

Linear regression, active learning

We arrived at the logistic regression model when trying to explicitly model the uncertainty about the labels in a linear classifier. The same general modeling approach permits us to use linear predictions in various other contexts as well. The simplest of them is regression where the goal is to predict a continuous response $y_t \in \mathcal{R}$ to each example vector. Here too focusing on linear predictions won't be inherently limiting as linear predictions can be easily extended (next lecture).

So, how should we model continuous responses? The linear function of the input already produces a “mean prediction” or $\underline{\theta}^T \underline{x} + \theta_0$. By treating this as a mean prediction more formally, we are stating that the expected value of the response variable, conditioned on \underline{x} , is $\underline{\theta}^T \underline{x} + \theta_0$. More succinctly, we say that $E\{y|\underline{x}\} = \underline{\theta}^T \underline{x} + \theta_0$. It remains to associate a distribution over the responses around such mean prediction. The simplest symmetric distribution is the normal (Gaussian) distribution. In other words, we say that the responses y follow the normal pdf

$$N(y; \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(y - \mu)^2\right) \quad (1)$$

where $\mu = \underline{\theta}^T \underline{x} + \theta_0$. Our response model is therefore defined as

$$P(y|\underline{x}, \underline{\theta}, \theta_0) = N(y; \underline{\theta}^T \underline{x} + \theta_0, \sigma^2) \quad (2)$$

So, when the input is 1-dimensional, we predict a mean response that is a line in (x, y) space, and assume that noise in y is normally distributed with zero mean and variance σ^2 . Note that the noise variance σ^2 in the model does not depend on the input x . Moreover, we only model variation in the y -direction while expecting to know x with perfect precision. Taking into account the effect of potential noise in x on the responses y would tie parameters $\underline{\theta}$ and θ_0 to the noise variance σ^2 , potentially in an input dependent manner. The specifics of this coupling depend on the form of noise added to x . We will discuss this in a bit more detail later on.

We can also write the linear regression model in another way to explicate how exactly the additive noise appears in the responses:

$$y = \underline{\theta}^T \underline{x} + \theta_0 + \epsilon \quad (3)$$

where $\epsilon \sim N(0, \sigma^2)$ (meaning that noise ϵ is distributed normally with mean zero and variance σ^2). Clearly for this model $E\{y|\underline{x}\} = \underline{\theta}^T \underline{x} + \theta_0$ since ϵ has zero mean. Moreover, adding Gaussian noise to a deterministic prediction $\underline{\theta}^T \underline{x} + \theta_0$ makes y normally distributed

with mean $\underline{\theta}^T \underline{x} + \theta_0$ and variance σ^2 , as before. So, in particular, for the training inputs $\underline{x}_1, \dots, \underline{x}_n$ and outputs y_1, \dots, y_n , the model relating them is

$$y_t = \underline{\theta}^T \underline{x}_t + \theta_0 + \epsilon_t, \quad t = 1, \dots, n \quad (4)$$

where $e_t \sim N(0, \sigma^2)$ and e_i is independent of e_j for any $i \neq j$.

Regardless of how we choose to write the model (both forms are useful) we can find the parameter estimates by maximizing the conditional likelihood. Similarly to the logistic regression case, the conditional likelihood is written as

$$L(\underline{\theta}, \theta_0, \sigma^2) = \prod_{t=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(y_t - \underline{\theta}^T \underline{x}_t - \theta_0)^2\right) \quad (5)$$

Note that σ^2 is also a parameter we have to estimate. It accounts for errors not captured by the linear model. In terms of the log-likelihood, we try to maximize

$$l(\underline{\theta}, \theta_0, \sigma^2) = \sum_{t=1}^n \log \left[\frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(y_t - \underline{\theta}^T \underline{x}_t - \theta_0)^2\right) \right] \quad (6)$$

$$= \sum_{t=1}^n \left[-\frac{1}{2} \log(2\pi) - \frac{1}{2} \log \sigma^2 - \frac{1}{2\sigma^2} (y_t - \underline{\theta}^T \underline{x}_t - \theta_0)^2 \right] \quad (7)$$

$$= \text{const.} - \frac{n}{2} \log \sigma^2 - \frac{1}{2\sigma^2} \sum_{t=1}^n (y_t - \underline{\theta}^T \underline{x}_t - \theta_0)^2 \quad (8)$$

where 'const.' absorbs terms that do not depend on the parameters. Now, the problem of estimating the parameters $\underline{\theta}$ and θ_0 is nicely decoupled from estimating σ^2 . In other words, we can find the *maximizing* $\hat{\underline{\theta}}$ and $\hat{\theta}_0$ by simply *minimizing* the mean squared error

$$\sum_{t=1}^n (y_t - \underline{\theta}^T \underline{x}_t - \theta_0)^2 \quad (9)$$

It is perhaps easiest to write the solution based on a bit of matrix calculation. Let \mathbf{X} be a matrix whose rows, indexed by training examples, are given by $[\underline{x}_t^T, 1]$ (\underline{x}_t turned into a row vector and 1 added at the end). In terms of this matrix, the minimization problem

becomes

$$\sum_{t=1}^n \left(y_t - \begin{bmatrix} \theta \\ \theta_0 \end{bmatrix}^T \begin{bmatrix} \underline{x}_t \\ 1 \end{bmatrix} \right)^2 = \sum_{t=1}^n \left(y_t - [\underline{x}_t^T, 1] \begin{bmatrix} \theta \\ \theta_0 \end{bmatrix} \right)^2 \quad (10)$$

$$= \left\| \begin{bmatrix} y_1 \\ \dots \\ y_n \end{bmatrix} - \begin{bmatrix} \underline{x}_1^T, 1 \\ \dots \\ \underline{x}_n^T, 1 \end{bmatrix} \begin{bmatrix} \theta \\ \theta_0 \end{bmatrix} \right\|^2 \quad (11)$$

$$= \left\| \mathbf{y} - \mathbf{X} \begin{bmatrix} \theta \\ \theta_0 \end{bmatrix} \right\|^2 \quad (12)$$

$$= \mathbf{y}^T \mathbf{y} - 2 \begin{bmatrix} \theta \\ \theta_0 \end{bmatrix}^T \mathbf{X}^T \mathbf{y} + \begin{bmatrix} \theta \\ \theta_0 \end{bmatrix}^T \mathbf{X}^T \mathbf{X} \begin{bmatrix} \theta \\ \theta_0 \end{bmatrix} \quad (13)$$

where $\mathbf{y} = [y_1, \dots, y_n]^T$ is a vector of training responses. Solving it yields

$$\begin{bmatrix} \hat{\theta} \\ \hat{\theta}_0 \end{bmatrix} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \quad (14)$$

Note that the optimal parameter values are linear functions of the observed responses \mathbf{y} . We will make use of this property later on. The dependence on the training inputs $\underline{x}_1, \dots, \underline{x}_n$ (or the matrix \mathbf{X}) is non-linear, however.

The noise variance can be subsequently set to account for the remaining prediction errors. Indeed, the the maximizing value of σ^2 is given by

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{t=1}^n (y_t - \hat{\theta}^T \underline{x}_t - \hat{\theta}_0)^2 \quad (15)$$

which is the average squared prediction error. Note that we cannot compute $\hat{\sigma}^2$ before knowing how well the linear model explains the responses.

Bias and variance of the parameter estimates

We can make use of the closed form parameter estimates in Eq.(14) to analyze how good these estimates are. For this purpose let's make the strong assumption that the actual relation between \underline{x} and y follows a linear model of the same type that we are estimating (we just don't know the correct parameter values $\underline{\theta}^*$, θ_0^* , and σ^{*2}). We can therefore describe the observed responses y_t as

$$y_t = \underline{\theta}^{*T} \underline{x}_t + \theta_0^* + \epsilon_t, \quad t = 1, \dots, n \quad (16)$$

where $\epsilon_t \sim N(0, \sigma^{*2})$. In a matrix form

$$\mathbf{y} = \mathbf{X} \begin{bmatrix} \theta^* \\ \theta_0^* \end{bmatrix} + \mathbf{e} \quad (17)$$

where $\mathbf{e} = [\epsilon_1, \dots, \epsilon_n]^T$, $E\{\mathbf{e}\} = 0$ and $E\{\mathbf{e}\mathbf{e}^T\} = \sigma^{*2}\mathbf{I}$. The noise vector \mathbf{e} is also independent of the inputs or \mathbf{X} . Plugging this form of responses into Eq.(14) we get

$$\begin{bmatrix} \hat{\theta} \\ \hat{\theta}_0 \end{bmatrix} = (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T(\mathbf{X} \begin{bmatrix} \theta^* \\ \theta_0^* \end{bmatrix} + \mathbf{e}) \quad (18)$$

$$= (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{X} \begin{bmatrix} \theta^* \\ \theta_0^* \end{bmatrix} + (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{e} \quad (19)$$

$$= \begin{bmatrix} \theta^* \\ \theta_0^* \end{bmatrix} + (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{e} \quad (20)$$

In other words, our parameter estimates can be decomposed into the sum of correct underlying parameters and estimates based on noise alone (i.e., based on \mathbf{e}). Thus, on average with fixed inputs

$$E\left\{\begin{bmatrix} \hat{\theta} \\ \hat{\theta}_0 \end{bmatrix} \middle| \mathbf{X}\right\} = \begin{bmatrix} \theta^* \\ \theta_0^* \end{bmatrix} + (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T E\{\mathbf{e}|\mathbf{X}\} = \begin{bmatrix} \theta^* \\ \theta_0^* \end{bmatrix} \quad (21)$$

Our parameter estimates are therefore *unbiased* or correct on average when averaging is over possible training sets we could generate. The averaging here is conditioned on the specific training inputs.

Using Eq.(20) and Eq.(21) we can also evaluate the conditional co-variance of the parameter estimates where the expectation is again over the noise in the outputs:

$$Cov\left\{\begin{bmatrix} \hat{\theta} \\ \hat{\theta}_0 \end{bmatrix} \middle| \mathbf{X}\right\} = E\left\{\left(\begin{bmatrix} \hat{\theta} \\ \hat{\theta}_0 \end{bmatrix} - \begin{bmatrix} \theta^* \\ \theta_0^* \end{bmatrix}\right)\left(\begin{bmatrix} \hat{\theta} \\ \hat{\theta}_0 \end{bmatrix} - \begin{bmatrix} \theta^* \\ \theta_0^* \end{bmatrix}\right)^T \middle| \mathbf{X}\right\} \quad (22)$$

$$= E\left\{((\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{e})((\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{e})^T \middle| \mathbf{X}\right\} \quad (23)$$

$$= E\left\{(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{e}\mathbf{e}^T\mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1} \middle| \mathbf{X}\right\} \quad (24)$$

$$= (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T E\{\mathbf{e}\mathbf{e}^T|\mathbf{X}\} \mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1} \quad (25)$$

$$= (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T (\sigma^{*2}\mathbf{I}) \mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1} \quad (26)$$

$$= \sigma^{*2}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1} \quad (27)$$

$$= \sigma^{*2}(\mathbf{X}^T\mathbf{X})^{-1} \quad (28)$$

So the way in which the parameters vary in response to noise in the outputs is a function of the inputs or \mathbf{X} . We will use this property in the next section to select inputs so as to improve the quality of the parameter estimates or to reduce the variance of predictions.

Based on the bias and variance calculations we can evaluate the mean squared error of the parameter estimates. To this end, we use the fact that the expectation of the squared norm of any vector valued random variable can be decomposed into a bias and variance components as follows:

$$E\left\{\|\mathbf{z} - \mathbf{z}^*\|^2\right\} = E\left\{\|\mathbf{z} - E\{\mathbf{z}\} + E\{\mathbf{z}\} - \mathbf{z}^*\|^2\right\} \quad (29)$$

$$= E\left\{\|\mathbf{z} - E\{\mathbf{z}\}\|^2 + 2(\mathbf{z} - E\{\mathbf{z}\})^T(E\{\mathbf{z}\} - \mathbf{z}^*) + \|E\{\mathbf{z}\} - \mathbf{z}^*\|^2\right\} \quad (30)$$

$$= E\left\{\|\mathbf{z} - E\{\mathbf{z}\}\|^2\right\} + 2E\left\{(\mathbf{z} - E\{\mathbf{z}\})^T\right\}(E\{\mathbf{z}\} - \mathbf{z}^*) + \|E\{\mathbf{z}\} - \mathbf{z}^*\|^2$$

$$= \underbrace{E\left\{\|\mathbf{z} - E\{\mathbf{z}\}\|^2\right\}}_{\text{variance}} + \underbrace{\|E\{\mathbf{z}\} - \mathbf{z}^*\|^2}_{\text{bias}^2} \quad (31)$$

where we have assumed that \mathbf{z}^* is fixed. Make sure you understand how this decomposition is derived. We will further elaborate the variance part to better use the result in our context:

$$E\left\{\|\mathbf{z} - E\{\mathbf{z}\}\|^2\right\} = E\left\{(\mathbf{z} - E\{\mathbf{z}\})^T(\mathbf{z} - E\{\mathbf{z}\})\right\} \quad (32)$$

$$= E\left\{Tr\left[(\mathbf{z} - E\{\mathbf{z}\})^T(\mathbf{z} - E\{\mathbf{z}\})\right]\right\} \quad (33)$$

$$= E\left\{Tr\left[(\mathbf{z} - E\{\mathbf{z}\})(\mathbf{z} - E\{\mathbf{z}\})^T\right]\right\} \quad (34)$$

$$= Tr\left[E\left\{(\mathbf{z} - E\{\mathbf{z}\})(\mathbf{z} - E\{\mathbf{z}\})^T\right\}\right] \quad (35)$$

$$= Tr[Cov\{\mathbf{z}\}] \quad (36)$$

where $Tr[\cdot]$ is the matrix *trace*, the sum of its diagonal components, and therefore a linear operation (exchangeable with the expectation). We have also used the fact that $Tr[\mathbf{a}^T\mathbf{b}] = Tr[\mathbf{a}\mathbf{b}^T]$ for any vectors \mathbf{a} and \mathbf{b} .

Now, adapting the result to our setting, we get

$$\begin{aligned}
 E \left\{ \left\| \begin{bmatrix} \hat{\theta} \\ \hat{\theta}_0 \end{bmatrix} - \begin{bmatrix} \theta^* \\ \theta_0^* \end{bmatrix} \right\|^2 \mid \mathbf{X} \right\} &= \overbrace{Tr \left[Cov \left\{ \begin{bmatrix} \hat{\theta} \\ \hat{\theta}_0 \end{bmatrix} \mid \mathbf{X} \right\} \right]}^{\text{variance}} + \overbrace{\left\| E \left\{ \begin{bmatrix} \hat{\theta} \\ \hat{\theta}_0 \end{bmatrix} \mid \mathbf{X} \right\} - \begin{bmatrix} \theta^* \\ \theta_0^* \end{bmatrix} \right\|^2}_{\text{bias}^2=0} \\
 &= \sigma^{*2} Tr \left[(\mathbf{X}^T \mathbf{X})^{-1} \right] \tag{37}
 \end{aligned}$$

Let's understand this result a bit further. How does it depend on n , the number of training examples? In other words, how quickly does the mean squared error decrease as the number of training examples increases, assuming the input examples \underline{x} are sampled independently from some underlying distribution $P(\underline{x})$? To answer this let's start by analyzing what happens to the matrix $\mathbf{X}^T \mathbf{X}$:

$$\mathbf{X}^T \mathbf{X} = \sum_{t=1}^n \begin{bmatrix} \underline{x}_t \\ 1 \end{bmatrix} [\underline{x}_t^T, 1] \tag{38}$$

$$= n \cdot \frac{1}{n} \sum_{t=1}^n \begin{bmatrix} \underline{x}_t \\ 1 \end{bmatrix} [\underline{x}_t^T, 1] \tag{39}$$

$$\approx n \cdot E_{\underline{x} \sim P} \left\{ \begin{bmatrix} \underline{x} \\ 1 \end{bmatrix} [\underline{x}^T, 1] \right\} = n \cdot \mathbf{C} \tag{40}$$

where for large n the average will be close to the corresponding expected value. For large n the mean squared error of the parameter estimates will therefore be close to

$$\frac{\sigma^{*2}}{n} \cdot Tr[\mathbf{C}^{-1}] \tag{41}$$

The variance of simply averaging the (noise in the) outputs would behave as σ^{*2}/n . Since we are estimating $d+1$ parameters where d is the input dimension, this dependence would have to be in $Tr[\mathbf{C}^{-1}]$. Indeed it is. This term, a trace of a $(d+1) \times (d+1)$ matrix \mathbf{C}^{-1} , is directly proportional to $d+1$.

Penalized log-likelihood and Ridge regression

When the number of training examples is small, i.e., not too much larger than the number of parameters (dimension of the inputs), it is often beneficial to *regularize* the parameter estimates. We will derive the form of regularization here by assigning a prior distribution over the parameters $P(\underline{\theta}, \theta_0)$. The purpose of the prior is to prefer small parameter values

(predict values close to zero) in the absence of data. Specifically, we will look at simple zero mean Gaussian distributions

$$P(\underline{\theta}, \theta_0; \sigma'^2) = N\left(\begin{bmatrix} \underline{\theta} \\ \theta_0 \end{bmatrix}; \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \sigma'^2 \mathbf{I}\right) = N(\theta_0; 0, \sigma'^2) \prod_{j=1}^d N(\underline{\theta}_j; 0, \sigma'^2) \quad (42)$$

where the variance parameter σ'^2 in the prior distribution specifies how strongly we wish to bias the parameters towards zero.

By combining the log-likelihood criterion with the prior we obtain a *penalized log-likelihood* function (penalized by the prior):

$$\begin{aligned} l'(\underline{\theta}, \theta_0, \sigma^2) &= \sum_{t=1}^n \log \left[\frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(y_t - \underline{\theta}^T \underline{x}_t - \theta_0)^2\right) \right] + \log P(\underline{\theta}, \theta_0; \sigma'^2) \quad (43) \\ &= \text{const.} - \frac{n}{2} \log \sigma^2 - \frac{1}{2\sigma^2} \sum_{t=1}^n (y_t - \underline{\theta}^T \underline{x}_t - \theta_0)^2 \\ &\quad - \frac{1}{2\sigma'^2} (\theta_0^2 + \sum_{j=1}^d \theta_j^2) - \frac{d+1}{2} \log \sigma'^2 \quad (44) \end{aligned}$$

It is convenient to tie the prior variance σ'^2 to the noise variance σ^2 according to $\sigma'^2 = \sigma^2/\lambda$. This has the effect that if the noise variance σ^2 is large, we penalize the parameters very little (permit large deviations from zero by assuming a large σ'^2). On the other hand, if the noise variance is small, we could be *over-fitting* the linear model. This happens, for example, when the number of training examples is small. In this case most of the responses can be explained directly by the linear model making the noise variance very small. In such cases our penalty for the parameters will be larger as well (prior variance is smaller).

Incorporating this parameter tie into the penalized log-likelihood function gives

$$\begin{aligned} l'(\underline{\theta}, \theta_0, \sigma^2) &= \text{const.} - \frac{n}{2} \log \sigma^2 - \frac{1}{2\sigma^2} \sum_{t=1}^n (y_t - \underline{\theta}^T \underline{x}_t - \theta_0)^2 \\ &\quad - \frac{\lambda}{2\sigma^2} (\theta_0^2 + \sum_{j=1}^d \theta_j^2) - \frac{d+1}{2} \log(\sigma^2/\lambda) \quad (45) \end{aligned}$$

$$= \text{const.} - \frac{n+d+1}{2} \log \sigma^2 + \frac{d+1}{2} \log \lambda \quad (46)$$

$$- \frac{1}{2\sigma^2} \left[\sum_{t=1}^n (y_t - \underline{\theta}^T \underline{x}_t - \theta_0)^2 + \lambda (\theta_0^2 + \sum_{j=1}^d \theta_j^2) \right] \quad (47)$$

where again the estimation of $\underline{\theta}$ and θ_0 separates from setting the noise variance σ^2 . Note that this separation is achieved because we tied the prior and noise variance parameters. The above regularized problem of finding the parameter estimates $\hat{\underline{\theta}}$ and $\hat{\theta}_0$ is known as *Ridge regression*.

As before, we can get closed form estimates for the parameters (we omit the analogous derivation):

$$\begin{bmatrix} \hat{\underline{\theta}} \\ \hat{\theta}_0 \end{bmatrix} = (\lambda \mathbf{I} + \mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \quad (48)$$

It is now useful to understand how the properties of these parameter estimates depend on λ . For example, are the parameter estimates *unbiased*? No, they are not:

$$E \left\{ \begin{bmatrix} \hat{\underline{\theta}} \\ \hat{\theta}_0 \end{bmatrix} \middle| \mathbf{X} \right\} = (\lambda \mathbf{I} + \mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{X} \begin{bmatrix} \underline{\theta}^* \\ \theta_0^* \end{bmatrix} \quad (49)$$

$$= (\lambda \mathbf{I} + \mathbf{X}^T \mathbf{X})^{-1} (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I} - \lambda \mathbf{I}) \begin{bmatrix} \underline{\theta}^* \\ \theta_0^* \end{bmatrix} \quad (50)$$

$$= \begin{bmatrix} \underline{\theta}^* \\ \theta_0^* \end{bmatrix} \overbrace{-\lambda (\lambda \mathbf{I} + \mathbf{X}^T \mathbf{X})^{-1} \begin{bmatrix} \underline{\theta}^* \\ \theta_0^* \end{bmatrix}}^{\text{bias}} \quad (51)$$

$$= (\mathbf{I} - \lambda (\lambda \mathbf{I} + \mathbf{X}^T \mathbf{X})^{-1}) \begin{bmatrix} \underline{\theta}^* \\ \theta_0^* \end{bmatrix} \quad (52)$$

It is straightforward to check that $(\mathbf{I} - \lambda (\lambda \mathbf{I} + \mathbf{X}^T \mathbf{X})^{-1})$ is a positive definite matrix with eigenvalues all less than one. The parameter estimates are therefore shrunk towards zero and more so the larger the value of λ . This is what we would expect since we explicitly favored small parameter values with the prior penalty. What do we gain from such *biased* parameter estimates? Let's evaluate the mean squared error, starting with the covariance:

$$Cov \left\{ \begin{bmatrix} \hat{\underline{\theta}} \\ \hat{\theta}_0 \end{bmatrix} \middle| \mathbf{X} \right\} = \sigma^{*2} (\lambda \mathbf{I} + \mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{X} (\lambda \mathbf{I} + \mathbf{X}^T \mathbf{X})^{-1} \quad (53)$$

$$= \sigma^{*2} (\lambda \mathbf{I} + \mathbf{X}^T \mathbf{X})^{-1} (\lambda \mathbf{I} + \mathbf{X}^T \mathbf{X} - \lambda \mathbf{I}) (\lambda \mathbf{I} + \mathbf{X}^T \mathbf{X})^{-1} \quad (54)$$

$$= \sigma^{*2} (\lambda \mathbf{I} + \mathbf{X}^T \mathbf{X})^{-1} - \lambda \sigma^{*2} (\lambda \mathbf{I} + \mathbf{X}^T \mathbf{X})^{-2} \quad (55)$$

The mean squared error in the parameters is therefore given by (we again omit the deriva-

tion that can be obtained similarly to previous expressions):

$$E \left\{ \left\| \begin{bmatrix} \hat{\theta} \\ \hat{\theta}_0 \end{bmatrix} - \begin{bmatrix} \theta^* \\ \theta_0^* \end{bmatrix} \right\|^2 \mid \mathbf{X} \right\} = \sigma^{*2} \cdot \text{Tr} [(\lambda \mathbf{I} + \mathbf{X}^T \mathbf{X})^{-1} - \lambda(\lambda \mathbf{I} + \mathbf{X}^T \mathbf{X})^{-2}] \\ + \lambda^2 \begin{bmatrix} \theta^* \\ \theta_0^* \end{bmatrix}^T (\lambda \mathbf{I} + \mathbf{X}^T \mathbf{X})^{-2} \begin{bmatrix} \theta^* \\ \theta_0^* \end{bmatrix} \quad (56)$$

Can this be smaller than the mean squared error corresponding to the unregularized estimates $\sigma^{*2} \cdot \text{Tr} [(\mathbf{X}^T \mathbf{X})^{-1}]$? Yes, it can. This is indeed the benefit from regularization: we can reduce large variance at the cost of introducing a bit of bias. We will get back to this trade-off in the context of *model selection*.

Let's exemplify the effect of λ on the mean squared error in a context of a very simple 1-dimensional example. Suppose, we have observed responses for only two points, $x = -1$ and $x = 1$. In this case,

$$\mathbf{X} = \begin{bmatrix} -1 & 1 \\ 1 & 1 \end{bmatrix}, \quad \mathbf{X}^T \mathbf{X} = \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix}, \quad (\lambda \mathbf{I} + \mathbf{X}^T \mathbf{X})^{-1} = \begin{bmatrix} 1/(2+\lambda) & 0 \\ 0 & 1/(2+\lambda) \end{bmatrix} \quad (57)$$

The expression for the mean squared error therefore becomes

$$E \left\{ \left\| \begin{bmatrix} \hat{\theta} \\ \hat{\theta}_0 \end{bmatrix} - \begin{bmatrix} \theta^* \\ \theta_0^* \end{bmatrix} \right\|^2 \mid \mathbf{X} \right\} = \sigma^{*2} \left(\frac{2}{(2+\lambda)} - \frac{2\lambda}{(2+\lambda)^2} \right) + \frac{\lambda^2}{(2+\lambda)^2} (\underline{\theta}^{*2} + \theta_0^{*2}) \\ = \frac{4\sigma^{*2}}{(2+\lambda)^2} + \frac{\lambda^2}{(2+\lambda)^2} (\underline{\theta}^{*2} + \theta_0^{*2}) \quad (58)$$

We should compare this to $\sigma^{*2} \text{Tr} [(\mathbf{X}^T \mathbf{X})^{-1}] = \sigma^{*2}$ obtained without regularization (corresponds to setting $\lambda = 0$). In the noisy case $\sigma^{*2} > \underline{\theta}^{*2} + \theta_0^{*2}$ we can set $\lambda = 2$ and obtain

$$E \left\{ \left\| \begin{bmatrix} \hat{\theta} \\ \hat{\theta}_0 \end{bmatrix} - \begin{bmatrix} \theta^* \\ \theta_0^* \end{bmatrix} \right\|^2 \mid \mathbf{X} \right\} = \frac{4\sigma^{*2}}{16} + \frac{4}{16} (\underline{\theta}^{*2} + \theta_0^{*2}) < \frac{8\sigma^{*2}}{16} = \frac{1}{2} \sigma^{*2} \quad (59)$$

The mean squared error of the parameters is therefore clearly smaller than without regularization.

Active learning

We can use the expressions for the mean squared error to actively select input points x_1, \dots, x_n , when possible, so as to reduce the resulting estimation error. This is an *active*

learning (experiment design) problem. By letting the method guide the selection of the training examples (inputs), we will generally need far fewer examples in comparison to selecting them at random from some underlying distribution, database, or trying available experiments at random.

To develop this further, recall that we continue to assume that the responses y come from some linear model $y = \underline{\theta}^{*T} \underline{x} + \theta_0^* + \epsilon$ where $\epsilon \sim N(0, \sigma^{*2})$. Nothing is assumed about the distribution of \underline{x} as the choice of the inputs is in our control. For any given set of inputs, $\underline{x}_1, \dots, \underline{x}_n$, we derived last time an expression for the mean squared error of the maximum likelihood parameter estimates $\hat{\underline{\theta}}$ and $\hat{\theta}_0$:

$$E \left\{ \left\| \begin{bmatrix} \hat{\underline{\theta}} \\ \hat{\theta}_0 \end{bmatrix} - \begin{bmatrix} \underline{\theta}^* \\ \theta_0^* \end{bmatrix} \right\|^2 \mid \mathbf{X} \right\} = \sigma^{*2} \text{Tr} [(\mathbf{X}^T \mathbf{X})^{-1}] \quad (60)$$

where the expectation is relative to the responses generated from the underlying linear model, i.e., over different training sets generated from the linear model. We do not know the noise variance σ^{*2} for the correct model but it only appears as a multiplicative constant in the above expression and therefore won't affect how we should choose the inputs. When the choice of inputs is indeed up to us (e.g., which experiments to carry out) we can select them so as to minimize $\text{Tr} [(\mathbf{X}^T \mathbf{X})^{-1}]$. One caveat of this approach is that it relies on the underlying relationship between the inputs and the responses to be linear. When this is no longer the case we may end up with clearly suboptimal selections.

Given the selection criterion, how should we find say n input examples $\underline{x}_1, \dots, \underline{x}_n$ that minimize it? One simple approach is to select them one after the other, merely optimizing the selection of the next one in light of what we already have. Let's assume then that we already have \mathbf{X} and thus have $\mathbf{A} = (\mathbf{X}^T \mathbf{X})^{-1}$ (assuming it is already invertible). We are trying to select another input example \underline{x} that adds a row $[\underline{x}^T, 1]$ to \mathbf{X} . In an applied context we are typically constrained by what \underline{x} can be (e.g., due to the experimental setup). We will discuss simple constraints below. Let's now evaluate the effect of adding a new (valid) row:

$$\begin{bmatrix} \mathbf{X} \\ \underline{x}^T \ 1 \end{bmatrix}^T \begin{bmatrix} \mathbf{X} \\ \underline{x}^T \ 1 \end{bmatrix} = (\mathbf{X}^T \mathbf{X}) + \begin{bmatrix} \underline{x} \\ 1 \end{bmatrix} \begin{bmatrix} \underline{x} \\ 1 \end{bmatrix}^T = \mathbf{A}^{-1} + \underline{v} \underline{v}^T \quad (61)$$

where $\underline{v} = [\underline{x}^T, 1]^T$. We would like to find a valid \underline{v} that minimizes

$$\text{Tr} [(\mathbf{A}^{-1} + \underline{v} \underline{v}^T)^{-1}] \quad (62)$$

The matrix inverse can actually be carried out in closed form (easy enough to check)

$$(\mathbf{A}^{-1} + \underline{v} \underline{v}^T)^{-1} = \mathbf{A} - \frac{1}{(1 + \underline{v}^T \mathbf{A} \underline{v})} \mathbf{A} \underline{v} \underline{v}^T \mathbf{A} \quad (63)$$

so that the trace becomes

$$\text{Tr} [(\mathbf{A}^{-1} + \underline{v}\underline{v}^T)^{-1}] = \text{Tr} [\mathbf{A}] - \frac{1}{(1 + \underline{v}^T \mathbf{A} \underline{v})} \text{Tr} [\mathbf{A} \underline{v} \underline{v}^T \mathbf{A}] \quad (64)$$

$$= \text{Tr} [\mathbf{A}] - \frac{1}{(1 + \underline{v}^T \mathbf{A} \underline{v})} \text{Tr} [\underline{v}^T \mathbf{A} \mathbf{A} \underline{v}] \quad (65)$$

$$= \text{Tr} [\mathbf{A}] - \frac{\underline{v}^T \mathbf{A} \mathbf{A} \underline{v}}{(1 + \underline{v}^T \mathbf{A} \underline{v})} \quad (66)$$

Note that since $\text{Tr}[\mathbf{A}] = \text{Tr}[(\mathbf{X}^T \mathbf{X})^{-1}]$ is the mean squared error before adding the new example, any choice of \underline{v} , i.e., any additional example \underline{x} will reduce the mean squared error. We are interested in finding the one that reduces the error the most. This is the example that maximizes

$$\frac{\underline{v}^T \mathbf{A} \mathbf{A} \underline{v}}{(1 + \underline{v}^T \mathbf{A} \underline{v})} \quad (67)$$

How much can we possibly reduce the squared error? The above term is bounded by the largest eigenvalue of \mathbf{A} . In other words, with each new example we can at most remove one degree of freedom from the parameter space. If we assume no constraints on the choice of \underline{v} , the maximizing vector would be of infinite length and proportional to the eigenvector of \mathbf{A} with the largest eigenvalue (all the eigenvalues of \mathbf{A} are positive as it is an inverse of a positive definite matrix $\mathbf{X}^T \mathbf{X}$). It is indeed advantageous in linear regression to have the input points as far from each other as possible (see Figure 1). If we constrain $\|\underline{v}\| \leq c$, then the maximizing \underline{v} is the normalized eigenvector of \mathbf{A} with the largest eigenvalue, normalized such that $\|\underline{v}\| = c$. Note that it may not be possible to select this eigenvector since $\underline{v} = [\underline{x}^T, 1]^T$. Other constraints on \underline{x} will further restrict \underline{v} .

Let's take a simple example to illustrate the criterion. Suppose we have a 1-dimensional regression model $y = \theta x + \theta_0 + \epsilon$ where x is constrained to lie within $[-1, 1]$. Assume we have already observed responses for $x_1 = 1$ and $x_2 = -1$. Thus

$$\mathbf{X} = \begin{bmatrix} 1 & 1 \\ -1 & 1 \end{bmatrix}, \quad \mathbf{X}^T \mathbf{X} = \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix}, \quad \mathbf{A} = (\mathbf{X}^T \mathbf{X})^{-1} = \frac{1}{2} \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \quad (68)$$

$\underline{v} = [x, 1]^T$ and therefore $\underline{v}^T \mathbf{A} \underline{v} = (x^2 + 1)/2$ and $\underline{v}^T \mathbf{A} \mathbf{A} \underline{v} = (x^2 + 1)/4$. The criterion to be maximized becomes

$$\frac{\underline{v}^T \mathbf{A} \mathbf{A} \underline{v}}{(1 + \underline{v}^T \mathbf{A} \underline{v})} = \frac{(x^2 + 1)/4}{1 + (x^2 + 1)/2} \quad (69)$$

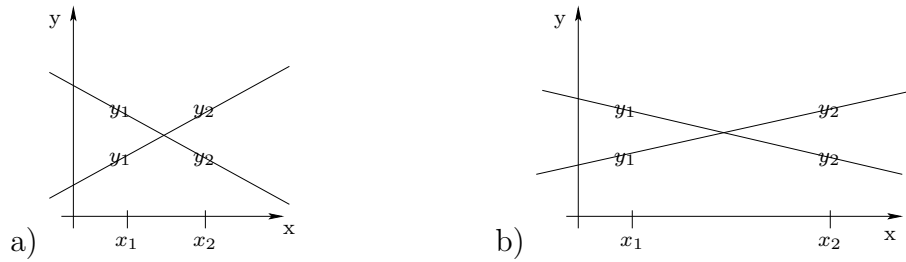


Figure 1: a) The effect of noise in the responses has a large effect on the parameters of the linear model when the corresponding inputs are close to each other; b) the effect is smaller when the inputs are further away.

Since $z/(1+z)$ is an increasing function of z , it follows that the criterion is maximized when $(x^2 + 1)/2$ is maximized. Given the constraints, the maximizing point is $x = 1$ or $x = -1$. Either choice would do but, after selecting one, the other one would be preferred at the next step. The result is consistent with the intuition that for linear models the inputs should be as far from each other as possible (cf. Figure 1).

We have so far used the mean squared error in the parameters as the selection criterion. What about the variance in the predictions? Let's try to find the point \underline{x} whose response we are the most uncertain about. We again write $\underline{v} = [\underline{x}^T, 1]^T$ so that

$$\text{Var}\{y|\mathbf{X}, \underline{x}\} = E \left\{ \left(\hat{\theta}^T \underline{x} + \hat{\theta}_0 - \underline{\theta}^{*T} \underline{x} - \underline{\theta}_0^* \right)^2 \mid \mathbf{X}, \underline{x} \right\} \tag{70}$$

$$= E \left\{ \begin{bmatrix} \underline{x} \\ 1 \end{bmatrix}^T \left(\begin{bmatrix} \hat{\theta} \\ \hat{\theta}_0 \end{bmatrix} - \begin{bmatrix} \theta^* \\ \theta_0^* \end{bmatrix} \right) \left(\begin{bmatrix} \hat{\theta} \\ \hat{\theta}_0 \end{bmatrix} - \begin{bmatrix} \theta^* \\ \theta_0^* \end{bmatrix} \right)^T \begin{bmatrix} \underline{x} \\ 1 \end{bmatrix} \mid \mathbf{X}, \underline{x} \right\} \tag{71}$$

$$= \begin{bmatrix} \underline{x} \\ 1 \end{bmatrix}^T \sigma^{*2} (\mathbf{X}^T \mathbf{X})^{-1} \begin{bmatrix} \underline{x} \\ 1 \end{bmatrix} \tag{72}$$

$$= \sigma^{*2} \cdot \underline{v}^T \mathbf{A} \underline{v} \tag{73}$$

where the expectation is over responses for existing training examples, again assuming that there is a correct underlying linear model. So the largest variance corresponds to the input \underline{x} that maximizes $\underline{v}^T \mathbf{A} \underline{v}$ where $\underline{v} = [\underline{x}^T, 1]^T$. In the unconstrained case where there are few or no restrictions on \underline{v} , this maximizing point is exactly the one we would query according to the previous selection criterion.