

Lecture 16: Kernel Ridge Regression

In the last lecture we introduced the idea of a kernel & the associated RKHS.

- A function $K: \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ is called a kernel. If it is non-negative definite & symm. then we have

$$K(\underline{x}, \underline{x}') = \sum_{j=0}^{\infty} F_j(\underline{x}) F_j(\underline{x}')$$

where $F_j: \mathbb{R}^N \rightarrow \mathbb{R}$ are the features of K .

- The space H_K , RKHS of K is defined as all $f: \mathbb{R}^d \rightarrow \mathbb{R}$ of the form $f(\underline{x}) = \sum_{j=0}^{\infty} c_j F_j(\underline{x})$ s.t $\sum_{j=0}^{\infty} c_j < +\infty$.
- Equivalently such f has the form $f(\underline{x}) = \sum_{j=0}^{\infty} a_j k(\underline{x}_j, \underline{x})$ for some a_j & \underline{x}_j .

We also looked at a simple kernel interpolation problem. The goal of this lecture is to extend the kernel methods to ridge regression.

16.1 Kernel Ridge Regression

Recall our original approach to ridge

was to solve problems of the form

(RR) minimize $\beta \cdot \|A\beta - Y\|^2 + \lambda \|\beta\|^2$

where A was our feature matrix often of the form

(P.s. we used γ_j instead of F_j in context of Ridge reg.) $A = \begin{bmatrix} F_0(\underline{x}_0) & F_1(\underline{x}_0) & \dots \\ F_0(\underline{x}_1) & F_1(\underline{x}_1) & \dots \\ \vdots & \vdots & \vdots \\ F_0(\underline{x}_{N-1}) & F_1(\underline{x}_{N-1}) & \dots \end{bmatrix}$

* essentially this meant that we would pick the features $F_j : \mathbb{R}^d \rightarrow \mathbb{R}$ leading to a model of the form

(*) $f(\underline{x}) = \sum_{j=0}^{J-1} \beta_j F_j(\underline{x}) \approx y(\underline{x})$

Our knowledge of Kernels now reveals to us that underneath this model is the kernel

(++) $K(\underline{x}, \underline{x}') = \sum_{j=0}^{J-1} F_j(\underline{x}) F_j(\underline{x}')$

with features $\gamma_j : \mathbb{R}^d \rightarrow \mathbb{R}$ & where RKHS,

$$\mathcal{H}_K := \{f : \mathbb{R}^d \rightarrow \mathbb{R} \mid f(\underline{x}) \text{ has form } (*)\}$$

and

$$\|f\|_{\mathcal{H}_K}^2 = \sum_{j=0}^{J-1} \beta_j^2 = \|\beta\|^2$$

In other words (KR) is equivalent to the problem

$$\text{minimize}_{f \in \mathcal{H}_K} \|f(x) - y\|^2 + \gamma \|f\|_{\mathcal{H}_K}^2$$

with \mathcal{H}_K induced by the kernel K above.

Observe, this (KR) problem is now defined for functions $f \in \mathcal{H}_K$ & not just vectors of parameters p .

Problem (KR) is the abstract formulation of Kernel Ridge (KR) Regression.

Let us now consider kernels that, unlike (KR) do not have finitely many features.

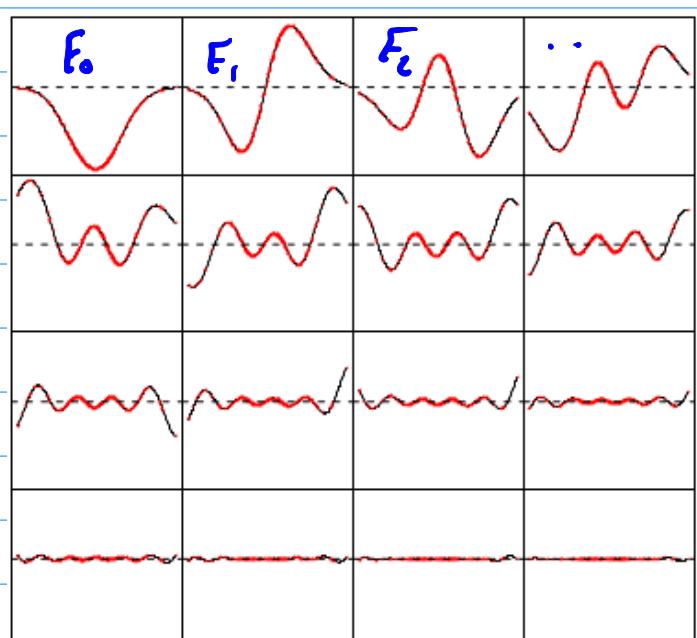
ex. the Gaussian / RBF / Squared exponential kernel

$$k(\underline{x}, \underline{x}') =$$

$$\exp(-\gamma \|\underline{x} - \underline{x}'\|^2)$$

$$= \exp\left(-\frac{\|\underline{x} - \underline{x}'\|^2}{2\sigma^2}\right)$$

$$= \sum_{j=0}^{\infty} F_j(\underline{x}) F_j(\underline{x}')$$



The same ideas that were applied to the interpolation problem last week still hold in the setting of KR. i.e., the problem KR) can be solved via a simple finite dimensional problem, even when K has infinitely many features.

This is a fundamental result in Kernel methods called a "Representer theorem".

Thm (Representer)

A function $\hat{f} \in \mathcal{H}_K$ is a minimizer of

$$\underset{f \in \mathcal{H}_K}{\text{minimize}} \|f(x) - Y\|^2 + \lambda \|f\|_{\mathcal{H}_K}^2$$

if & only if

$$\hat{f} = \sum_{n=0}^{N-1} \hat{\alpha}_n K(\underline{\alpha}_n, \underline{\alpha})$$

with $\hat{\underline{\alpha}} = (\alpha_0, \dots, \alpha_{N-1}) \in \mathbb{R}^N$ being a minimizer of

$$\underset{\underline{\alpha} \in \mathbb{R}^N}{\text{minimize}} \|\Theta_{\underline{\alpha}} - Y\|^2 + \lambda \underline{\alpha}^T \Theta \underline{\alpha}$$

where we recall $\Theta_{ij} = K(\underline{\alpha}_i, \underline{\alpha}_j)$

Similar to the case of penalized least squares this problem has an exact solution

$$\hat{\alpha} = (\Theta + \lambda I)^{-1} Y.$$

- The representer thm is a profound result & it is the reason why Kernel methods are so prevalent.

- Even though K might have infinitely many features (ex Gaussian kernel) Rep. Thm. tells us that the minimizer \hat{f} is still identified by a finite-dimensional problem for $\hat{\alpha} \in \mathbb{R}^N$.

16.2 Kernel Ridge Reg demo

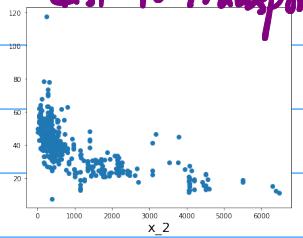
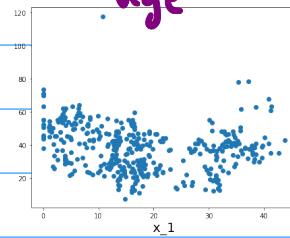
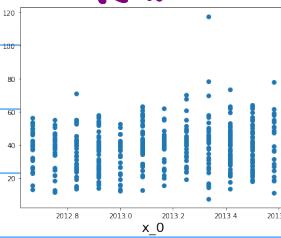
Before proceeding further with theoretical discussions let us consider an application of KR.

Housing prices in Taiwan

Year

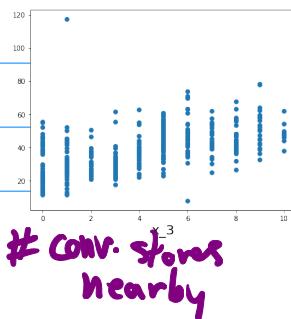
age

dist. to transport

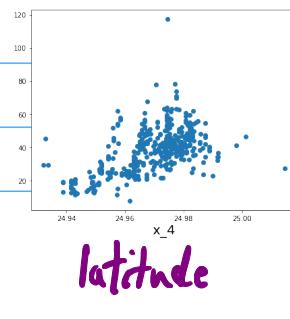


from UCI Repo. Price
Real state Valuation
dataset

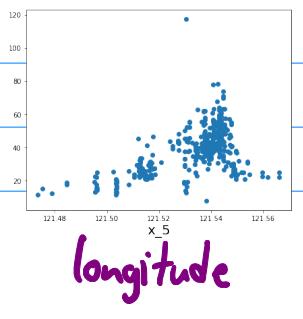
Price



conv. stores
nearby



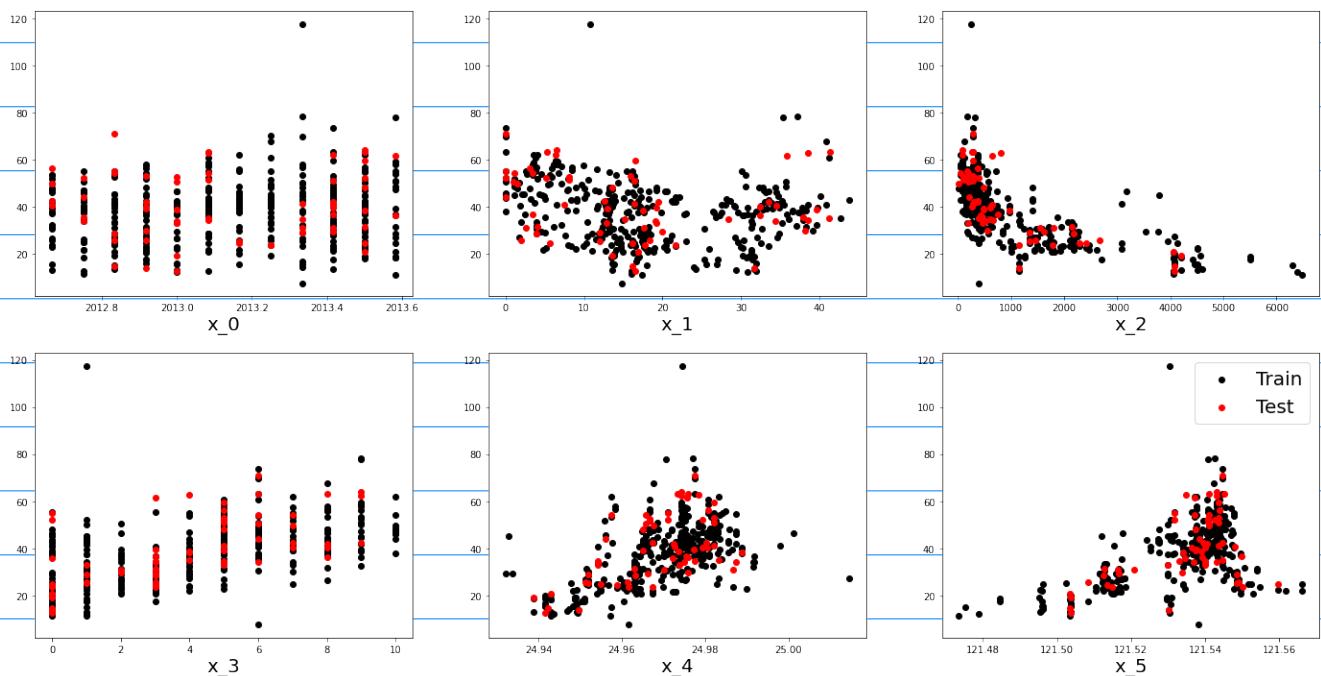
latitude



longitude

Goal: Find $\hat{f}(x_0, \dots, x_5)$ that predicts the price of a house on the market.

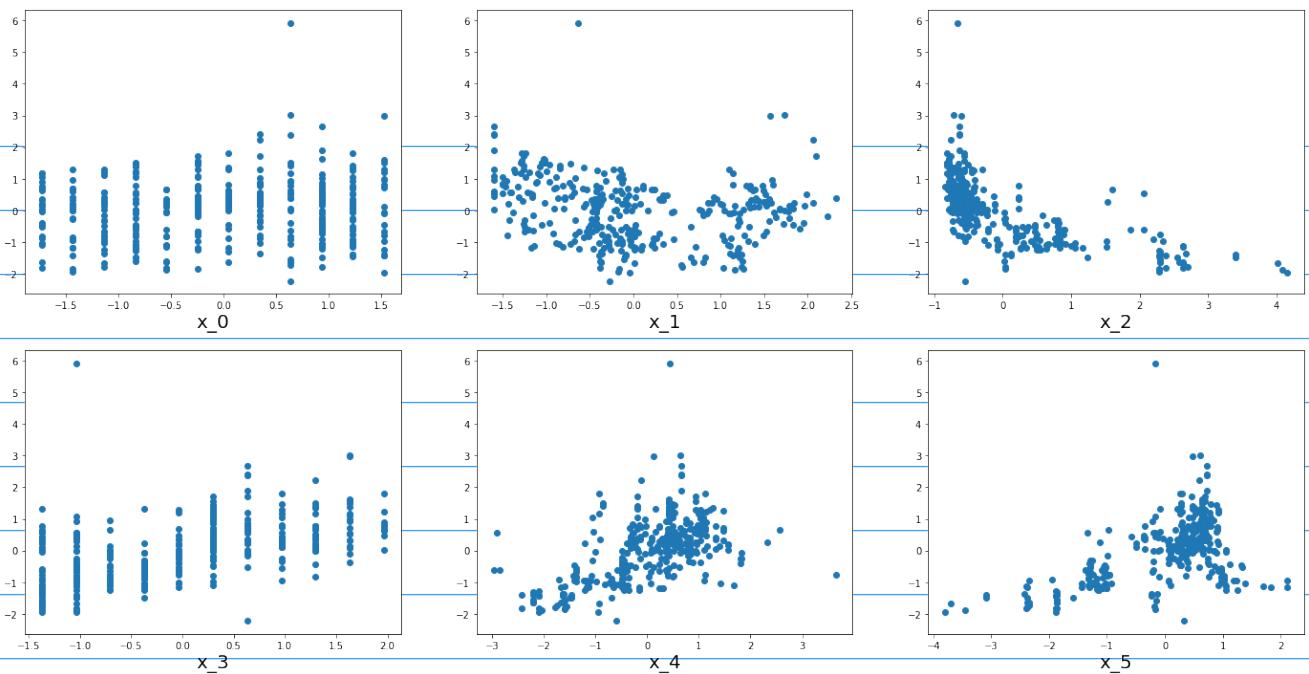
First thing we do is to split the data set ($N=414$ instances) into training & test sets. Luckily Sklearn has a function that does this for us.



Only the black pts will be used for training (finding \hat{g}).

Next we normalize & center the training data set so that each column of X , ie, $\underline{\alpha}_j$ have mean 0 & standard deviation 1. This is not crucial in theory but is an important step in practice!

Note, for ex., that $x_0 \approx 2000$'s while $x_3 \in (0, 10)$ while $y \in (0, 120)$. This makes it hard to equally weight the $\underline{\alpha}_j$.



We are now ready to train our KR model on our normalized training set.

We have some parameters to choose:

(a) kernel $K(\underline{a}, \underline{a}')$ & its possible parameters.

(b) The penalty parameter λ .

Sklearn accommodates various kernels:

$$\text{"linear"} \quad K(\underline{a}, \underline{a}') = \underline{a}^T \underline{a}'$$

$$\text{"poly"} \quad K(\underline{a}, \underline{a}') = (\gamma \underline{a}^T \underline{a}' + c_0)^d$$

$$\text{"sigmoid"} \quad K(\underline{a}, \underline{a}') = \tanh(\gamma \underline{a}^T \underline{a}' + c_0)$$

$$\text{"rbf"} \quad K(\underline{a}, \underline{a}') = \exp(-\gamma \|\underline{a} - \underline{a}'\|^2)$$

$$\text{"Laplace"} \quad K(\underline{a}, \underline{a}') = \exp(-\gamma \|\underline{a} - \underline{a}'\|)$$

You can define your own kernel as well but we will use RBF (Gaussian). We have two parameters to pass to `Sklearn`, γ & λ .

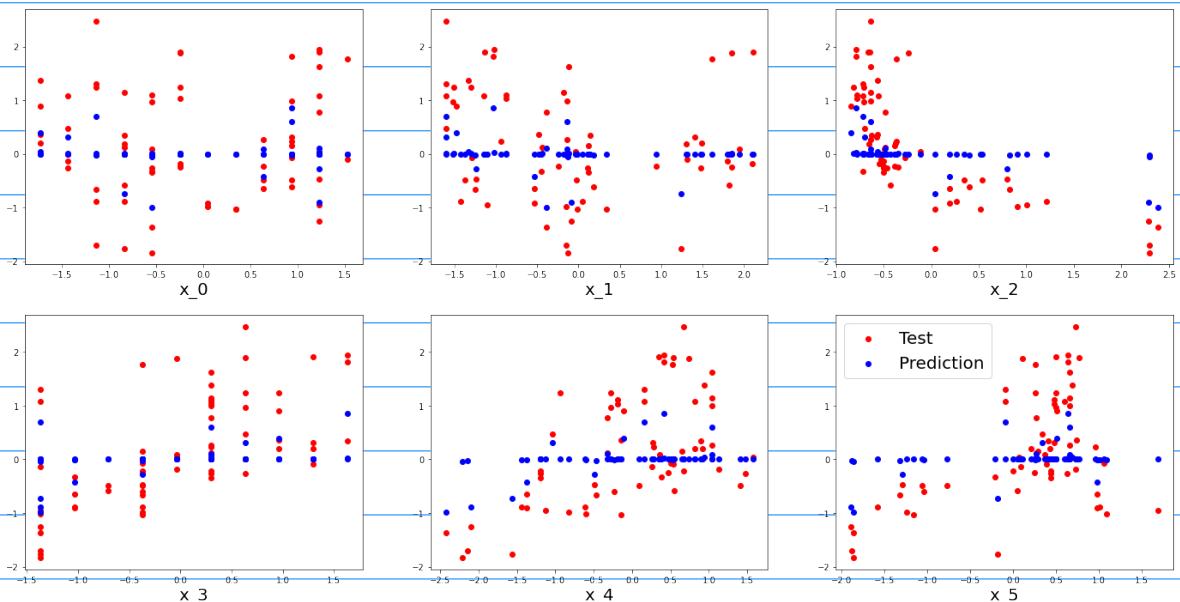
$\text{lengthscale } \gamma \quad \uparrow \text{reg. parameter.}$
 of kernel

Instead of working with γ I prefer to define $\sigma = \sqrt{2\gamma}$ so that the kernel looks like $\exp\left(-\frac{\|x - x'\|^2}{2\sigma^2}\right)$, a Gaussian.

The choice of σ & λ has a profound impact!

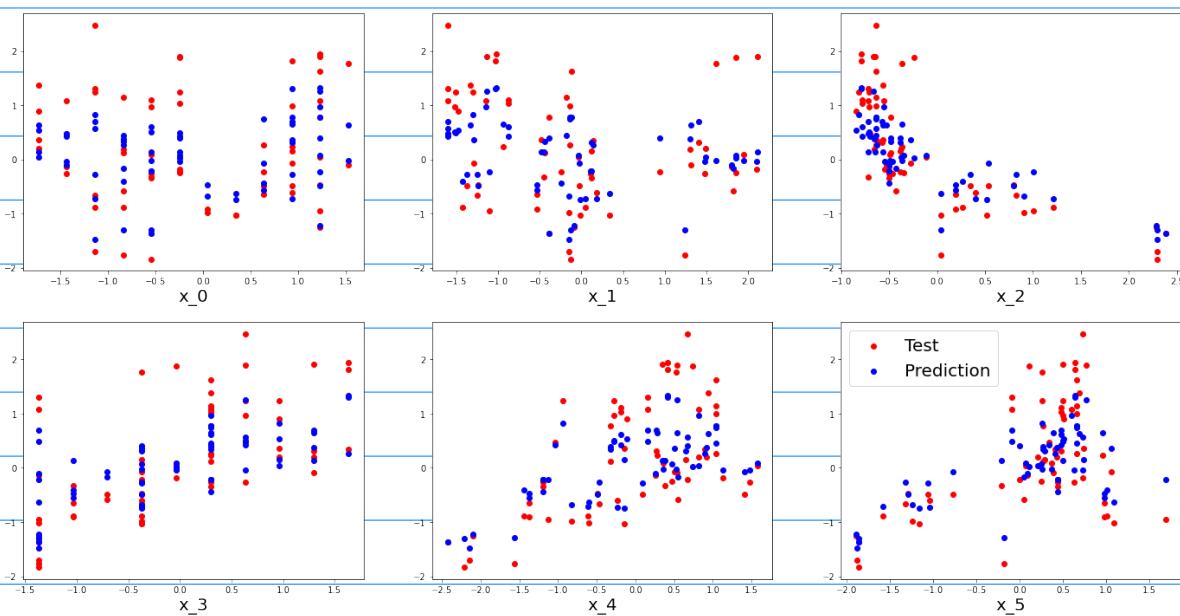
$\sigma = 0.1$

Bad results



$\sigma = 0.5$

Good!



So, we need a strategy to pick σ & λ . We will use 5-fold CV. But since we have two parameters, σ & λ there are no automatic functions to give us the optimal choice. This is where your job as a data scientist is to make an informed decision.

