\mu$ be absolutely continuous wrt to the Lebesgue measure (ignore if you haven't heard the terms). Let the pdf of $\mu$ be $f_\mu$. Then

$$F(x - y) = \int_{\nu \in \mathbb{R}^d} e^{-j2\pi\nu^T(x-y)} f_\mu(\nu) d\nu$$

is a valid kernel.

we interpret the kernel $k(x, y) = F(x - y)$.

we call $f_\mu$ the kernel spectral measure

## Bochner's kernels

Familiar examples:

if measure is normal, radial basis kernel

similarly for Matern kernels

## Bochner's kernels

Familiar examples:

     if measure is normal, radial basis kernel

     similarly for Matern kernels

Interestingly, these are also universal
    any compactly supported function arbitrarily approximated

University *of* Hawai'i
MĀNOA

## Computational speedups

Bochner's theorem can also speed up computations (stationary kernels)

From Bochner's theorem

$$k(\mathsf{x}, \mathsf{y}) = \mathbb{E} \exp\left(j 2\pi \nu^T (\mathsf{x} - \mathsf{y})\right)$$

$\nu$ random $d-$vector $\sim$ kernel spectral measure $f_\mu(\nu)$
$\mathbb{E}$ denotes expectation

# Random Fourier Features

$z_\nu(x) = \cos(2\pi\nu^T x + b)$, $\nu \sim f_\mu$ and $b$ uniform

## Random Fourier Features

$z_\nu(x) = \cos(2\pi \nu^T x + b)$, $\nu \sim f_\mu$ and $b$ uniform

Recall feature map $x \to \phi(x)$, $k(x, y) = \phi(x)^T \phi(y)$

    replace $\phi(x)$ with $z(x)$ with same property

    yet $z$ is a vector with $D$ coordinates ($D$ small)

    $z^T(x) = \begin{bmatrix} z_{\nu_1}(x) & , \ldots, z_{\nu_D}(x) \end{bmatrix}$

    $k(x, y) = \mathbb{E}_\nu z_\nu(x)^T z_\nu(y) \approx z(x)^T z(y)$