

Machine Learning Course - CS-433

Kernel Ridge Regression and the Kernel Trick

Oct 31, 2019

changes by Martin Jaggi 2019, changes by Rüdiger Urbanke 2018, changes by Martin Jaggi
2016, 2017 ©Mohammad Emtiyaz Khan 2015

Last updated on: October 31, 2019



Motivation

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

Ridge regression

Recall the ridge regression problem

$$\min_{\mathbf{w}} \quad \frac{1}{2} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|^2 + \frac{\lambda}{2} \|\mathbf{w}\|^2$$

For its solution, we have that

$$\begin{aligned} \mathbf{w}^* &= (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}_D)^{-1} \mathbf{X}^\top \mathbf{y} \\ &= \mathbf{X}^\top (\mathbf{X} \mathbf{X}^\top + \lambda \mathbf{I}_N)^{-1} \mathbf{y} =: \mathbf{X}^\top \boldsymbol{\alpha}^*, \end{aligned}$$

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

This can be proved using the following identity: let \mathbf{P} be an $N \times D$ matrix while \mathbf{Q} be $D \times N$, and let both $\mathbf{PQ} + \mathbf{I}$ and $\mathbf{QP} + \mathbf{I}$ be invertible.

$$(\mathbf{PQ} + \mathbf{I}_N)^{-1} \mathbf{P} = \mathbf{P}(\mathbf{QP} + \mathbf{I}_D)^{-1}$$

What are the [computational complexities](#) for the above two ways of computing \mathbf{w}^* ?

With this, we know that $\mathbf{w}^* = \mathbf{X}^\top \boldsymbol{\alpha}^*$ lies in the column space of \mathbf{X}^\top ,

$$\text{where } \mathbf{X} = \begin{bmatrix} x_{11} & x_{12} & \dots & x_{1D} \\ x_{21} & x_{22} & \dots & x_{2D} \\ \vdots & \vdots & \ddots & \vdots \\ x_{N1} & x_{D2} & \dots & x_{ND} \end{bmatrix}$$

The representer theorem

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

$$\min_{\mathbf{w}} \sum_{n=1}^N \mathcal{L}_n(\mathbf{x}_n^\top \mathbf{w}, y_n) + \frac{\lambda}{2} \|\mathbf{w}\|^2$$

there exists $\boldsymbol{\alpha}^*$ such that $\mathbf{w}^* = \mathbf{X}^\top \boldsymbol{\alpha}^*$.

Such a general statement was originally proved by *Schölkopf, Herbrich and Smola (2001)*.

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:

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} \quad \frac{1}{2} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|^2 + \frac{\lambda}{2} \|\mathbf{w}\|^2$$

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

i.e. they both have the same optimal value. Also, we can always have the correspondence mapping $\mathbf{w} = \mathbf{X}^\top \boldsymbol{\alpha}$.

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

Note: To see the equivalence, you can show that we obtain equal optimal values for the two problems. Take the gradient of the second expression, to get $(\mathbf{X}\mathbf{X}^\top + \lambda \mathbf{I}_N) \boldsymbol{\alpha} - \mathbf{y}$. Setting this to $\mathbf{0}$ and solving for $\boldsymbol{\alpha}$ results in

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

If we combine this with the representer theorem $\mathbf{w}^* = \mathbf{X}^\top \boldsymbol{\alpha}^*$ we find back the dual solution.

Advantages of kernelized ridge regression

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

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

Third, working with $\boldsymbol{\alpha}$ is sometimes advantageous, and provides additional information for each datapoint (e.g. as in SVMs).

Kernel functions

The [linear kernel](#) is defined below:

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

Kernel with basis functions $\phi(\mathbf{x})$
with $\mathbf{K} := \boldsymbol{\Phi}^\top \boldsymbol{\Phi}$ is shown below:

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

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 $\kappa(\mathbf{x}, \mathbf{x}')$ and compute the (i, j) -th entry of \mathbf{K} as $K_{ij} = \kappa(\mathbf{x}_i, \mathbf{x}_j)$. A kernel function κ is usually associated with a feature map ϕ , such that

$$\kappa(\mathbf{x}, \mathbf{x}') := \phi(\mathbf{x})^\top \phi(\mathbf{x}') .$$

For example, for the linear kernel $\kappa(\mathbf{x}, \mathbf{x}') := \mathbf{x}^\top \mathbf{x}'$, the feature map is just the original features, $\phi(\mathbf{x}') = \mathbf{x}'$.

Another example: The kernel $\kappa(x, x') := x^2(x')^2$ corresponds to $\phi(x) = x^2$, and $\kappa(\mathbf{x}, \mathbf{x}') := (x_1x'_1 + x_2x'_2 + x_3x'_3)^2$ corresponds to

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

The good news is that the evaluation of a kernel is often faster when using κ instead of ϕ .

Visualization

Why would we want such general feature maps?

See video explaining linear separation in the kernel space (where $\phi(\mathbf{x})$ maps to) corresponding to non-linear separation in the original \mathbf{x} -space: <https://www.youtube.com/watch?v=3liCbRZPrZA>

Examples of kernels

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

$$\kappa(\mathbf{x}, \mathbf{x}') = \exp \left[-\frac{1}{2}(\mathbf{x} - \mathbf{x}')^\top (\mathbf{x} - \mathbf{x}') \right]$$

See more examples in Section 14.2 of Murphy's book.

A natural question is the following: how can we ensure that there exists a ϕ corresponding to a given kernel \mathbf{K} ? The answer is: as long as the kernel satisfies 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.
 $\kappa(\mathbf{x}, \mathbf{x}') = \kappa(\mathbf{x}', \mathbf{x})$.
2. For any arbitrary input set $\{\mathbf{x}_n\}$ and all N , \mathbf{K} should be positive semi-definite.

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

Exercises

1. Understand the relationship $\mathbf{w}^* = \mathbf{X}^\top \boldsymbol{\alpha}^*$. Understand the statement of the representer theorem.
2. Show that the primal and dual formulations of ridge regression are equivalent. Hint: show that the optimization problems corresponding to \mathbf{w} 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 Murphy's book for examples of kernels.
4. Revise and understand the difference between positive-definite and positive-semi-definite matrices.
5. If you are interested in more details about kernels, read about Mercer and Matern kernels from Kevin Murphy's Section 14.2. There is also a small note by Matthias Seeger on the git repository under [lectures/07](#), in particular for the case of infinite dimensional ϕ .