

## Kernalized linear regression

If we are using basis functions, we are transforming the input space.

e.g. if we want polynomial of degree  $n$  then the basis function becomes more and more complex

$x$	degree	$d$	$m$	$\phi(x)$
$[x_1]$	2	1	2	$[x_1^2, x_1]$
$[x_1, x_2]$	2	2	5	$[x_1^2, x_2^2, x_1 x_2, x_1, x_2]$
$[x_1, x_2, x_3]$	2	3	9	$[x_1^2, x_2^2, x_3^2, x_1 x_2, x_2 x_3, x_1 x_3, \dots]$
$k$	$d$		$k+d$	$C_d - 1$

The amount of computation grows with combination, in terms of factorial.

We are doing the computation on millions of datapoints

e.g.  $n: 1000, d = 3$  (cubic)  
 $d = 5$

then  $x = 1000 \times 5 = 5000$  element matrix

$\Phi = {}^{5+3}C_5 \times 1000 = 56,000$  element matrix

Time to compute the kernel is high

Time to compute OLS from such a huge matrix also is high

What if we want to fit a polynomial with degree 15?  
 the degree of  $\Phi$  would be 816!

## Another form of OLS

$$W = (x^T x)^{-1} x^T y, \quad p \text{ is a datapoint}$$

$$\hat{y} = p^T W = p^T (x^T x)^{-1} x^T y$$

$1 \times d \quad d \times d \quad d \times 1$

$$\text{Using identity } (x^T x)^{-1} x^T = x^T (x x^T)^{-1}$$

$$\hat{y} = p^T x^T (x x^T)^{-1} y$$

$1 \times d \quad d \times n \quad n \times d \quad d \times 1$

But we can write  $p^T x^T$  as

$$p^T x^T = [p_1, p_2, \dots, p_d] \begin{bmatrix} x_{11} & x_{12} & \dots & x_{1d} \\ x_{21} & x_{22} & & x_{2d} \\ \vdots & \vdots & & \vdots \\ x_{n1} & x_{n2} & \dots & x_{nd} \end{bmatrix}$$

$$= [p_1, p_2, \dots, p_d] \begin{bmatrix} 1 & 1 & \dots & 1 \\ x_1 & x_2 & \dots & x_n \\ 1 & 1 & \dots & 1 \end{bmatrix}$$

where  $x_i \in \mathbb{R}^d$  is a datapoint  $i \in (0, n)$

$$\therefore p^T x^T = \left[ [p_1, p_2, \dots, p_d] \begin{bmatrix} 1 \\ x_1 \\ 1 \end{bmatrix}, [p_1, p_2, \dots, p_d] \begin{bmatrix} 1 \\ x_2 \\ 1 \end{bmatrix}, \dots \right]$$

$$= [p^T x_1, p^T x_2, \dots, p^T x_n]$$

$$\therefore \hat{y} = [p^T x_1, p^T x_2, \dots, p^T x_n] \cdot (x x^T)^{-1} y$$

$n \times d \times d \quad d \times 1$

Let  $(x x^T)^{-1} y$  be  $\alpha$ .  $\alpha \in \mathbb{R}^{n \times 1}$

$$\therefore \hat{y} = [p^T x_1, p^T x_2, \dots, p^T x_n] \begin{bmatrix} \alpha_1 \\ \alpha_2 \\ \vdots \\ \alpha_n \end{bmatrix}$$

$$= \sum_{i=0}^n p^T x_i \alpha_i, \quad \text{Both terms are scalar}$$

$$\hat{y} = \sum_{i=0}^n \alpha_i (p^T x_i)$$

$$\text{Now, } \alpha = (x x^T)^{-1} y$$

$$x x^T = \begin{bmatrix} -x_1 - \\ -x_2 - \\ \vdots \\ -x_n - \end{bmatrix} \begin{bmatrix} 1 & 1 & \dots & 1 \\ x_1 & x_2 & \dots & x_n \\ 1 & 1 & \dots & 1 \end{bmatrix}$$

$$= \begin{bmatrix} x_1^T x_1 & x_1^T x_2 & \dots & x_1^T x_n \\ x_2^T x_1 & x_2^T x_2 & \dots & x_2^T x_n \\ \vdots & \vdots & & \vdots \\ x_n^T x_1 & x_n^T x_2 & \dots & x_n^T x_n \end{bmatrix}$$

Notice that whenever  $x$  comes, it comes in the form of  $x_i^T x_j$

# The Kernel trick

If we use  $x = \phi(x)$  then we get

$$\hat{y} = \sum_i^n \alpha_i \phi(p)^T \phi(x_i) \quad x_i \in \mathbb{R}^d$$

$$\phi: \mathbb{R}^d \rightarrow \mathbb{R}^m$$

$$\alpha = (\Phi^T \Phi)^{-1} y$$

Instead of calculating  $\Phi$  everytime, why don't we directly calculate  $\alpha$

Instead of ever computing  $\phi(x_i)$  we directly calculate  $\phi(x_i)^T \phi(x_i)$

That is we directly calculate a bivariate function

$$k(x_i, x_j) = \phi(x_i)^T \phi(x_j)$$

$$\therefore \hat{y} = \sum_i^n \alpha_i k(p, x_i)$$

$$\alpha = (\Phi^T \Phi)^{-1} y$$

$$\text{But } \Phi^T \Phi = \begin{bmatrix} \Phi_1^T \Phi_1 & \Phi_1^T \Phi_2 & \dots & \Phi_1^T \Phi_n \\ \Phi_2^T \Phi_1 & \Phi_2^T \Phi_2 & \dots & \Phi_2^T \Phi_n \\ \vdots & \vdots & \ddots & \vdots \\ \Phi_n^T \Phi_1 & \Phi_n^T \Phi_2 & \dots & \Phi_n^T \Phi_n \end{bmatrix}$$

$$= \begin{bmatrix} \phi(x_1)^T \phi(x_1) & \phi(x_1)^T \phi(x_2) & \dots & \phi(x_1)^T \phi(x_n) \\ \vdots & \vdots & \ddots & \vdots \\ \phi(x_n)^T \phi(x_1) & \phi(x_n)^T \phi(x_2) & \dots & \phi(x_n)^T \phi(x_n) \end{bmatrix}$$

$$= \begin{bmatrix} k(x_1, x_1) & k(x_1, x_2) & \dots & k(x_1, x_n) \\ k(x_2, x_1) & k(x_2, x_2) & \dots & k(x_2, x_n) \\ \vdots & \vdots & \ddots & \vdots \\ k(x_n, x_1) & k(x_n, x_2) & \dots & k(x_n, x_n) \end{bmatrix} = K$$

$$\therefore \alpha = K^{-1} y \quad K \in \mathbb{R}^{n \times n}$$

So instead of calculating  $\Phi$ , then  $\Phi^T \Phi$ , we can directly calculate  $K$

Eg.  $\rightarrow$

$$\phi(x) = [x_1 \ x_2 \ x_1 x_2 \ x_1^2 \ x_2^2]^T$$

$$\phi(x)^T \phi(x') = x_1 x'_1 + x_2 x'_2 + x_1 x'_2 x'_1 x'_2 + x_1^2 x'_1^2 + x_2^2 x'_2^2$$

$$k(x, x') = (x^T x' + 1)^2$$

$$= (x_1 x'_1 + x_2 x'_2 + x_1 x'_2 x'_1 x'_2 + x_1^2 x'_1^2 + x_2^2 x'_2^2)^2$$

Both give the same terms (constants don't matter) but number of computations for  $k(x, x')$  is much lesser (only one squaring)

## kernel trick

The advantage of using the kernel trick is that we don't really ever calculate  $\phi(x)$ . Hence we never need to find it in closed form.

This allows for complex cases like the RGF kernel which has  $\infty$  dimensions.

Properties of kernel function  $\rightarrow$

① Symmetry  $k(x_1, x_2) = k(x_2, x_1)$

② Positive semidefiniteness of Gram matrix  $z^T K z \geq 0 \quad \forall z \in \mathbb{R}^n$

Gram Matrix

$$K = \Phi^T \Phi = \begin{bmatrix} \Phi_1^T \Phi_1 & \Phi_1^T \Phi_2 & \dots & \Phi_1^T \Phi_n \\ \Phi_2^T \Phi_1 & \Phi_2^T \Phi_2 & \dots & \Phi_2^T \Phi_n \\ \vdots & \vdots & \ddots & \vdots \\ \Phi_n^T \Phi_1 & \Phi_n^T \Phi_2 & \dots & \Phi_n^T \Phi_n \end{bmatrix}$$

By Markov's theorem, if the above two properties are satisfied, then there exists a feature map such that

$$k(x, x') = \Phi(x)^T \Phi(x')$$

---

Common types of kernels

① stationary kernels

$$k(x, x') = \text{function of } (x - x')$$

Invariant to translations in input space

② Homogeneous / Radial basis functions

$$k(x, x') = \text{function of } \|x - x'\|$$

depend only on the magnitude of the distance

# Property of $\kappa$

If  $\kappa(x_i, x_j)$  can be broken down as  $\phi(x_i)^\top \phi(x_j)$   
then it is a kernel function.

$\kappa(x_i, x_j)$  is a kernel function iff

$$\sum_{i=1}^n \sum_{j=1}^n a_i a_j \kappa(x_i, x_j) \geq 0 \quad \forall a_i, a_j \in \mathbb{R}$$

Note  $\rightarrow \kappa$  is symmetric  $\kappa(a, b) = \kappa(b, a)$

Proof  $\rightarrow A \rightarrow B$

$$\text{if } \kappa(x) = \phi(x_i)^\top \phi(x_j)$$

$$\text{then } \sum_i^n \sum_j^n a_i a_j \kappa(x_i, x_j) = \sum_i^n \sum_j^n a_i a_j \phi(x_i)^\top \phi(x_j)$$

$$= \left( \sum_j^n a_j \phi(x_j) \right) (a_0 \phi(x_0)) + \left( \sum_j^n a_j \phi(x_j) \right) (a_1 \phi(x_1)) + \dots + \left( \sum_j^n a_j \phi(x_j) \right) a_n \phi(x_n)$$

$$= \left( \sum_j^n a_j \phi(x_j) \right)^\top \left( \sum_{i=0}^n a_i \phi(x_i) \right) = \left\| \sum_{j=1}^n a_j \phi(x_j) \right\|_2^2$$

This is always  $\geq 0$ .

Hence if  $\kappa(x) = \phi(x_i)^\top \phi(x_j)$  then

$$\sum_{i=1}^n \sum_{j=1}^n a_i a_j \kappa(x_i, x_j) \geq 0 \quad \forall a_i, a_j \in \mathbb{R}$$

Proof  $B \rightarrow A$

We can write  $\sum a_i a_j \kappa(x_i, x_j)$  as

$$\begin{bmatrix} a_1 & a_2 & \dots & a_n \end{bmatrix} \begin{bmatrix} \kappa(x_1, x_1) & \kappa(x_1, x_2) & \dots & \kappa(x_1, x_n) \\ \kappa(x_2, x_1) & \kappa(x_2, x_2) & \dots & \vdots \\ \vdots & & & \kappa(x_n, x_1) \end{bmatrix} \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_n \end{bmatrix} \geq 0$$

$$a^\top K a \geq 0$$

$\therefore K$  is P.S.D.,  $\therefore K$  can be factorized as  $\phi \phi^\top$

$\therefore \kappa(x_i, x_j)$  can be written as  $\phi(x_i)^\top \phi(x_j)$

## Corollary Properties

$$\textcircled{1} \quad k_{ii}^2 \leq k_{ii} \quad k_{jj}$$

$$\text{Proof} \rightarrow k_{ii}^2 = (\phi^\top(x_i) \phi(x_i))^2$$

$$k_{ii} = \phi^\top(x_i) \phi(x_i) = \|\phi(x_i)\|^2$$

$$k_{jj} = \phi^\top(x_j) \phi(x_j) = \|\phi(x_j)\|^2$$

By Cauchy-Schwarz

$$(\phi^\top(x_i) \phi(x_j))^2 \leq \|\phi(x_i)\|^2 \|\phi(x_j)\|^2$$

$$\textcircled{2} \quad \text{Diagonal positivity}$$

$$k = V \Lambda V^\top \quad (\text{eigendecomposition})$$

$$\Lambda \geq 0 \quad (\text{P.S.D.})$$

$V \rightarrow \text{eigenvectors}$

$\phi$  can be constructed as

$$\phi(x_i) = \sqrt{\lambda_i} V_i^\top$$

$$k_{ii} = \lambda_i \|V_i\|_2^2$$

$$\text{Proof} \rightarrow \phi = \begin{bmatrix} \sqrt{\lambda_1} V_1^\top \\ \sqrt{\lambda_2} V_2^\top \\ \vdots \\ \sqrt{\lambda_n} V_n^\top \end{bmatrix}$$

$$\phi^\top \phi = \begin{bmatrix} \lambda_1 & V_1 V_1^\top \\ \lambda_2 & V_2 V_2^\top \\ \vdots & \\ \lambda_n & V_n V_n^\top \end{bmatrix}$$

$$= \begin{bmatrix} \lambda_1 V_1 & \lambda_2 V_2 & \cdots & \lambda_n V_n \end{bmatrix} \begin{bmatrix} V_1^\top \\ V_2^\top \\ \vdots \\ V_n^\top \end{bmatrix}$$

$$= [V_1 \ V_2 \ \cdots \ V_n] \begin{bmatrix} \lambda_1 & 0 & 0 & \cdots \\ 0 & \lambda_2 & 0 & \cdots \\ \vdots & & & \\ 0 & 0 & \cdots & \lambda_n \end{bmatrix} \begin{bmatrix} V_1^\top \\ V_2^\top \\ \vdots \\ V_n^\top \end{bmatrix}$$

=  $V \Lambda V^\top$  Hence proved

Note  $\rightarrow$  because  $k$  is symmetric & P.S.D.

SVD gives  $V \Sigma V^\top$  with  $V = U \quad \left\{ \begin{array}{l} \Sigma = \Lambda \\ \Sigma = I \end{array} \right\}$  same as eigendecomposition

$\infty$  dimensional vectors

Consider  $x = [a, a^2, \dots, a^\infty]^T$ ,  $a \in \mathbb{R}$

$x \in \mathbb{R}^\infty$  is an  $\infty$  dimensional vector

what is  $x^T x$ ?  $x^T x = \sum_{n=1}^{\infty} a^{2n} = \frac{1}{1-2a}$

Even if we can't represent an  $\infty$  dimensional vector, we can represent its dot product with another  $\infty$  dimensional vector!

Even if we don't have closed form for  $\phi(x)$  we can have a closed form for  $\phi(x) \cdot \phi(y)$

That means even if we don't know  $\phi(x)$  but we know  $\phi(x) \phi(y)$  it's okay.

Or rather we don't need to know  $\phi(x)$

# Composition of Kernels

Given 2 valid kernels  $k_1(x, x')$  &  $k_2(x, x')$   
the following is also a kernel

①  $c k_1(x, x')$

②  $f(x) k_1(x, x') f(x')$

③  $q(k_1(x, x'))$

( $q \rightarrow$  polynomial with non negative coefficients)

④  $\exp(k_1(x, x'))$

⑤  $k_1(x, x') + k_2(x, x')$

⑥  $k_1(x, x') k_2(x, x')$

⑦  $k(f(x), f(x'))$

$$f: \mathbb{R}^d \rightarrow \mathbb{R}^m$$

For more details on "Kernel Engineering", see

shawe-taylor & Cristianini (2004)

## Kernalized Ridge Regression

We have

$$w^* = \underset{w}{\operatorname{argmin}} \|y - x w\|_2^2 + \lambda \|w\|_2^2$$

$x w$  is an  $N$  dimensional vector

$$x w = \begin{bmatrix} x_1^T w \\ x_2^T w \\ \vdots \\ x_N^T w \end{bmatrix} \quad \begin{array}{l} x_i \in \mathbb{R}^d \\ w \in \mathbb{R}^d \end{array}$$

By orthogonal decomposition theorem, we can write  $w$  as

$$w = w_{||} + w_{\perp}$$

Parallel component      Perpendicular

$$w_{||} \in \text{Span of rows of } x$$

$$w_{\perp} \perp \text{Span of rows of } x$$

$$x w_{\perp} = 0$$

$$w_{||} \in \mathbb{R}^d, w_{\perp} \in \mathbb{R}^d$$

$$\therefore x w = x(w_{||} + w_{\perp}) = x w_{||} + x w_{\perp} = x w_{||}$$

By Pythagoras theorem

$$\|w\|_2^2 = \|w_{||}\|_2^2 + \|w_{\perp}\|_2^2$$

Hence objective becomes

$$\underset{w_{||}, w_{\perp}}{\operatorname{argmin}} \|y - x w_{||}\|_2^2 + \|w_{||}\|_2^2 + \|w_{\perp}\|_2^2$$

Now we have 2 separate variables  $w_{||}$  &  $w_{\perp}$  to optimize over

Since  $\|w_{\perp}\|^2 \geq 0$ , the objective is strictly increasing. Hence the minimum occurs when  $w_{\perp} = 0$

That is, the optimal  $w_{\perp} = 0$ .

Hence the optimal  $w^*$  has  $w_{\perp} = 0$ , i.e.  $w^*$  lies entirely in span of  $x$

Since  $w^* \in \text{span}(x_1, x_2, \dots, x_N)$

$$w^* = x^T \alpha \text{ for some } \alpha \in \mathbb{R}^N$$

(by span property)

This is called as "representer trick"

$$\begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_d \end{bmatrix} = \begin{bmatrix} x_{11} & x_{12} & \dots & \dots & x_{1N} \\ x_{21} & \dots & - & & x_{2N} \\ \vdots & & & & \vdots \\ x_{d1} & \dots & \dots & & x_{dN} \end{bmatrix} \begin{bmatrix} \alpha_1 \\ \alpha_2 \\ \vdots \\ \alpha_N \end{bmatrix}$$

$$= \begin{bmatrix} \sum_{i=0}^N x_{1i} \alpha_i \\ \sum_{i=0}^N x_{2i} \alpha_i \\ \vdots \\ \sum_{i=0}^N x_{di} \alpha_i \end{bmatrix} \rightarrow \text{2nd dimension}$$

$$\text{i.e. } w_k = \sum_{i=0}^N x_{ki} \alpha_i$$

$\downarrow k^{\text{th}}$  dimension of  $x$

## Representer Theorem

Whenever  $w$  comes in dot products, we can use the kernel trick. This is called the representer theorem.

The training data lives in subspace spanned by  $\{\phi(x_1), \phi(x_2), \dots, \phi(x_n)\}$

Any component outside the subspace doesn't affect the loss.

Hence  $w^*$  too lies in the span

Hence  $w^*$  can be written as a linear combination of mapped training points

$$w = \sum_i \alpha_i \phi(x_i)$$

## Kernalized Ridge Regression

$$\text{loss} = \|y - xw\|_2^2 + \lambda \|w\|_2^2 \quad \text{let } w = x^\top \alpha$$

$$\text{loss} = (y - x x^\top \alpha)^\top (y - x x^\top \alpha) + \lambda (x^\top \alpha)^\top (x^\top \alpha)$$

at minima

$$\frac{\partial \text{loss}}{\partial \alpha} = 0$$

$$\frac{\partial \text{loss}}{\partial \alpha} = -2x^\top (y - x x^\top \alpha) + 2\lambda x^\top x = 0$$

$$\therefore (x^\top x + \lambda I) \alpha = y$$

$$\therefore \alpha = (x^\top x + \lambda I)^{-1} y$$

as we have seen  $k = x x^\top$

$$\therefore \alpha = (k + \lambda I)^{-1} y$$

for inference on point  $p$

$$\hat{y} = p^\top w = \sum_{i=1}^d p_i w_i$$

$$\hat{y} = p^\top x^\top \alpha$$

$(d \times d) \quad d \times N \quad N \times 1$

$$= [p_1 \ p_2 \ \dots \ p_d] \begin{bmatrix} x_{11} & x_{12} & \dots & \dots & x_{1N} \\ x_{21} & \dots & - & & x_{2N} \\ \vdots & & & & \vdots \\ x_{d1} & \dots & & & x_{dN} \\ - & & & & \end{bmatrix} \begin{bmatrix} \alpha_1 \\ \alpha_2 \\ \vdots \\ \alpha_N \end{bmatrix}$$

$$= \left[ \sum_{i=1}^d p_i x_{i1} \quad \sum_{i=1}^d p_i x_{i2} \quad \dots \quad \sum_{i=1}^d p_i x_{iN} \right] \alpha$$

$$= [p^\top x_1 \ p^\top x_2 \ \dots \ p^\top x_N] \alpha$$

$$= [k(p, x_1) \ k(p, x_2) \ \dots \ k(p, x_N)] \alpha$$

thus we can avoid all dot products and write in form of  $k$