

## Kernels

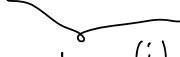
Supervised learning framework:

- Given training data  $D = \{x^{(i)}, y^{(i)}\}_{i=1}^n$ , we want to find  $f(x)$  s.t.  $y^{(i)} = f(x^{(i)})$ .
- Use Empirical Risk Minimization:
  - Decide a function class  $F$ .
  - Define an empirical loss function  $L(f, D)$  for  $f \in F$ .
  - Solve the optimization problem:  $\hat{f} = \underset{f \in F}{\operatorname{argmin}} L(f, D)$ .

For example, for linear regression:

The function class is:

$$F \triangleq \left\{ f_{\theta}(x) = \sum_{l=1}^d \theta_l x_l + \theta_0 \mid \theta_l \in \mathbb{R} \right\}, \quad x = \begin{bmatrix} x_1 \\ \vdots \\ x_d \end{bmatrix}$$

  
each  $x^{(i)}$   
has  $d$  features.

The objective function is:

$$\text{Min } L(f_{\theta}, D) = \frac{1}{n} \sum_{i=1}^n (y^{(i)} - f_{\theta}(x^{(i)}))^2$$

## Intro Kernels

A way to incorporate non-linearities into most linear classifiers, by applying basis function (feature transformations) on the input feature vectors.

$$x \rightarrow \underbrace{\phi(x)}_{\in \mathbb{R}^D}$$

,  $D \gg d$ , bc we add dimensions that capture non-linear interactions among original features.

Advantage: the problem stays simple, convex, well-behaved.

Disadvantage:  $\phi(x)$  might be very high dimensional.

Consider the following example:  $x = \begin{pmatrix} x^{(1)} \\ x^{(2)} \\ \vdots \\ x^{(d)} \end{pmatrix} \rightarrow \phi(x) = \begin{pmatrix} 1 & x^{(1)} \\ & \vdots \\ & x^{(d)} \\ & \vdots \\ & x^{(k-1)} & x^{(d)} \\ & \vdots \\ x^{(1)} & x^{(2)} & \dots & x^{(d)} \end{pmatrix}$

## Linear Regression Using Kernels:

Linear Regression:  $f(x, \theta) = \sum_{l=1}^d \theta_l x_l = \theta^\top X$ ,  $X = \begin{bmatrix} x^{(1)} \\ \vdots \\ x^{(d)} \end{bmatrix}$

Objective:  $\min_{\theta} L(f_\theta, D) = \underbrace{\frac{1}{d} \sum_{i=1}^d (y^{(i)} - f_\theta(x^{(i)}))^2}_{\text{Mean Square Loss}}$

However, This set-up cannot capture non-linear relations.

Solution:

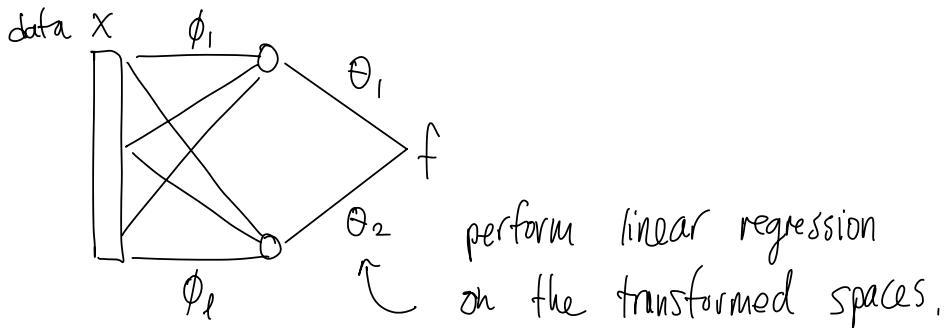
Transform linear function into non-linear function:

Linear:  $f(x, \theta) = \sum_{l=1}^d \theta_l x_l = \theta^\top X$   $x_l \rightarrow \phi_l$

Nonlinear:  $f(x, \theta) = \sum_{l=1}^m \theta_l \phi_l(x) = \theta^\top \phi(x)$

↑ coefficient vector      ↑ set of basis functions

$$\theta = [\theta_1, \dots, \theta_m]^\top \quad [\phi_1(x), \dots, \phi_m(x)]^\top$$



here, # of  $\phi$   $>>$  # of data points.

$\therefore$  You need to manually decide what basis function to use, which is not ideal. You want this to be automatic and adaptive.

$\Rightarrow$  use adaptive basis function, of which kernel method is one.

How to construct the basis functions?

Adaptive basis functions :

Kernels

Neural networks

## Kernel Method :

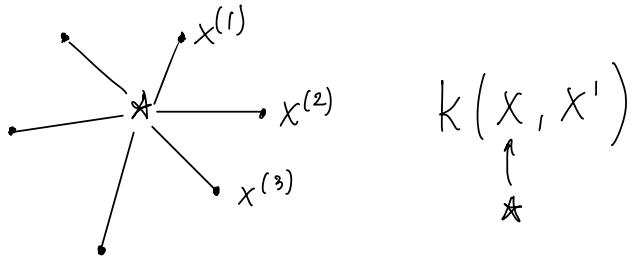
let  $k(x, x')$  be the similarity measure btwn  $x$  and  $x'$ .

$$k(x, x') : X \times X \mapsto \mathbb{R},$$

let  $k(x, x') = k(x', x)$ . // two variable symmetric function.

This similarity measure is called a kernel function.

Geometric interpretation: Given  $D = \{x^{(i)}, y^{(i)}\}_{i=1}^n$



## Example Kernel Functions

A typical kernel function is the Gaussian radial basis function (RBF) kernel:

$$k(x, x') = \exp\left[-\frac{1}{2h^2} \|x - x'\|_2^2\right],$$

where

$h$  = bandwidth  $\in \mathbb{R}^+$ . It is a scaling parameter, that should be in the range of the data points.

Ex: variance.

Another Example:  $k(x, x') = \exp\left[-\frac{1}{h} \|x - x'\|\right]$  Laplace kernel.

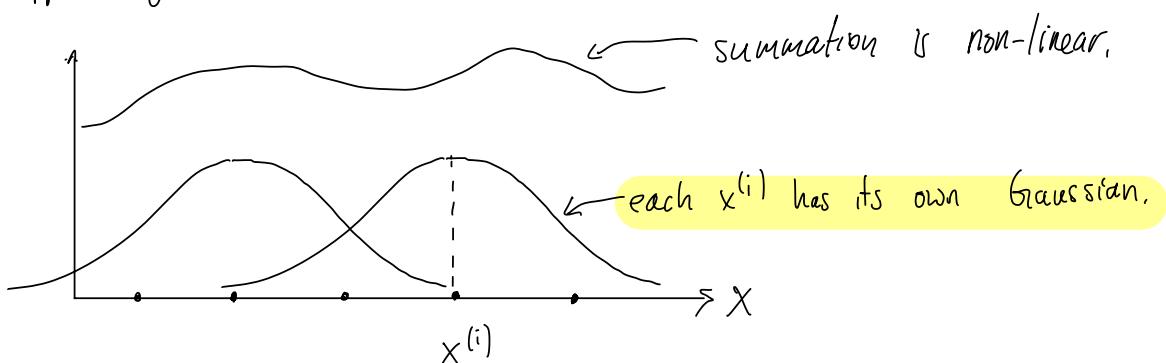
Given a dataset  $D = \{x^{(i)}\}_{i=1}^n$ , we may construct a kernel representation of a point  $x$  by comparing it with each observed data point in the dataset:

$$\Theta = \begin{bmatrix} \theta_1 \\ \vdots \\ \theta_n \end{bmatrix}, \quad \phi(x) = \begin{bmatrix} k(x, x^{(1)}) \\ \vdots \\ k(x, x^{(n)}) \end{bmatrix} \text{ where } \{x^{(i)}\}_s \text{ are the fixed observed data points, } x \text{ is the variable of the function: } \phi_i = k(x, x^{(i)}).$$

$$f_\Theta(x) = \Theta^\top \phi(x) = \sum_{i=1}^n \theta_i \phi_i(x) = \sum_{i=1}^n \theta_i k(x, x_i)$$

Example using RBF:

Suppose you have a set of 1D data:



We want:

$F$  to be adaptive with data  $D$

$$\dim(F) = n \rightarrow +\infty.$$

The class of function that has dimension that grows with the # of data points is called non-parametric methods. Kernel method is a special case of non-parametric method.

Loss function:

$$L(\theta) = \sum_{j=1}^n \left[ y^{(j)} - \sum_{i=1}^n \theta_i \underbrace{k(x^{(j)}, x^{(i)})}_{\phi_i(x)} \right]^2$$

$$= \|y - K\theta\|_2^2, \text{ where}$$

$$Y = \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(n)} \end{bmatrix}_{n \times 1} \quad K = \begin{bmatrix} k(x_1, x_1) & \dots & k(x_1, x_n) \\ \vdots & \ddots & \vdots \\ k(x_n, x_1) & \dots & k(x_n, x_n) \end{bmatrix}_{n \times n} \quad \text{gram matrix, } n \times n.$$

or  $K = [k(x^{(i)}, x^{(j)})]_{i,j=1}^n$

Objective =  $\min L(\theta)$

$$= \min_{\theta} \|y - K\theta\|_2^2$$

$$\nabla L(\theta) = 2K^T(K\theta - y) \quad // \text{shown without derivation}$$

$$\Rightarrow K^T K \theta = K^T y.$$

$$\Rightarrow \boxed{\theta = K^{-1} y} \quad \text{Need to assume } K \text{ is invertible.}$$

This regression can perfectly fit all the data points. However, it would be at risk for overfitting.

Regularization : use regularization to avoid overfitting.

Solve a regularized version:

$$\min_{\theta} \left\| y - K\theta \right\|_2^2 + \alpha \Phi(\theta)$$

↑  
hyperparameter      ↑ regularization. Often L2 Norm.  
that we choose.

Instead L2 regularization ( $\Phi(\theta) = \|\theta\|_2^2$ ), use k-norm:

$$\Phi(\theta) = \|\theta\|_k^2 \equiv \theta^T K \theta$$

↑ Gram / Kernel Matrix, size  $n \times n$ .

---

Why  $\|\theta\|_k^2$  vs  $\|\theta\|_2^2$ ? Consider an Example:

$$D = \{x^{(1)}, x^{(2)}, x^{(3)}\}$$

assume equal:  $x^{(1)} = x^{(2)}$ .

$$x^{(1)} = x^{(2)} + x^{(3)}$$

$$\bullet x^{(1)} = x^{(2)}$$

$$x^{(3)}$$

The Kernel representation is:

$$\begin{aligned} f(x) &= \theta_1 k(x, x^{(1)}) + \theta_2 k(x, x^{(2)}) + \theta_3 k(x, x^{(3)}) \\ &= (\theta_1 + \theta_2) k(x, x^{(1)}) + \theta_3 k(x, x^{(3)}) \end{aligned}$$

$\Rightarrow \|\theta\|_2^2 = \theta_1^2 + \theta_2^2 + \theta_3^2$ , which regularizes each term,  
does not make sense. More reasonable to use:

$$\underbrace{(\theta_1 + \theta_2)^2}_{\text{combined.}} + \theta_3^2 = \|\theta\|_k^2 \text{ as the penalty.}$$

Kernel norm does what we want (grouping effect):

Consider the Gaussian kernel:

$$k(x, x') = \exp\left(\frac{-1}{2h^2} \|x - x'\|^2\right)$$

which has  $k(x, x) = 1 \neq x$ .

If  $x^{(1)}$  and  $x^{(3)}$  are very different, then  $k(x^{(1)}, x^{(3)}) \approx 0$

Therefore the Gram matrix is:

$$K \approx \begin{matrix} & x^{(1)} & x^{(2)} & x^{(3)} \\ x^{(1)} & 1 & 1 & 0 \\ x^{(2)} & 1 & 1 & 0 \\ x^{(3)} & 0 & 0 & 1 \end{matrix},$$

which suggests  $\underline{\Phi}(\theta) =$

$$\theta^\top K \theta = [\theta_1 \ \theta_2 \ \theta_3] \begin{bmatrix} 1 & 1 & 0 \\ 1 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} \theta_1 \\ \theta_2 \\ \theta_3 \end{bmatrix} = [\theta_1 + \theta_2]^2 + \theta_3^2$$

Regularized Objective Function:

$$\min_{\theta} \|\gamma - K \theta\|_2^2 + \alpha \theta^\top K \theta.$$

$\Rightarrow$  Solve in closed form:

$$\hat{\theta} = (K + \lambda I)^{-1} \gamma$$

Identity

## Kernel Set Up :

Classification:  $\{x^{(i)}, y^{(i)}\}$ ,  $y^{(i)} \in \{-1, +1\}$  //binary classification

$$L(\theta) = \sum_{i=1}^n \sigma(y^{(i)} f_\theta(x^{(i)}))$$

↑  
surrogate function, can be logistic loss,  
hinge loss, etc.

$$= \sum_{i=1}^n \sigma\left(y^{(i)} \sum_{j=1}^n \theta_j K(x^{(j)}, x^{(i)})\right) + \lambda \theta^\top K \theta$$

↑  
 $\alpha?$

Solve numerically.

## Kernel as infinite dimensional features:

Using Kernel  $\Leftrightarrow$  using infinite basis functions.

Kernel classification  $\Leftrightarrow$  solving an infinite dimensional optimization problem.

Define:  $\{\psi_\ell(x) : \ell = 0, 1, 2, \dots, \infty\}$

Example function:  $\psi_\ell(x) = x^\ell$

Set up:  $f(x) = \sum_{\ell=0}^{\infty} w_\ell \psi_\ell(x)$

↑  $\ell^{\text{th}}$  basis function of  
infinite sequence of  
basis functions.

weight of the  $\ell^{\text{th}}$  basis function, or other

The idea is to have infinite number of (polynomial) functions, and a linear combination of them can be used to approximate any function. Sort of like Taylor Series, or Fourier Series.

Define Loss function:

$$L(w) = \underset{\substack{\text{infinite} \\ \text{sequence} \\ \text{of } w_l}}{\mathbb{E}_D} \left[ \left( y - \sum_{l=0}^{\infty} w_l \psi_l(x) \right)^2 \right] + \sum_{l=0}^{\infty} \frac{w_l^2}{\lambda_l}$$

$\uparrow$   
expectation over  
the dataset.

$\lambda_l$  = regularization coefficient for  $w_l$ .  
Controls how much to penalize  $w_l$ .

$$\text{Define } \lambda_l := \frac{1}{\alpha_l} > 0. ?$$

$$\lambda_1 \geq \lambda_2 \geq \dots \geq 0$$

Note: If  $\lambda_l$  is small, then penalty is huge:  $w_l \rightarrow 0$ .

$\Rightarrow$  Higher order terms are penalized exponentially;  
they contribute much more to the loss function,

Result:

Solving this infinite dimensional optimization problem is equivalent to solving a finite dimensional kernel representation. This is due to the regularization term turning off the higher order terms.

In practice, the optimization problem is not infinite, because we don't have infinite data.

How is kernel related to infinite dimensional features?

Define  $\hat{w} = \underset{w}{\operatorname{argmin}} L(w)$

$$\underset{\substack{\text{optimal} \\ \text{solution.}}}{=} \underset{w}{\operatorname{argmin}} \sum_{i=1}^n \left( y^{(i)} - \sum_{l=0}^{\infty} w_l \psi_l(x^{(i)}) \right)^2 + \sum_{l=1}^{\infty} \frac{w_l^2}{\lambda_l}$$

vector of weights      vector of infinite basis functions

$$= \underset{w}{\operatorname{argmin}} \hat{E}_D \left[ |y - \sum_{l=0}^{\infty} w_l \psi_l(x)|^2 \right] + \underbrace{\sum_{l=0}^{\infty} \frac{w_l^2}{\lambda_l}}_{\text{regularization term; prevents norm of } \{w_l\} \text{ from being too large.}}$$

define as  $\hat{f}(x)$

Define  $\hat{f}(x) = \sum_{l=0}^{\infty} \hat{w}_l \psi_l(x)$   $\begin{pmatrix} \text{infinite sequence} \\ \text{representation} \\ \text{of the data} \end{pmatrix}$   
 Corresponding optimal function for  $\hat{w}$

Then:  $\hat{f}(x) = \sum_{i=1}^n \theta_i \underbrace{k(x, x^{(i)})}_{\substack{\text{replaced } \psi \\ \text{replaced } w}}$  for some  $\theta_i$   $\begin{pmatrix} \text{finite kernel} \\ \text{representation} \\ \text{of the data} \end{pmatrix}$

and 
$$k(x, x') \equiv \sum_{l=0}^{\infty} \lambda_l \psi_l(x) \psi_l(x')$$
  $\begin{pmatrix} \text{definition of} \\ \text{kernel function} \end{pmatrix}$

and 
$$\sum_{l=0}^{\infty} \frac{w_l^2}{\lambda_l} = \theta^\top K \theta$$
 Note:  $\theta_i = y^{(i)} - \hat{f}(x^{(i)})$  ?

The infinite representation then reduces to:

$\Leftrightarrow \hat{\theta} = \underset{\theta}{\operatorname{argmin}} \hat{E}_D \left[ \left( y - \sum_{i=1}^n \theta_i k(x, x^{(i)}) \right)^2 \right] + \overbrace{\theta^\top K \theta}^{\substack{\text{no longer infinite} \\ \text{regularization term is = kernel loss.}}}$

$\Rightarrow \hat{\theta} = \underset{\theta}{\operatorname{argmin}} \hat{E}_D \left[ \left( y - \sum_{i=1}^n \theta_i \right)^2 \right]$

proof: Find the zero gradient condition, which implies this relation at optimal point.

Objective is to solve this optimization problem:

$$\min_{\hat{w}} \hat{E}_D \left[ \underbrace{\left( y - \sum_{l=0}^{\infty} w_l \psi_l(x) \right)^2}_{f(x)} \right] + \sum_{l=0}^{\infty} \frac{w_l^2}{\lambda_l} \quad \begin{array}{l} \text{(infinite sequence)} \\ \text{representation.} \end{array}$$

Define  $\delta(x) = y - \sum_{l=0}^{\infty} w_l \psi_l(x)$

take gradient  $\rightarrow \nabla_{w_l} L(\hat{w}) = 2 \hat{E}_D \left[ -\delta(x) \psi_l(x) \right] + \frac{2 \hat{w}_l}{\lambda_l} \stackrel{\text{at optimal point, gradient } = 0}{=} 0$

$$\Rightarrow \frac{2 w_l}{\lambda_l} = 2 \hat{E}_D \left[ \delta(x) \psi_l(x) \right] \Rightarrow \hat{w}_l = \lambda_l \hat{E}_D \left[ \delta(x) \psi_l(x) \right]$$

$$\Rightarrow \hat{w}_l = \frac{\lambda_l}{n} \sum_{i=1}^n \left[ \delta(x^{(i)}) \psi_l(x^{(i)}) \right]$$

Then substitute  $\hat{w}$  into  $f(x)$ :

summation over all data points,

$$\hat{f}(x) = \sum_{l=0}^{\infty} \hat{w}_l \psi_l(x) = \sum_{l=0}^{\infty} \left[ \underbrace{\frac{\lambda_l}{n} \sum_{i=1}^n \delta(x^{(i)}) \psi_l(x^{(i)})}_{\text{summation over all basis functions}} \right] \psi_l(x)$$

$$= \sum_{l=0}^{\infty} \sum_{i=1}^n \frac{\lambda_l}{n} \delta(x^{(i)}) \psi_l(x^{(i)}) \psi_l(x) \quad \text{// take summation out.}$$

$$\begin{aligned}
&= \sum_{i=1}^n \sum_{l=0}^{\infty} \frac{\lambda_l}{n} \delta(x^{(i)}) \psi_l(x^{(i)}) \psi_l(x) \quad // \text{switched order of summation.} \\
&= \sum_{i=1}^n \frac{\delta(x^{(i)})}{n} \underbrace{\sum_{l=0}^{\infty} \lambda_l \psi_l(x^{(i)}) \psi_l(x)}_{\text{kernel function } k(x^{(i)}, x)} \quad // \text{move out of summation, bc } \delta(x^{(i)}) \text{ is independent of } l. \\
&= \boxed{\sum_{i=1}^n \hat{\theta}_i k(x^{(i)}, x) \quad \text{where } \hat{\theta}_i \triangleq \frac{\delta(x^{(i)})}{n}}
\end{aligned}$$

= linear combo of kernel functions.

This proves the first part of the result.

The second part is in textbook.

Two views that are equivalent:

$$\begin{aligned}
1) \hat{f}(x) &= \sum_{i=1}^n \hat{\theta}_i k(x^{(i)}, x) \quad // \text{finite linear combo of kernel functions evaluated at data point.} \\
2) \hat{f}(x) &= \sum_{l=0}^{\infty} \hat{\omega}_l \psi_l(x) \quad // \text{infinite basis functions. Theoretical purpose.}
\end{aligned}$$

Connected by:  $k(x, x') = \underbrace{\sum_{l=0}^{\infty} \lambda_l \psi_l(x) \psi_l(x')}_{\text{any kernel function}} \quad \underbrace{\sum_{l=0}^{\infty} \lambda_l \psi_l(x) \psi_l(x')}_{\text{infinite basis functions}}$

$\lambda_l \geq 0$ .  
 $\lambda_l$  inverse regularization coefficient.

$\therefore$  Any kernel function can be represented using some basis functions.

Mercer Theorem: any continuous positive definite kernel can be represented by  $k(x, x') = \sum_{\ell=0}^{\infty} \lambda_{\ell} \psi_{\ell}(x) \psi_{\ell}(x')$  with  $\lambda_{\ell} \geq 0$ .

1) If  $k(x, x')$  is positive definite  $\Leftrightarrow k = [k(x^{(i)}, x^{(j)})]_{i,j=1}^n \geq 0$  for any  $n$  and any data point  $\{x^{(i)}\}$ .

2)  $\Leftrightarrow \exists \psi_{\ell}, \lambda_{\ell} \geq 0$ .

The typical kernels are positive definite, so don't worry!

Example: Gaussian RBF kernel:

$$\begin{aligned} k(x, x') &= \exp[-\gamma \|x - x'\|^2] \quad // \text{definition,} \\ &= \exp[-\gamma \|x\|^2 + 2\gamma x^T x' - \gamma \|x'\|^2] \quad // \text{expanded form} \\ &= \underbrace{\exp(-\gamma \|x\|^2)}_{\text{depends on } x} \underbrace{\exp(-\gamma \|x'\|^2)}_{\text{depends on } x'} \underbrace{\exp(2\gamma x^T x')}_{\text{depends on both}}, \end{aligned}$$

$$\text{Taylor Expansion: } \exp(t) = 1 + t + \frac{t^2}{2} + \dots + \frac{t^n}{n!} = \sum_{\ell=0}^{\infty} \frac{t^{\ell}}{\ell!}$$

$$\exp(2\gamma x^T x') = \sum_{\ell=0}^{\infty} \frac{(2\gamma x^T x')^{\ell}}{\ell!}, \text{ etc, for the other 2 terms.}$$

$$\Rightarrow k(x, x') = \sum_{\ell=0}^{\infty} \lambda_{\ell} \psi_{\ell}(x) \psi_{\ell}(x'), \quad \text{Assume } x \in \mathbb{R}.$$

$$\text{where } \psi_{\ell}(x) = \exp(-\gamma x^2) x^{\ell},$$

$$\lambda_{\ell} = \frac{(2\gamma)^{\ell}}{\ell!}.$$

$\Rightarrow \lambda$  is larger than 0, therefore RBF is positive definite.