

Multiple Perspectives on Linear Regression

Jayanth Koushik
Carnegie Mellon University
jkoushik@cs.cmu.edu

1 Introduction

Linear regression is one of the most common regression methods. It is simple to understand, and easy to implement making it a suitable baseline to evaluate. Further, by using basis functions, it can actually be used to model very complex relations. This is a tutorial on simple linear regression and ridge regression. However, we will go beyond simply deriving results. We will consider multiple perspectives to gain a deeper understanding of linear regression, and see how these different perspectives are related. We also consider the relation between ordinary least squares and ridge regression.

2 Problem setting

We will consider a very general problem setting independent of any particular method. We have a set of n data points $\mathcal{D} \equiv \{(\mathbf{x}_i, y_i)\}_{i=1}^n$ where $\mathbf{x}_i \in \mathbb{R}^d$, and $y_i \in \mathbb{R}$. This set is referred to as the training set. Let $\mathbb{X} \equiv [\mathbf{x}_1 \ \dots \ \mathbf{x}_n]^T$, and $\mathbf{y} \equiv [y_1 \ \dots \ y_n]^T$. The goal is to learn a function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ to make predictions for points from \mathbb{R}^d not in \mathcal{D} . Given a data point (\mathbf{x}, y) , the error in predicting $f(\mathbf{x})$ is the squared difference $(y - f(\mathbf{x}))^2$.

3 Best predictor

We can assume that the training set points are iid samples of a random variable pair (\mathbf{X}, Y) distributed according to some unknown joint probability distribution. We define the risk of f as the expected error $\mathbb{E}[(Y - f(\mathbf{X}))^2]$. With this, it is natural to seek the function with minimum risk. Define the regression function $m(\mathbf{x}) \equiv \mathbb{E}[Y|\mathbf{X} = \mathbf{x}]$ and write the risk as

$$\begin{aligned} \mathbb{E}[(Y - f(\mathbf{X}))^2] &= \mathbb{E}[\mathbb{E}[(Y - f(\mathbf{x}))^2|\mathbf{X}]] \\ &= \mathbb{E}[\mathbb{E}[(Y - m(\mathbf{X}) + m(\mathbf{X}) - f(\mathbf{X}))^2|\mathbf{X}]] \\ &= \mathbb{E}[\mathbb{E}[(Y - m(\mathbf{X}))^2|\mathbf{X}]] + \mathbb{E}[\mathbb{E}[(m(\mathbf{X}) - f(\mathbf{X}))^2|\mathbf{X}]] + \\ &\quad 2\mathbb{E}[\mathbb{E}[(Y - m(\mathbf{X}))(m(\mathbf{X}) - f(\mathbf{X}))|\mathbf{X}]] \end{aligned}$$

$$\begin{aligned}
&= \mathbb{E}[V[Y|\mathbf{X}]] + \mathbb{E}[(m(\mathbf{X}) - f(\mathbf{X}))^2] + \\
&\quad 2\mathbb{E}[(m(\mathbf{X}) - f(\mathbf{X}))(\mathbb{E}[Y|\mathbf{X}] - m(\mathbf{X}))] \\
&= \mathbb{E}[V[Y|\mathbf{X}]] + \mathbb{E}[(m(\mathbf{X}) - f(\mathbf{X}))^2] \tag{1}
\end{aligned}$$

Note that the first term is independent of f , and the second term is minimum (and equal to 0) if $f(\mathbf{x}) = m(\mathbf{x}) \forall \mathbf{x}$. So the regression function is the ideal predictor in some sense; but since it depends on the unknown data distribution, we cannot actually compute it. This result is only theoretical.

4 Ordinary least squares

Ordinary least squares (OLS) is a very simple approach to regression. We model f as a linear function of \mathbf{x} i.e. $f(\mathbf{x}) = \boldsymbol{\theta}^T \mathbf{x}$ for some $\boldsymbol{\theta} \in \mathbb{R}^d$. We do not add a separate bias variable, since we can always just add a constant 1 to each \mathbf{x} . This linear model reduces the problem from finding f to finding $\boldsymbol{\theta}$. We can consider this from several perspectives, which all lead to the same solution.

4.1 Minimization of squared error sum

Under the linear model, the error corresponding to a point (\mathbf{x}, y) is $(y - \boldsymbol{\theta}^T \mathbf{x})^2$. So we can seek to find the $\boldsymbol{\theta}$ that minimizes the sum of errors over the training set.

$$\boldsymbol{\theta}_* = \arg \min_{\boldsymbol{\theta}} \sum_{i=1}^n (y_i - \boldsymbol{\theta}^T \mathbf{x}_i)^2 = \arg \min_{\boldsymbol{\theta}} (\mathbf{y} - \mathbb{X}\boldsymbol{\theta})^T (\mathbf{y} - \mathbb{X}\boldsymbol{\theta}) \tag{2}$$

The derivative of this error with respect to $\boldsymbol{\theta}$ is $-2\mathbb{X}^T \mathbf{y} + 2\mathbb{X}^T \mathbb{X}\boldsymbol{\theta}$, so

$$2\mathbb{X}^T \mathbb{X}\boldsymbol{\theta}_* = 2\mathbb{X}^T \mathbf{y} \implies \boldsymbol{\theta}_* = (\mathbb{X}^T \mathbb{X})^{-1} \mathbb{X}^T \mathbf{y} \tag{3}$$

provided $\mathbb{X}^T \mathbb{X}$ is not singular. In case of singularity, we say that the least squares solution does not exist. The second derivative of the error is $2\mathbb{X}^T \mathbb{X}$. As shown in Lemma 1, $\mathbf{A}^T \mathbf{A}$ is positive semi-definite for any matrix \mathbf{A} ; so the solution given by Equation 3 is indeed a minima. With this expression for $\boldsymbol{\theta}$, the prediction for any \mathbf{x} is given by

$$f(\mathbf{x}) = \boldsymbol{\theta}^T \mathbf{x} = \mathbf{x}^T \boldsymbol{\theta} = \mathbf{x}^T (\mathbb{X}^T \mathbb{X})^{-1} \mathbb{X}^T \mathbf{y} = \mathbf{h}^T \mathbf{y} = \sum_{i=1}^n h_i y_i \tag{4}$$

where $\mathbf{h}^T \equiv \mathbf{x}^T (\mathbb{X}^T \mathbb{X})^{-1} \mathbb{X}^T$. We see that the prediction is just a linear combination of the training set labels. Such predictors are called linear smoothers.

4.2 Generalization of linear fit

Suppose we want to fit a line through the data points. So we have $y_i = \boldsymbol{\theta}^T \mathbf{x}_i$ for $i \in \{1, \dots, n\}$. This system of equations can be expressed concisely as

$$\mathbb{X}\boldsymbol{\theta} = \mathbf{y} \tag{5}$$

From this equation, it is clear that \mathbb{X} must be square for a solution to exist; this corresponds to the number of equations being equal to the number of parameters we are trying to estimate. Furthermore, \mathbb{X} must be invertible in which case the solution of Equation (5) is given by

$$\boldsymbol{\theta} = \mathbb{X}^{-1}\mathbf{y} \quad (6)$$

Now we can see the OLS solution given by Equation (3) as a generalization of the above solution with $(\mathbb{X}^T\mathbb{X})^{-1}\mathbb{X}^T$ being the Moore-Penrose pseudoinverse of \mathbb{X} .

4.3 Maximum likelihood estimation

Let us now look at OLS from a probabilistic viewpoint. Specifically, we assume $y = \boldsymbol{\theta}^T \mathbf{x} + \epsilon$, where ϵ is independent Gaussian noise with 0 mean and variance σ^2 . So given \mathbf{x} , y is a random variable distributed as $\mathcal{N}(\boldsymbol{\theta}^T \mathbf{x}, \sigma^2)$ where $\mathcal{N}(\mu, \sigma^2)$ denotes a normal distribution with mean μ and variance σ^2 . This is a parametric model for y characterized by $\boldsymbol{\theta}$, and we can find the likelihood of $\boldsymbol{\theta}$ given the training data as

$$\begin{aligned} L(\boldsymbol{\theta}|\mathbb{X}, \mathbf{y}) &\equiv p(\mathbf{y}|\mathbb{X}, \boldsymbol{\theta}) \propto \prod_{i=1}^n \exp\left(-\frac{1}{2\sigma^2}(y_i - \boldsymbol{\theta}^T \mathbf{x}_i)^2\right) \\ &= \exp\left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \boldsymbol{\theta}^T \mathbf{x}_i)^2\right) \\ &= \exp\left(-\frac{1}{2\sigma^2} (\mathbf{y} - \mathbb{X}\boldsymbol{\theta})^T (\mathbf{y} - \mathbb{X}\boldsymbol{\theta})\right) \end{aligned} \quad (7)$$

The maximum likelihood solution for $\boldsymbol{\theta}$ is the value that maximizes L , or equivalently minimizes $l \equiv -2\sigma^2 \log L$. But this is the same as Equation (2) which shows that the maximum likelihood solution is the same as that obtained by minimizing the sum of squared errors (Equation (3)).

4.4 Consistent estimation of best linear predictor

In the probabilistic view of Section 3, we made no assumptions about the relation between \mathbf{X} and Y . Particularly, we did not assume it to be linear. However, we can still talk about the *best* linear predictor i.e. the linear predictor with the least risk.

$$\boldsymbol{\theta}_* = \arg \min_{\boldsymbol{\theta}} \mathbb{E}[(Y - \boldsymbol{\theta}^T \mathbf{X})^2] = \arg \min_{\boldsymbol{\theta}} \mathbb{E}[(Y - \mathbf{X}^T \boldsymbol{\theta})^2] \quad (8)$$

The derivative of the risk with respect to $\boldsymbol{\theta}$ is $-2\mathbb{E}[\mathbf{X}(Y - \mathbf{X}^T \boldsymbol{\theta})]$. So $\boldsymbol{\theta}_*$ is given by

$$\mathbb{E}[\mathbf{X}\mathbf{X}^T]\boldsymbol{\theta}_* = \mathbb{E}[\mathbf{X}Y] \implies \boldsymbol{\theta}_* = (\mathbb{E}[\mathbf{X}\mathbf{X}^T])^{-1}\mathbb{E}[\mathbf{X}Y] \quad (9)$$

provided the inverse exists. The second derivative is $2E[\mathbf{X}\mathbf{X}^T]$ which is positive semi-definite since, from Lemma 1, $\mathbf{b}\mathbf{b}^T$ is positive semi-definite for any vector \mathbf{b} . But like before, this value depends on the unknown data distribution, and cannot be computed.

However, we can try to estimate the expression in Equation (9). Let $S \equiv \mathbb{X}^T\mathbb{X}/n$. Let A_i^j denote the element in row i , column j of matrix A . For any $i, j \in \{1, \dots, d\}$, we have by the law of large numbers:

$$S_i^j = \frac{1}{n} \sum_{k=1}^n \mathbf{x}_k^i \mathbf{x}_k^j \xrightarrow{P} E[\mathbf{X}^i \mathbf{X}^j] \implies S \xrightarrow{P} E[\mathbf{X}\mathbf{X}^T] \quad (10)$$

So by the continuous mapping theorem,

$$S^{-1} = n(\mathbb{X}^T\mathbb{X})^{-1} \xrightarrow{P} (E[\mathbf{X}\mathbf{X}^T])^{-1} \quad (11)$$

Similarly define $\mathbf{q} \equiv \mathbb{X}^T\mathbf{y}/n$, and note that for any $i \in \{1, \dots, d\}$, we have

$$\mathbf{q}^i = \frac{1}{n} \sum_{k=1}^n \mathbf{x}_k^i \mathbf{y}_k \xrightarrow{P} E[\mathbf{X}^i Y] \implies \mathbf{q} \xrightarrow{P} E[\mathbf{X}Y] \quad (12)$$

Finally we can combine the results of Equations (11) and (12) with the continuous mapping theorem to get

$$S^{-1}\mathbf{q} = (\mathbb{X}^T\mathbb{X})^{-1}\mathbb{X}^T\mathbf{y} \xrightarrow{P} (E[\mathbf{X}\mathbf{X}^T])^{-1}E[\mathbf{X}Y] = \boldsymbol{\theta}_* \quad (13)$$

$S^{-1}\mathbf{q}$ is the same value given by Equation (3) which shows that the OLS solution is a consistent estimator of the best linear predictor.

5 Ridge regression

One problem with the OLS solution (Equation (3)) is that it depends on $\mathbb{X}^T\mathbb{X}$ being invertible. The original motivation for ridge regression was to directly fix this problem. As shown in Lemma 2, for any positive semi-definite matrix A , and any positive number λ , the matrix $A + \lambda I$ is invertible. So the ridge regression solution modifies the OLS solution as

$$\boldsymbol{\theta}_* = (\mathbb{X}^T\mathbb{X} + \lambda I)^{-1}\mathbb{X}^T\mathbf{y} \quad (14)$$

for some positive λ . This solution always exists, independent of \mathbb{X} and \mathbf{y} . The constant λ is a hyperparameter that is usually picked by cross-validation. However, in the following subsections, we will look at ridge regression from two different perspectives to gain some intuition about this hyperparameter, and the overall solution in general. Note that setting λ to 0 reduces this solution to the OLS solution. Given a new point \mathbf{x} , the prediction is given by

$$\mathbf{x}^T\boldsymbol{\theta}_* = \mathbf{x}^T(\mathbb{X}^T\mathbb{X} + \lambda I)^{-1}\mathbb{X}^T\mathbf{y} \quad (15)$$

Note that this too is a linear smoother.

5.1 Tikhonov regularization

When the OLS solution does not exist, the linear regression problem is said to be ill-posed. This can happen if $n < d$, or if the columns of \mathbb{X} are highly correlated. A common strategy for solving ill-posed problems is Tikhonov regularization. It seeks to make the problem well posed by adding further constraints on the model. In the case of linear regression, this is achieved by adding another term to the OLS objective (Equation (2)) which penalizes the norm of $\boldsymbol{\theta}$:

$$\boldsymbol{\theta}_* = \arg \min_{\boldsymbol{\theta}} \|\mathbf{y} - \mathbb{X}\boldsymbol{\theta}\|^2 + \lambda \|\boldsymbol{\theta}\|^2 = \arg \min_{\boldsymbol{\theta}} (\mathbf{y} - \mathbb{X}\boldsymbol{\theta})^T (\mathbf{y} - \mathbb{X}\boldsymbol{\theta}) + \lambda \boldsymbol{\theta}^T \boldsymbol{\theta} \quad (16)$$

Here λ is a hyperparameter that controls the amount regularization. The derivative of the above expression with respect to $\boldsymbol{\theta}$ is $-2\mathbb{X}^T \mathbf{y} + 2\mathbb{X}^T \mathbb{X} \boldsymbol{\theta} + 2\lambda \boldsymbol{\theta}$. So,

$$(\mathbb{X}^T \mathbb{X} + \lambda \mathbf{I}) \boldsymbol{\theta}_* = \mathbb{X}^T \mathbf{y} \implies \boldsymbol{\theta}_* = (\mathbb{X}^T \mathbb{X} + \lambda \mathbf{I})^{-1} \mathbb{X}^T \mathbf{y} \quad (17)$$

While this view motivates ridge regression as a way to make the problem well posed, the role of λ is not very clear. This will be made more explicit in the next section.

5.2 Maximum a posteriori estimation

In Section 4.3, we assumed Gaussian noise, and found $\boldsymbol{\theta}$ by maximizing the likelihood. Let us now take this further with a Bayesian approach. We will put a prior distribution on $\boldsymbol{\theta}$; a normal distribution with mean $\mathbf{0}$ and variance Σ . With the Bayesian approach, we are interested in the posterior distribution of $\boldsymbol{\theta}$ i.e. the distribution given the data. We will assume that the covariates \mathbf{x}_i are fixed, so by Bayes' theorem

$$p(\boldsymbol{\theta} | \mathbb{X}, \mathbf{y}) = \frac{p(\mathbf{y} | \mathbb{X}, \boldsymbol{\theta}) p(\boldsymbol{\theta})}{p(\mathbf{y} | \mathbb{X})} \propto L(\boldsymbol{\theta} | \mathbb{X}, \mathbf{y}) p(\boldsymbol{\theta}) \quad (18)$$

Here L is the same likelihood calculated in Equation (7). The denominator $p(\mathbf{y} | \mathbb{X})$ can be ignored since it does not depend on $\boldsymbol{\theta}$ – it is just a normalizing constant to ensure that $p(\boldsymbol{\theta} | \mathbb{X}, \mathbf{y})$ is a valid distribution. Multiplying the likelihood and the prior, we see that the posterior distribution is proportional to

$$\begin{aligned} & \exp\left(-\frac{1}{2\sigma^2} (\mathbf{y} - \mathbb{X}\boldsymbol{\theta})^T (\mathbf{y} - \mathbb{X}\boldsymbol{\theta})\right) \exp\left(-\frac{1}{2} \boldsymbol{\theta}^T \Sigma^{-1} \boldsymbol{\theta}\right) \\ &= \exp\left(-\frac{1}{2} \left(\frac{1}{\sigma^2} \mathbf{y}^T \mathbf{y} - \frac{2}{\sigma^2} \mathbf{y}^T \mathbb{X} \boldsymbol{\theta} + \frac{1}{\sigma^2} \boldsymbol{\theta}^T \mathbb{X}^T \mathbb{X} \boldsymbol{\theta} + \boldsymbol{\theta}^T \Sigma^{-1} \boldsymbol{\theta}\right)\right) \\ &\propto \exp\left(-\frac{1}{2} \left(\boldsymbol{\theta}^T (\Sigma^{-1} + \frac{1}{\sigma^2} \mathbb{X}^T \mathbb{X}) \boldsymbol{\theta} - \frac{2}{\sigma^2} \boldsymbol{\theta}^T \mathbb{X}^T \mathbf{y}\right)\right) \end{aligned} \quad (19)$$

Let $\mathbf{A} \equiv \Sigma^{-1} + \mathbb{X}^T \mathbb{X} / \sigma^2$ and $\mathbf{b} \equiv \mathbb{X}^T \mathbf{y} / \sigma^2$. We can use a very handy completion of squares trick to simplify the above expression as

$$\exp\left(-\frac{1}{2} (\boldsymbol{\theta}^T \mathbf{A} \boldsymbol{\theta} - 2\boldsymbol{\theta}^T \mathbf{b})\right)$$

$$\begin{aligned}
&= \exp\left(-\frac{1}{2}(\boldsymbol{\theta}^T \mathbf{A} \boldsymbol{\theta} - \boldsymbol{\theta}^T \mathbf{A} \mathbf{A}^{-1} \mathbf{b} - \mathbf{b}^T \mathbf{A}^{-1} \mathbf{A} \boldsymbol{\theta})\right) \\
&\propto \exp\left(-\frac{1}{2}(\boldsymbol{\theta}^T \mathbf{A} \boldsymbol{\theta} - \boldsymbol{\theta}^T \mathbf{A} \mathbf{A}^{-1} \mathbf{b} - \mathbf{b}^T \mathbf{A}^{-1} \mathbf{A} \boldsymbol{\theta} + \mathbf{b}^T \mathbf{A}^{-1} \mathbf{A} \mathbf{A}^{-1} \mathbf{b})\right) \\
&= \exp\left(-\frac{1}{2}(\boldsymbol{\theta} - \mathbf{A}^{-1} \mathbf{b})^T \mathbf{A} (\boldsymbol{\theta} - \mathbf{A}^{-1} \mathbf{b})\right) \tag{20}
\end{aligned}$$

Note that this expression has the form of a normal distribution with mean $\mathbf{A}^{-1} \mathbf{b}$ and variance \mathbf{A}^{-1} . This is not a coincidence. The prior distribution of $\boldsymbol{\theta}$ was chosen specifically for this property, and such distributions are called conjugate priors. The Bayesian approach does not give us a single estimate of $\boldsymbol{\theta}$; rather, it gives us a whole distribution. This can be very useful, and in fact, Bayesian regression directly uses this distribution to make predictions. However, for now we are just interested in a point estimate. This is typically obtained by taking the mean or mode of the posterior distribution. In the case of a Gaussian the two are equivalent, so our estimate is given by

$$\boldsymbol{\theta}_* = \mathbf{A}^{-1} \mathbf{b} = \left(\boldsymbol{\Sigma}^{-1} + \frac{1}{\sigma^2} \mathbb{X}^T \mathbb{X}\right)^{-1} \frac{1}{\sigma^2} \mathbb{X}^T \mathbf{y} = (\mathbb{X}^T \mathbb{X} + \sigma^2 \boldsymbol{\Sigma}^{-1})^{-1} \mathbb{X}^T \mathbf{y} \tag{21}$$

This expression bears a lot of resemblance to the ridge regression solution (Equation (14)); in fact if we assume a diagonal variance for $\boldsymbol{\theta}$ i.e. $\boldsymbol{\Sigma} = \nu^2 \mathbf{I}$, then the expression becomes equivalent to Equation (14) with $\lambda = \sigma^2/\nu^2$. This captures the role of the regularization hyperparameter as the amount of noise relative to the uncertainty in $\boldsymbol{\theta}$. ν controls the breadth of the prior distribution, and signifies our uncertainty in $\boldsymbol{\theta}$. A large value indicates that we do not have strong prior beliefs, and in the limit $\nu \rightarrow \infty$ produces what is known as an “uninformative” prior. It causes λ to shrink, and the solution goes closer to the OLS solution. Conversely, a small value of ν indicates a strong prior belief that $\boldsymbol{\theta}$ is close to $\mathbf{0}$ – it increases λ and makes $\boldsymbol{\theta}$ remain close to $\mathbf{0}$.

6 Kernel ridge regression

Linear regression tries to fit a very simple model to the data, and so will not perform very well if the data deviates strongly from a line. However, it is possible to implicitly fit more complex models by making use of basis functions. Instead of finding a line as a function of \mathbf{x} , we can first transform \mathbf{x} with a function $\boldsymbol{\phi} : \mathbb{R}^d \rightarrow \mathbb{R}^m$, and find a line as a function of $\boldsymbol{\phi}(\mathbf{x})$. Note that m can be much larger than d or even infinite. This simple change allows for much more complex functions to be fit through the data. As a simple example, consider unidimensional covariates x , and let $\phi(x) \equiv x^2$. Then a line in ϕ -space is a quadratic curve in x -space. The regression solution remains the same. Let $\boldsymbol{\Phi} \equiv [\boldsymbol{\phi}(\mathbf{x}_1) \ \dots \ \boldsymbol{\phi}(\mathbf{x}_n)]^T$. Then the ridge regression solution in $\boldsymbol{\phi}$ -space is simply

$$\boldsymbol{\theta}_* = (\boldsymbol{\Phi}^T \boldsymbol{\Phi} + \lambda \mathbf{I})^{-1} \boldsymbol{\Phi}^T \mathbf{y} \tag{22}$$

However there is a problem with this solution. It requires explicitly transforming points with $\boldsymbol{\phi}$ and working with $\boldsymbol{\Phi}$. This is computationally expensive if m is

very large, and impossible if it is infinite. So we will make use of the *kernel trick* which will let us make *implicit* transformations. Let $k : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ be any symmetric positive semi-definite function. Then by Mercer's theorem,

$$k(\mathbf{x}, \mathbf{x}') = \langle \Phi(\mathbf{x}), \Phi(\mathbf{x}') \rangle \quad (23)$$

for some $\Phi : \mathbb{R}^d \rightarrow \mathbb{R}^m$. Mercer's theorem allows us to create an implicit mapping by using a kernel function. So, in our regression we will assume that Φ is generated by some Mercer kernel k . A common choice is the Gaussian kernel $k(\mathbf{x}, \mathbf{x}') \equiv \exp(-\|\mathbf{x}-\mathbf{x}'\|^2/2\sigma^2)$ for some constant σ . Let \mathbb{K} be a matrix such that $\mathbb{K}_i^j \equiv \Phi(\mathbf{x}_i)^T \Phi(\mathbf{x}_j) = \langle \Phi(\mathbf{x}_i), \Phi(\mathbf{x}_j) \rangle = k(\mathbf{x}_i, \mathbf{x}_j)$. Observe that $\mathbb{K} = \Phi\Phi^T$. Next we will make use of the following identity from Lemma 3. If P and R are positive definite matrices, then

$$(\mathbf{P}^{-1} + \mathbf{B}^T \mathbf{R}^{-1} \mathbf{B})^{-1} \mathbf{B}^T \mathbf{R}^{-1} = \mathbf{P} \mathbf{B}^T (\mathbf{B} \mathbf{P} \mathbf{B}^T + \mathbf{R})^{-1} \quad (24)$$

where B is such that the inverses are well defined. Taking \mathbf{P}^{-1} equal to $\lambda \mathbf{I}$, R equal to I, and B equal to Φ , we get

$$(\lambda \mathbf{I} + \Phi^T \Phi)^{-1} \Phi^T = \frac{1}{\lambda} \Phi^T \left(\frac{1}{\lambda} \Phi \Phi^T + \mathbf{I} \right)^{-1} = \Phi^T (\Phi \Phi^T + \lambda \mathbf{I})^{-1} \quad (25)$$

For any new point \mathbf{x} , let $\mathbf{k}_x \equiv [k(\mathbf{x}, \mathbf{x}_1) \dots k(\mathbf{x}, \mathbf{x}_n)]^T = \Phi \Phi(\mathbf{x})$. From Equations (22) and (25), we have that the prediction at \mathbf{x} is given by

$$\Phi(\mathbf{x})^T \theta_* = \Phi(\mathbf{x})^T \Phi^T (\Phi \Phi^T + \lambda \mathbf{I})^{-1} \mathbf{y} = \mathbf{k}_x^T (\mathbb{K} + \lambda \mathbf{I})^{-1} \mathbf{y} \quad (26)$$

\mathbb{K} and \mathbf{k}_x both only make use of the kernel function k ; so the above expression allows us to use kernel ridge regression without explicitly mapping points using Φ . Observe that kernel ridge regression is also a form of linear smoothing.

7 Conclusion

We considered different approaches to simple linear regression and ridge regression. These different approaches, while leading to the same solution, provide different perspectives on these techniques, and show how they are related.

A Appendix

Lemma 1. *Let A be any matrix, and \mathbf{b} be any vector. $A^T A$ and $\mathbf{b}\mathbf{b}^T$ are positive semi-definite.*

Proof. For any vector \mathbf{x} , we have

$$\mathbf{x}^T A^T A \mathbf{x} = (A\mathbf{x})^T A\mathbf{x} = \|A\mathbf{x}\|^2 \geq 0$$

and

$$\mathbf{x}^T \mathbf{b}\mathbf{b}^T \mathbf{x} = (\mathbf{b}^T \mathbf{x})^T \mathbf{b}^T \mathbf{x} = \|\mathbf{b}^T \mathbf{x}\|^2 \geq 0$$

So $A^T A$ and $\mathbf{b}\mathbf{b}^T$ are positive semi-definite. \square

Lemma 2. *Let A be any positive semi-definite matrix, and ϵ be any positive scalar. $A + \epsilon I$ is invertible.*

Proof. Let \mathbf{v} be an eigenvector of A with corresponding eigenvalue λ . Since A is positive semi-definite, $\lambda \geq 0$. We have

$$(A + \epsilon I)\mathbf{v} = A\mathbf{v} + \epsilon\mathbf{v} = (\lambda + \epsilon)\mathbf{v}$$

So \mathbf{v} is also an eigenvector of $A + \epsilon I$ with eigenvalue $\lambda + \epsilon$. But $\lambda + \epsilon > 0$, and since the choice of \mathbf{v} was arbitrary, this implies that all eigenvalues of $A + \epsilon I$ are positive. Hence it is invertible. \square

Lemma 3. *Let P and R be positive definite matrices. Then*

$$(P^{-1} + B^T R^{-1} B)^{-1} B^T R^{-1} = P B^T (B P B^T + R)^{-1}$$

for any B such that the inverses are well defined.

Proof. We have

$$\begin{aligned} & (P^{-1} + B^T R^{-1} B)^{-1} B^T R^{-1} (B P B^T + R) \\ &= (P^{-1} + B^T R^{-1} B)^{-1} (B^T R^{-1} B P B^T + B^T R) \\ &= (P^{-1} + B^T R^{-1} B)^{-1} (B^T R^{-1} B + P^{-1}) P B^T \\ &= P B^T \end{aligned}$$

Now multiplying both sides by $(B P B^T + R)^{-1}$ gives the desired result. \square