

# Kernel Ridge Regression

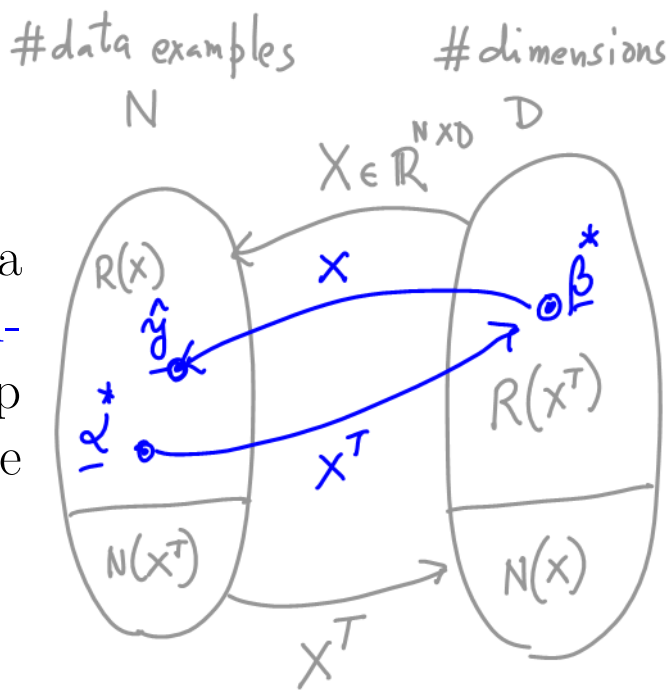
Mohammad Emtiyaz Khan  
EPFL

Oct 27, 2015



ÉCOLE POLYTECHNIQUE  
FÉDÉRALE DE LAUSANNE

©Mohammad Emtiyaz Khan 2015



## Motivation

The ridge solution  $\beta^* \in \mathbb{R}^D$  has a counterpart  $\alpha^* \in \mathbb{R}^N$ . Using [duality](#), we will establish a relationship between  $\beta^*$  and  $\alpha^*$  which leads the way to [kernels](#).

## Ridge regression

Throughout, we assume that there is no intercept term  $\beta_0$  to make the math easier.

The following is true for ridge reg:

$$\beta^* = \underbrace{(\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}_D)^{-1} \mathbf{X}^T \mathbf{y}}_{\text{A}} = \mathbf{X}^T \underbrace{(\mathbf{X} \mathbf{X}^T + \lambda \mathbf{I}_N)^{-1} \mathbf{y}}_{\text{B} \quad \alpha^* \in \mathbb{R}^N} := \mathbf{X}^T \alpha^*,$$

where  $\alpha^* := (\mathbf{X} \mathbf{X}^T + \lambda \mathbf{I}_N)^{-1} \mathbf{y}$ .

This can be proved using the following identity: let  $\mathbf{P}$  be an  $N \times M$  matrix while  $\mathbf{Q}$  be an  $M \times N$  matrix,

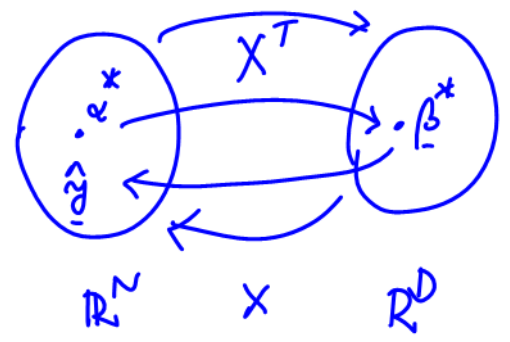
$$\underbrace{(\mathbf{X}^T \mathbf{X} + \mathbf{I}_D)^{-1}}_{\text{A}} \mathbf{X}^T \underbrace{\mathbf{X} (\mathbf{X} \mathbf{X}^T + \mathbf{I}_N)^{-1}}_{\text{B}} = \mathbf{P} (\mathbf{Q} \mathbf{P} + \mathbf{I}_M)^{-1}$$

by  
Multiply  $(\mathbf{I}_N + \mathbf{P} \mathbf{Q})^{-1}$   
&  $(\mathbf{I}_M + \mathbf{Q} \mathbf{P})^{-1}$

$$\boxed{[\mathbf{P} + \mathbf{P} \mathbf{Q} \mathbf{P}]^{-1}} = \boxed{[(\mathbf{I}_N + \mathbf{P} \mathbf{Q}) \mathbf{P}]^{-1}} = \boxed{[\mathbf{P} (\mathbf{I}_M + \mathbf{Q} \mathbf{P})]^{-1}}$$

What are the computational complexities for these two ways of computing  $\beta^*$ ?

For **A**:  $O(D^2 N + D^3)$ , For **B**:  $O(N^2 D + N^3)$



With this, we know that  $\beta^* = \mathbf{X}^T \alpha^*$  lies in the row space of  $\mathbf{X}$ . Previously, we have seen that  $\hat{\mathbf{y}} = \mathbf{X}\beta^*$  lies in the column space of  $\mathbf{X}$ . In other words,

$$\beta^* = \sum_{n=1}^N \alpha_n^* \mathbf{x}_n, \quad \hat{\mathbf{y}} = \sum_{d=1}^D \beta_d^* \bar{\mathbf{x}}_d$$

$$\text{where } \mathbf{X} = \begin{bmatrix} x_{11} & x_{12} & \dots & x_{1D} \\ x_{21} & x_{22} & \dots & x_{2D} \\ \vdots & \vdots & \ddots & \vdots \\ x_{N1} & x_{N2} & \dots & x_{ND} \end{bmatrix} = \begin{bmatrix} \mathbf{x}_1^T \\ \mathbf{x}_2^T \\ \vdots \\ \mathbf{x}_N^T \end{bmatrix} \\ = [\bar{\mathbf{x}}_1 \ \bar{\mathbf{x}}_2 \ \dots \ \bar{\mathbf{x}}_D]$$

## The representer theorem

The representer theorem generalizes this result: for a  $\beta^*$  minimizing the following function for any  $\mathcal{L}$ ,

$$\min_{\beta} \sum_{n=1}^N \mathcal{L}(y_n, \mathbf{x}_n^T \beta) + \sum_{j=1}^D \lambda \beta_j^2$$

there exists  $\alpha^*$  such that  $\beta^* = \mathbf{X}^T \alpha^*$ .

See even more general statement in [Wikipedia](#), originally proved in Scholkopf, Herbrich and Smola (2001).

$\beta^* = \mathbf{X}^T \alpha$  implies for Ridge regression that

$$\mathbf{x}_*^T \beta^* = \sum_{n=1}^N \alpha_n \mathbf{x}_*^T \mathbf{x}_n$$

Optimal predictor (pointing to  $\beta^*$ ), test input (pointing to  $\mathbf{x}_*$ ), train i/p (pointing to  $\mathbf{x}_n$ )

train inputs (pointing down to the next equation)

In general,

$$f^*(\mathbf{x}_*) = \sum_{n=1}^N \alpha_n k(\mathbf{x}_*, \mathbf{x}_n)$$

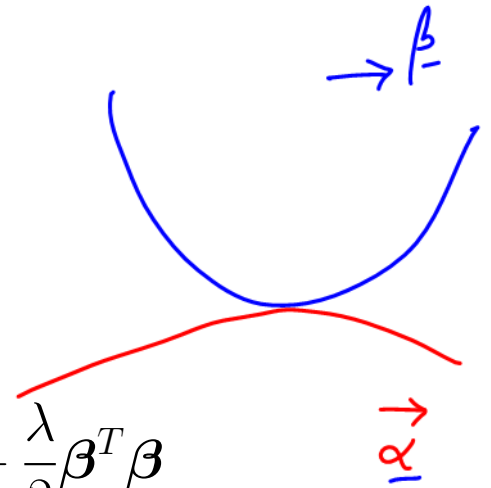
optimal predictor (pointing to  $f^*$ ), test input (pointing to  $\mathbf{x}_*$ ), kernel (pointing to  $k$ )

2 Therefore, no need to work with  $\mathbf{X}$ .

# Kernelized ridge regression

The representer theorem allows us to write an equivalent optimization problem in terms of  $\boldsymbol{\alpha}$ . For example, for ridge regression, the following two problems are equivalent:

$$\boldsymbol{\beta}^* = \arg \min_{\boldsymbol{\beta}} \quad \frac{1}{2}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T(\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) + \frac{\lambda}{2}\boldsymbol{\beta}^T\boldsymbol{\beta}$$
$$\boldsymbol{\alpha}^* = \arg \max_{\boldsymbol{\alpha}} \quad -\frac{1}{2}\boldsymbol{\alpha}^T(\mathbf{X}\mathbf{X}^T + \lambda\mathbf{I}_N)^T\boldsymbol{\alpha} + \boldsymbol{\alpha}^T\mathbf{y}$$



i.e. they both return the same optimal value and there is a one-to-one mapping between  $\boldsymbol{\alpha}^*$  and  $\boldsymbol{\beta}^*$ . Note that the optimization over  $\boldsymbol{\alpha}$  is a *maximization* problem.

Most importantly, the second problem is expressed in terms of the matrix  $\mathbf{X}\mathbf{X}^T$ . This is our first example of a [kernel](#) matrix.

Note: We don't give a detailed derivation of the second problem, but to show the equivalence, you can show that we obtain equal optimal values for the two problems. We will do a derivation later using the duality principle. You can see a derivation here [http://www.ics.uci.edu/~welling/classnotes/papers\\_class/Kernel-Ridge.pdf](http://www.ics.uci.edu/~welling/classnotes/papers_class/Kernel-Ridge.pdf).

# Advantages of kernelized ridge regression

First, it might be computationally efficient in some cases when solving the system of equations.

$$O(N D^2 + D^3), \quad N \gg D$$

$$O(N^2 D + N^3), \quad N \ll D$$

Second, by defining  $\mathbf{K} = \mathbf{X}\mathbf{X}^T$ , we can work directly with  $\mathbf{K}$  and never have to worry about  $\mathbf{X}$ . This is the kernel trick.

②

→ explained in the following sections.

Third, working with  $\boldsymbol{\alpha}$  is sometimes advantageous (e.g. in SVMs many entries of  $\boldsymbol{\alpha}$  will be zero).

③

→ We will see this in SVM

## Kernel functions

The linear kernel is defined below:

$$\mathbf{K} = \mathbf{X}\mathbf{X}^T = \begin{bmatrix} \mathbf{x}_1^T \mathbf{x}_1 & \mathbf{x}_1^T \mathbf{x}_2 & \dots & \mathbf{x}_1^T \mathbf{x}_N \\ \mathbf{x}_2^T \mathbf{x}_1 & \mathbf{x}_2^T \mathbf{x}_2 & \dots & \mathbf{x}_2^T \mathbf{x}_N \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{x}_N^T \mathbf{x}_1 & \mathbf{x}_N^T \mathbf{x}_2 & \dots & \mathbf{x}_N^T \mathbf{x}_N \end{bmatrix} \cdot \mathbf{X} = \begin{bmatrix} \mathbf{x}_1^T \\ \mathbf{x}_2^T \\ \vdots \\ \mathbf{x}_N^T \end{bmatrix}$$

Kernel with basis functions  $\phi(\mathbf{x}) : \mathbb{R}^D \rightarrow \mathbb{R}^M$  with  $\mathbf{K} := \boldsymbol{\Phi}\boldsymbol{\Phi}^T$  is shown below:

$$\begin{bmatrix} \phi(\mathbf{x}_1)^T \phi(\mathbf{x}_1) & \phi(\mathbf{x}_1)^T \phi(\mathbf{x}_2) & \dots & \phi(\mathbf{x}_1)^T \phi(\mathbf{x}_N) \\ \phi(\mathbf{x}_2)^T \phi(\mathbf{x}_1) & \phi(\mathbf{x}_2)^T \phi(\mathbf{x}_2) & \dots & \phi(\mathbf{x}_2)^T \phi(\mathbf{x}_N) \\ \vdots & \vdots & \ddots & \vdots \\ \phi(\mathbf{x}_N)^T \phi(\mathbf{x}_1) & \phi(\mathbf{x}_N)^T \phi(\mathbf{x}_2) & \dots & \phi(\mathbf{x}_N)^T \phi(\mathbf{x}_N) \end{bmatrix} \cdot$$

## The kernel trick

A big advantage of using kernels is that we do not need to specify  $\phi(\mathbf{x})$  explicitly, since we can work directly with  $\mathbf{K}$ .

We will use a kernel function  $k(\mathbf{x}, \mathbf{x}')$  and compute the  $(i, j)$ th entry of  $\mathbf{K}$  as follows:  $K_{ij} = k(\mathbf{x}_i, \mathbf{x}_j)$ . For example, for linear kernel and basis function expansion, the kernel function is the following:

$$k(\mathbf{x}, \mathbf{x}') := \mathbf{x}^T \mathbf{x}', \quad k(\mathbf{x}, \mathbf{x}') := \phi(\mathbf{x})^T \phi(\mathbf{x}')$$

However, a kernel function  $k$  is usually associated with a  $\phi$ , e.g.  $k(x, x') = x^2(x')^2$  corresponds to  $\phi(x) = x^2$  and  $k(\mathbf{x}, \mathbf{x}') = (x_1x'_1 + x_2x'_2 + x_3x'_3)^2$  corresponds to

$$\phi(\mathbf{x})^T = [x_1^2, x_2^2, x_3^2, \sqrt{2}x_1x_2, \sqrt{2}x_1x_3, \sqrt{2}x_2x_3]$$

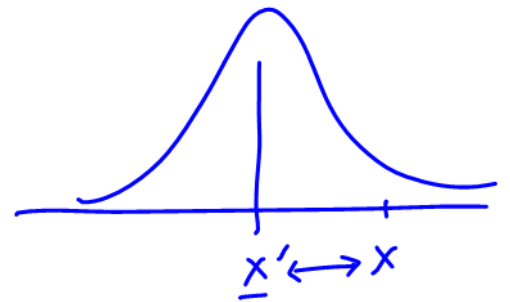
The good news is that the evaluation of a kernel is usually faster with  $k$  than with  $\phi$ .

Ex! Show that this  $\phi$  leads to this  $k$ .

↑ Check it in the above example-

## Examples of kernels

The above kernel is an example of the [polynomial kernel](#). Another example is the [Radial Basis Function \(RBF\) kernel](#).



$$k(\mathbf{x}, \mathbf{x}') = \exp \left[ -\frac{1}{2}(\mathbf{x} - \mathbf{x}')^T (\mathbf{x} - \mathbf{x}') \right]$$

what is  $\phi$ ? well, we don't care!

↳ See more examples in Section [14.2](#) of the KPM book.

A natural question is the following: how can we ensure that there exists a  $\phi$  corresponding to a given kernel  $\mathbf{K}$ ? The answer is: as long as a kernel satisfy certain properties.

## Properties of a kernel

A kernel function must be an inner-product in some feature space. Here are a few properties that ensure it is the case.

1.  $\mathbf{K}$  should be symmetric, i.e.  $k(\mathbf{x}, \mathbf{x}') = k(\mathbf{x}', \mathbf{x})$ .
2. For any arbitrary input set  $\{\mathbf{x}_n\}$  and all  $N$ ,  $\mathbf{K}$  should be positive semidefinite.

Use <sup>the</sup> definition of p.s.d.

$$\underline{t}^T \mathbf{K} \underline{t} \geq 0, \text{ to show}$$

$$\sum_{ij} t_i x_i^2 x_j^2 t_j > 0.$$

An important subclass is the [positive-definite kernel functions](#), giving rise to infinite-dimensional feature spaces.

Read about Mercer and Matern kernels from Kevin Murphy's Section 14.2. There is a small note about [Reproducing kernel Hilbert Space](#) in the website (written by Matthias Seeger), please read that as well.

## To do

1. Clearly understand the relationship  $\boldsymbol{\beta}^* = \mathbf{X}^T \boldsymbol{\alpha}^*$ . Understand the statement of the representer theorem.
2. Show that ridge regression and kernel ridge regression are equivalent. Hint: show that the optimization problems corresponding to  $\boldsymbol{\beta}$  and  $\boldsymbol{\alpha}$  have the same optimal value.
3. Get familiar with various examples of kernels. See Section 6.2 of Bishop on examples of kernel construction. Read Section 14.2 of KPM book for examples of kernels.
4. Revise and understand the difference between positive-definite and positive-semidefinite matrices.
5. If curious about infinite  $\phi$ , see Matthias Seeger's notes (uploaded on the website).