

Bias variance trade-off

Supplementary note for CSCI-UA 9473

Augustin Cosse

September 2018

1 Expected loss

Suppose our data set is given by $\mathcal{D} = \{(\mathbf{x}, t)\}$ where t are the labels associated to each point \mathbf{x} . We will assume that the data is distributed according to an ideal function $y(\mathbf{x})$ plus some perturbation ε with zero mean and variance σ^2 .

$$t = y(\mathbf{x}) + \varepsilon \quad (1)$$

t is thus the observed label. In order to study how well a given regression model $h(\mathbf{x})$ is fitting the data, one can look at the expected loss

$$\mathbb{E}[L] = \int \int (h(\mathbf{x}) - t)^2 p(\mathbf{x}, t) d\mathbf{x} dt \quad (2)$$

Note that given the objective (2), the optimal model one can choose for $h(\mathbf{x})$ is to solve the minimization problem

$$\min_{h(\mathbf{x})} \int \int (h(\mathbf{x}) - t)^2 d\mathbf{x} dt \quad (3)$$

To do this, we set the first order derivatives to 0, and get

$$\frac{\delta \mathbb{E}[L]}{\delta h(\mathbf{x})} = 2 \int (h(\mathbf{x}) - t) p(\mathbf{x}, t) d\mathbf{x} dt = 0 \quad (4)$$

from this we get

$$\int h(\mathbf{x}) p(\mathbf{x}) d\mathbf{x} = \int \int t p(t, \mathbf{x}) d\mathbf{x} dt \quad (5)$$

The t disappeared in the left handside because $h(\mathbf{x})$ does not depend on t so we can simply integrate $p(\mathbf{x}, t)$ over all values of t which gives $\int p(\mathbf{x}, t) dt = p(\mathbf{x})$.

This last equation tells us that in general we should at every value \mathbf{x} choose the model/prediction $h(\mathbf{x})$ corresponding to

$$h(\mathbf{x}) = \frac{\int t p(\mathbf{x}, t) dt}{p(\mathbf{x})} = \int tp(t|\mathbf{x}) dt = \mathbb{E}_t \{t|\mathbf{x}\}. \quad (6)$$

The last line follows from $p(\mathbf{x}, t) = p(t|\mathbf{x})p(\mathbf{x})$. This is however not always possible. On the other hand, we always have

$$\begin{aligned} (h(\mathbf{x}) - t)^2 &= (h(\mathbf{x}) - \mathbb{E}\{t|\mathbf{x}\} + \mathbb{E}\{t|\mathbf{x}\} - t)^2 \\ &= (h(\mathbf{x}) - \mathbb{E}\{t|\mathbf{x}\})^2 + (\mathbb{E}\{t|\mathbf{x}\} - t)^2 + 2(h(\mathbf{x}) - \mathbb{E}\{t|\mathbf{x}\})(\mathbb{E}\{t|\mathbf{x}\} - t) \end{aligned} \quad (7)$$

If you substitute this last expression into the expected Loss (2) you can see that the cross terms will disappear as it leads to the following four terms

$$2 \int h(\mathbf{x})\mathbb{E}\{t|\mathbf{x}\} p(\mathbf{x}, t) d\mathbf{x} dt = 2 \int h(\mathbf{x})\mathbb{E}\{t|\mathbf{x}\} p(\mathbf{x}) d\mathbf{x} \quad (9)$$

$$- 2 \int h(\mathbf{x})tp(\mathbf{x}, t) dt d\mathbf{x} = -2 \int \mathbb{E}\{t|\mathbf{x}\} h(\mathbf{x})p(\mathbf{x}) d\mathbf{x} \quad (10)$$

$$- 2 \int \mathbb{E}\{t|\mathbf{x}\}^2 p(\mathbf{x}, t) d\mathbf{x} dt = -2 \int \mathbb{E}\{t|\mathbf{x}\}^2 p(\mathbf{x}) d\mathbf{x} \quad (11)$$

$$2 \int \mathbb{E}\{t|\mathbf{x}\} tp(\mathbf{x}, t) d\mathbf{x} dt = 2 \int \mathbb{E}\{t|\mathbf{x}\}^2 p(\mathbf{x}) d\mathbf{x} \quad (12)$$

So that the average loss reduces to

$$\mathbb{E}[L] = \int (t - h(\mathbf{x}))^2 d\mathbf{x} dt = \int (h(\mathbf{x}) - \mathbb{E}\{t|\mathbf{x}\})^2 p(\mathbf{x}) d\mathbf{x} \quad (13)$$

$$+ \int (\mathbb{E}\{t|\mathbf{x}\} - t)^2 p(\mathbf{x}, t) d\mathbf{x} dt \quad (14)$$

$$= \int (h(\mathbf{x}) - \mathbb{E}\{t|\mathbf{x}\})^2 p(\mathbf{x}) d\mathbf{x} \quad (15)$$

$$+ \int \sigma_x^2 p(\mathbf{x}) d\mathbf{x} \quad (16)$$

Where I used

$$\int (\mathbb{E}\{t|\mathbf{x}\} - t)^2 p(\mathbf{x}, t) d\mathbf{x} dt = \int (\mathbb{E}\{t|\mathbf{x}\} - t)^2 p(t|\mathbf{x})p(\mathbf{x}) d\mathbf{x} dt \quad (17)$$

$$= \int \text{Var}[t|\mathbf{x}]p(\mathbf{x}) d\mathbf{x} \quad (18)$$

Once again, the t is integrated over in the first integral (13) because none of the functions $h(\mathbf{x})$ and $\mathbb{E}\{t|\mathbf{x}\}$ depend on t anymore.

The final expression for the average loss is then

$$\mathbb{E}[L] = \int (h(\mathbf{x}) - \mathbb{E}\{t|\mathbf{x}\})^2 p(\mathbf{x}) d\mathbf{x} + \int \sigma_x^2 p(\mathbf{x}) d\mathbf{x} \quad (19)$$

The last term is the average noise variance over all \mathbf{x} .

You see that there is a first term (MSE) that depends on the model (i.e how well the regression model captures the average data distribution $\mathbb{E}\{t|\mathbf{x}\} = y(x)$) and a second term that depends only on the data. Now we will show that the first term can be decomposed into a *variance* and a *bias* contributions.

2 Bias variance trade-off

In the derivations above, we haven't make any assumption on the regression model learned $h(\mathbf{x})$. In general, when we learn a regression model, we learn it from a subset of the data. We will denote this subset \mathcal{D}_i so that $\mathcal{D}_i \subset \mathcal{D}$ and I will now denote our regression model as $h(\mathbf{x}|\mathcal{D}_i)$ to emphasize the dependence of the model on the particular choice of \mathcal{D}_i . That being said, if we use $\bar{t}(\mathbf{x})$ to denote $\mathbb{E}\{t|\mathbf{x}\}$, we have the decomposition

$$(h(\mathbf{x}; \mathcal{D}_i) - \bar{t}(\mathbf{x}))^2 \quad (20)$$

$$= \left(h(\mathbf{x}; \mathcal{D}_i) - \mathbb{E}_{\mathcal{D}_i} \{h(\mathbf{x}; \mathcal{D}_i)\} + \mathbb{E}_{\mathcal{D}_i} \{h(\mathbf{x}; \mathcal{D}_i)\} - \bar{t}(\mathbf{x}) \right)^2 \quad (21)$$

$$= \left(h(\mathbf{x}; \mathcal{D}_i) - \mathbb{E}_{\mathcal{D}_i} \{h(\mathbf{x}; \mathcal{D}_i)\} \right)^2 + \left(\mathbb{E}_{\mathcal{D}_i} \{h(\mathbf{x}; \mathcal{D}_i)\} - \bar{t}(\mathbf{x}) \right)^2 \quad (22)$$

$$+ 2 \left(h(\mathbf{x}; \mathcal{D}_i) - \mathbb{E}_{\mathcal{D}_i} \{h(\mathbf{x}; \mathcal{D}_i)\} \right) \left(\mathbb{E}_{\mathcal{D}_i} \{h(\mathbf{x}; \mathcal{D}_i)\} - \bar{t}(\mathbf{x}) \right) \quad (23)$$

If we take the expectation over all possible datasets \mathcal{D}_i used to learn the model, we will get a measure on how well the family of models $h(\mathbf{x}, \mathcal{D}_i)$ performs on average for the different datasets. Below I use $\mathbb{E}_{\mathcal{D}_i}$ to denote the average when sampling random datasets \mathcal{D}_i from the full set of samples \mathcal{D} .

$$\mathbb{E}_{\mathcal{D}_i} \{h(\mathbf{x}; \mathcal{D}_i) - \bar{t}(\mathbf{x})\}^2 \quad (24)$$

$$= \mathbb{E}_{\mathcal{D}_i} \left(h(\mathbf{x}; \mathcal{D}_i) - \mathbb{E}_{\mathcal{D}_i} \{h(\mathbf{x}; \mathcal{D}_i)\} \right)^2 + \mathbb{E}_{\mathcal{D}_i} \left(\mathbb{E}_{\mathcal{D}_i} \{h(\mathbf{x}; \mathcal{D}_i)\} - \bar{t}(\mathbf{x}) \right)^2 \quad (25)$$

$$+ 2 \mathbb{E}_{\mathcal{D}_i} \left\{ \left(h(\mathbf{x}; \mathcal{D}_i) - \mathbb{E}_{\mathcal{D}_i} \{h(\mathbf{x}; \mathcal{D}_i)\} \right) \left(\mathbb{E}_{\mathcal{D}_i} \{h(\mathbf{x}; \mathcal{D}_i)\} - \bar{t}(\mathbf{x}) \right) \right\} \quad (26)$$

$$= \underbrace{\left(\mathbb{E}_{\mathcal{D}_i} \{y(\mathbf{x}; \mathcal{D}_i)\} - \bar{t}(\mathbf{x}) \right)^2}_{\text{bias}^2} + \underbrace{\mathbb{E}_{\mathcal{D}_i} \left[\left(h(\mathbf{x}; \mathcal{D}_i) - \mathbb{E}_{\mathcal{D}_i} \{h(\mathbf{x}; \mathcal{D}_i)\} \right)^2 \right]}_{\text{variance}} \quad (27)$$

The cross term (26) is zero because the second factor does not depend on \mathcal{D}_i (it can be taken out of the average, $\bar{t}(\mathbf{x})$ is the unknown regression function

	simpler models	complex models
squared bias	large	small
variance	small	large

Table 1: The bias variance trade-off and the model complexity. As the *model complexity increases*, the *variance* tends to *increase* and the *squared bias* tends to *decrease*. An additional illustration of this phenomenon can be found in Fig. 3.5 in [1] which is available (pdf) online.

and remains the same independently of the data we use to learn the regression model, $\bar{t}(\mathbf{x}) = y(\mathbf{x})$ and the first term has mean zero.

Complex models will usually have a large contribution to the variance term because within the class of complex models, the models will usually *vary* a lot when you change the training data. I.e if you have a very complex model with many parameters that perfectly fits the data in \mathcal{D}_i , then when you will take another dataset \mathcal{D}_j , the model that you build for \mathcal{D}_j will be completely different from the model you built with \mathcal{D}_i (think of a high degree polynomial). You will thus have large differences between models $h(\mathbf{x}; \mathcal{D}_i)$ and $h(\mathbf{x}; \mathcal{D}_j)$ and hence large deviations from the mean $h(\mathbf{x}; \mathcal{D}_i) - \mathbb{E}_{\mathcal{D}_i} \{h(\mathbf{x}; \mathcal{D}_i)\}$. So complex models will have larger variance.

The bias, on the other hand is a measure of how well the model "captures" the average behavior of the data $\mathbb{E} \{t|\mathbf{x}\}$. In general (except if the data is assumed to be very simple or linear which is the framework of the *Gauss Markov Theorem* that I discuss below), the simpler models will not be able to capture such behavior. As an example, imagine we have data of the form

$$t_\ell = y(x_\ell) + \varepsilon_\ell = \sin(x_\ell) + \varepsilon_\ell \quad (28)$$

Then without using non linear terms (and defining our new variables as $x' = \phi(x)$ where $\phi(x)$ encodes the non linearity), it is impossible to design a model of the form $h(x) = \beta_0 + \beta_1 x$ such that

$$h(x_\ell) = \sin(x_\ell) \quad (29)$$

This should illustrate the fact that simpler models have a large bias (*At least when we don't make assumptions on the data distribution*)

This relation is summarized in table 1.

3 The particular case of linearly generated data and the Gauss Markov Theorem

So far we haven't made any assumption on the data. When we don't make assumptions on the data, the relation of table 1 holds in general.

However, when we know that the data has a particular distribution, we can better describe the quality of approximation of given models. This is the framework of the *Gauss Markov Theorem* (GMT).

When the data is distributed according to a linear model, so that

$$t_\ell = y(\mathbf{x}_\ell) + \varepsilon = \langle \boldsymbol{\beta}, \mathbf{x}_\ell \rangle + \varepsilon_\ell, \quad (30)$$

with $\mathbb{E}\boldsymbol{\varepsilon} = \mathbb{E}[\varepsilon_1, \dots, \varepsilon_L] = \mathbf{0}$ and $\mathbb{E}\{\boldsymbol{\varepsilon}\boldsymbol{\varepsilon}^T\} = \sigma^2\mathbf{I}$. Then the linear model has no bias (note that this might seem counterintuitive compared to table 1 but this is because here we assumed that the data is linear).

Moreover, the *Gauss Markov theorem* states that among all possible linear estimators of the $\boldsymbol{\beta}$, the estimator $\hat{\boldsymbol{\beta}}_{\text{RSS}}$ which is computed by minimizing the residual sum of squares,

$$\hat{\boldsymbol{\beta}}_{\text{RSS}} = \underset{\boldsymbol{\beta}}{\operatorname{argmin}} \sum_{\ell=1}^L (t_\ell - \langle \boldsymbol{\beta}, \mathbf{x}_\ell \rangle)^2, \quad (31)$$

is the best estimator.

In the Gauss Markov framework, when we say data, we mean the t_ℓ . What we now call data are the noisy measurements $t_\ell = \langle \boldsymbol{\beta}, \mathbf{x}_\ell \rangle$. Because we've put ourselves in a framework where we now know how the relation between t_ℓ and \mathbf{x}_ℓ , $t_\ell(\mathbf{x}) = \beta_0 + \langle \boldsymbol{\beta}_1, \mathbf{x}_\ell \rangle + \varepsilon$, we can focus on denoising the t_ℓ . In this particular framework, a linear estimator for the coefficients $\boldsymbol{\beta}$ is an estimator of the form

$$\tilde{\boldsymbol{\beta}} = \mathbf{M}\mathbf{t} \quad (32)$$

where $\mathbf{t} = [t_1, t_2, \dots, t_N]$ encode the labels. The estimator provided by the linear regression (i.e. RSS) approach (31) is a particular such estimator as it can read as (see my slides)

$$\hat{\boldsymbol{\beta}}_{\text{RSS}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{t} \quad (33)$$

$$= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}) \quad (34)$$

where $\mathbf{t} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$, $\boldsymbol{\varepsilon} = (\varepsilon_1, \dots, \varepsilon_L)$ and \mathbf{X} is the matrix

$$\mathbf{X} = \begin{bmatrix} 1 & \mathbf{x}_1^T \\ 1 & \mathbf{x}_2^T \\ \vdots & \vdots \\ 1 & \mathbf{x}_L^T \end{bmatrix} \quad (35)$$

To prove the GMT, we consider any other linear estimator $\tilde{\boldsymbol{\beta}}$ of the coefficients $\boldsymbol{\beta} = (\beta_0, \boldsymbol{\beta}_1)$, $\tilde{\boldsymbol{\beta}} = \hat{\boldsymbol{\beta}}_{\text{RSS}} + \Delta\mathbf{t}$ where Δ is a perturbation.

We first require the linear estimator $\tilde{\beta}$ to be unbiased, which gives

$$\mathbb{E}\{\tilde{\beta}\} = \mathbb{E}\{\hat{\beta}\} + \mathbb{E}\{\Delta\mathbf{t}\} \quad (36)$$

$$= \mathbb{E}\{\Delta\mathbf{X}\beta + \Delta\varepsilon\} \quad (37)$$

$$= \Delta\mathbf{X}\beta + 0 \quad (38)$$

$$= 0 \quad (39)$$

The only source of randomness here is the noise ε and we assumed that $\mathbb{E}\varepsilon = 0$. As we haven't made any assumption on β , the last line implies that for the linear estimator to be unbiased, we must have

$$\Delta\mathbf{X} = 0 \quad (40)$$

Now, for the variance, first note that

$$\text{Var}\{\hat{\beta}_{\text{RSS}}\} = \mathbb{E}\{(\hat{\beta}_{\text{RSS}} - \mathbb{E}\hat{\beta}_{\text{RSS}})(\hat{\beta}_{\text{RSS}} - \mathbb{E}\hat{\beta}_{\text{RSS}})^T\} \quad (41)$$

$$= \mathbb{E}\{\hat{\beta}_{\text{RSS}}\hat{\beta}_{\text{RSS}}^T\} - \beta\beta^T \quad (42)$$

as the estimator is unbiased. For any other unbiased estimator $\tilde{\beta}$ of the form $\tilde{\beta} = \hat{\beta}_{\text{RSS}} + \Delta\mathbf{t}$ with zero bias, we have (the expectation is over the noise ε)

$$\text{Var}\{\tilde{\beta}\} \quad (43)$$

$$= \mathbb{E}\{(\tilde{\beta} - \mathbb{E}\tilde{\beta})(\tilde{\beta} - \mathbb{E}\tilde{\beta})^T\} \quad (44)$$

$$= \mathbb{E}\left\{\left(\hat{\beta}_{\text{RSS}} + \Delta\mathbf{t} - \beta - \mathbb{E}\{\Delta\mathbf{t}\}\right)\left(\hat{\beta}_{\text{RSS}} + \Delta\mathbf{t} - \beta - \mathbb{E}\{\Delta\mathbf{t}\}\right)^T\right\} \quad (45)$$

$$= \mathbb{E}\left\{\left(\hat{\beta}_{\text{RSS}} - \beta\right)\left(\hat{\beta}_{\text{RSS}} - \beta\right)^T\right\} \quad (46)$$

$$+ \mathbb{E}\left\{(\Delta\mathbf{t} - \mathbb{E}\{\Delta\mathbf{t}\})(\Delta\mathbf{t} - \mathbb{E}\{\Delta\mathbf{t}\})^T\right\} \quad (47)$$

$$- \mathbb{E}\left\{(\hat{\beta}_{\text{RSS}} - \beta)(\Delta\mathbf{t} - \mathbb{E}\{\Delta\mathbf{t}\})^T\right\} \quad (48)$$

$$- \mathbb{E}\left\{(\Delta\mathbf{t} - \mathbb{E}\{\Delta\mathbf{t}\})(\hat{\beta}_{\text{RSS}} - \beta)^T\right\} \quad (49)$$

The last two (cross) terms disappear simplify as

$$- \mathbb{E}\left\{(\hat{\beta}_{\text{RSS}} - \beta)(\Delta\mathbf{t} - \mathbb{E}\{\Delta\mathbf{t}\})^T\right\} - \mathbb{E}\left\{(\Delta\mathbf{t} - \mathbb{E}\{\Delta\mathbf{t}\})(\hat{\beta}_{\text{RSS}} - \beta)^T\right\} \quad (50)$$

$$= -\mathbb{E}\left\{(\hat{\beta}_{\text{RSS}} - \beta)(\Delta\mathbf{t})^T\right\} - \mathbb{E}\left\{(\Delta\mathbf{t})(\hat{\beta}_{\text{RSS}} - \beta)^T\right\} \quad (51)$$

$$= \mathbb{E}\left\{\hat{\beta}_{\text{RSS}}(\Delta\mathbf{t})^T\right\} - \mathbb{E}\left\{(\Delta\mathbf{t})\hat{\beta}_{\text{RSS}}^T\right\} \quad (52)$$

$$= -\mathbb{E}\left\{\hat{\beta}_{\text{RSS}}(\Delta\varepsilon)^T\right\} - \mathbb{E}\left\{(\Delta\varepsilon)\hat{\beta}_{\text{RSS}}^T\right\} \quad (53)$$

The second line (51) comes from $\mathbb{E}\{\Delta\mathbf{t}\} = \Delta\mathbf{X}\boldsymbol{\beta} = 0$, (52) comes from the fact that $\boldsymbol{\beta}$ does not depend on the noise ε and $\mathbb{E}\{\Delta\mathbf{t}\} = 0$, the third line (53) follows from developing $\Delta\mathbf{t} = \Delta\mathbf{X}\boldsymbol{\beta} + \Delta\varepsilon$ and using $\Delta\mathbf{X} = 0$. For the last line, developing the expression of $\hat{\boldsymbol{\beta}}_{RSS}$, we get

$$- \mathbb{E}\left\{\hat{\boldsymbol{\beta}}_{RSS}(\Delta\varepsilon)^T\right\} - \mathbb{E}\left\{(\Delta\varepsilon)\hat{\boldsymbol{\beta}}_{RSS}^T\right\} \quad (54)$$

$$= -\mathbb{E}(\mathbf{X}\mathbf{X})^{-1}\mathbf{X}^T(\mathbf{X}\boldsymbol{\beta} + \varepsilon)\varepsilon^T\Delta^T \quad (55)$$

$$- \mathbb{E}\Delta\varepsilon(\mathbf{X}\boldsymbol{\beta} + \varepsilon)^T\mathbf{X}(\mathbf{X}\mathbf{X})^{-T} \quad (56)$$

Then use $\mathbb{E}\varepsilon = 0$ as well as $\mathbb{E}\varepsilon\varepsilon^T = \sigma^2\mathbf{I}$ and finally $\Delta\mathbf{X} = \mathbf{X}^T\Delta^T = 0$. All of this gives

$$\text{Var}\left\{\tilde{\boldsymbol{\beta}}\right\} = \text{Var}\left\{\hat{\boldsymbol{\beta}}_{RSS}\right\} + \mathbb{E}\left\{(\Delta\mathbf{t})(\Delta\mathbf{t})^T\right\} \quad (57)$$

$$- \mathbb{E}\left\{(\Delta\mathbf{t})\right\}\mathbb{E}\left\{(\Delta\mathbf{t})\right\}^T \quad (58)$$

$$= \text{Var}\left\{\hat{\boldsymbol{\beta}}_{RSS}\right\} + \Delta\mathbf{X}\boldsymbol{\beta}\boldsymbol{\beta}^T\mathbf{X}^T\Delta^T \quad (59)$$

$$+ \Delta\varepsilon\varepsilon^T\Delta^T \quad (60)$$

$$= \text{Var}\left\{\hat{\boldsymbol{\beta}}_{RSS}\right\} + \sigma^2\Delta\Delta^T \quad (61)$$

As soon as $\Delta \neq 0$, the diagonal of $\Delta\Delta^T$ is non zero as well. This in particular means that there is one component of $\tilde{\boldsymbol{\beta}}$ for which we always have

$$\mathbb{E}_\varepsilon\left\{(\tilde{\beta}_k - \beta_k)^2\right\} > \mathbb{E}_\varepsilon\left\{(\hat{\beta}_k - \beta_k)^2\right\} \quad (62)$$

Both estimators have zero bias (we will overestimate and underestimate the truth equally) but the errors we make will often be larger when we don't use the RSS estimator.

In conclusion, you can thus see that any linear estimator with no bias will always increase the variance with respect to $\hat{\boldsymbol{\beta}}_{RSS}$ (This is because we are in the particular framework of linearly generated data)

For this reason, the estimator $\hat{\boldsymbol{\beta}}_{RSS}$ is sometimes called the BLUE (Best Linear Unbiased Estimator)

3.1 Increasing the bias to reduce generalization under the assumption of linearly generated data

That being said, the fact that the BLUE is the best estimator among all unbiased estimators with respect to linearly generated data, does not mean that it is not possible to do better in terms of future predictions. And this is because a similar *bias variance trade-off* as the one discussed in section 2 applies.

Imagine that we take a new point \mathbf{x}_0 and we apply our estimator $\tilde{\beta}$ (*now we don't assume that this estimator is unbiased anymore*, this is just a general linear estimator) to this point to get an estimated label $\langle \tilde{\beta}, \mathbf{x}_0 \rangle$ (which is a prediction for $\langle \beta, \mathbf{x}_0 \rangle$). We can then write the expected mean square prediction error as before

$$\mathbb{E}_\varepsilon \left\{ \left(\langle \tilde{\beta}, \mathbf{x}_0 \rangle - \langle \beta, \mathbf{x}_0 \rangle \right)^2 \right\} \quad (63)$$

$$= \mathbb{E}_\varepsilon \left\{ \left(\langle \tilde{\beta}, \mathbf{x}_0 \rangle - \mathbb{E} \left\{ \langle \tilde{\beta}, \mathbf{x}_0 \rangle \right\} + \mathbb{E} \left\{ \langle \tilde{\beta}, \mathbf{x}_0 \rangle \right\} - \langle \beta, \mathbf{x}_0 \rangle \right)^2 \right\} \quad (64)$$

$$= \mathbb{E}_\varepsilon \left\{ \left(\langle \tilde{\beta}, \mathbf{x}_0 \rangle - \mathbb{E} \left\{ \langle \tilde{\beta}, \mathbf{x}_0 \rangle \right\} \right)^2 \right\} \quad (65)$$

$$+ \mathbb{E}_\varepsilon \left\{ \left(\mathbb{E} \left\{ \langle \tilde{\beta}, \mathbf{x}_0 \rangle \right\} - \langle \beta, \mathbf{x}_0 \rangle \right)^2 \right\} \quad (66)$$

$$+ 2\mathbb{E}_\varepsilon \left\{ \left(\langle \tilde{\beta}, \mathbf{x}_0 \rangle - \mathbb{E}_\varepsilon \left\{ \tilde{\beta}, \mathbf{x}_0 \right\} \right) \left(\mathbb{E} \left\{ \langle \tilde{\beta}, \mathbf{x}_0 \rangle \right\} - \langle \beta, \mathbf{x}_0 \rangle \right) \right\} \quad (67)$$

The last term vanishes as the second factor is deterministic (i.e does not depend on ε) and the first factor has mean 0. We thus once again have

$$\mathbb{E}_\varepsilon \left\{ \left(\langle \tilde{\beta}, \mathbf{x}_0 \rangle - \langle \beta, \mathbf{x}_0 \rangle \right)^2 \right\} = \underbrace{\mathbb{E}_\varepsilon \left\{ \left(\langle \tilde{\beta}, \mathbf{x}_0 \rangle - \mathbb{E} \left\{ \langle \tilde{\beta}, \mathbf{x}_0 \rangle \right\} \right)^2 \right\}}_{\text{variance}} \quad (68)$$

$$+ \underbrace{\mathbb{E}_\varepsilon \left\{ \left(\mathbb{E} \left\{ \langle \tilde{\beta}, \mathbf{x}_0 \rangle \right\} - \langle \beta, \mathbf{x}_0 \rangle \right)^2 \right\}}_{\text{bias}^2} \quad (69)$$

This shows that even in the framework of linearly generated data, despite the fact that $\hat{\beta}$ is the BLUE estimator, it does not mean that this estimator will be the best at predicting a new value of $t = \langle \mathbf{x}, \beta \rangle$ (on average over the noise). There might be an estimator $\tilde{\beta}$ with a non zero bias (on linear data) that will have smaller variance.

If I generate many values $t_k = \langle \beta, \mathbf{x}_k \rangle + \varepsilon_k$, learn my estimator $\hat{\beta}$ from those values, then this estimator will be unbiased but when I will want to get the prediction for a new point x_0 , the average prediction might be better if I lower the variance and increase the bias a little. Again this might sound counter intuitive with respect to table 1 but remember that we assumed linearly distributed data here.

3.2 To go further: Ridge regression*

Another way to understand this bias-variance trade-off in the case of linear data, is to look at the variance of the RSS estimator $\hat{\beta}_{\text{RSS}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{X} \beta + \varepsilon)$ estimator. This variance (42) has the form

$$\mathbb{E} \left\{ \hat{\beta}_{\text{RSS}} \hat{\beta}_{\text{RSS}}^T \right\} - \mathbb{E} \left\{ \hat{\beta}_{\text{RSS}} \right\} \left(\mathbb{E} \left\{ \hat{\beta}_{\text{RSS}} \right\} \right)^T \quad (70)$$

$$= \mathbb{E} \left\{ (\mathbf{X} \mathbf{X}^T)^{-1} \mathbf{X}^T (\mathbf{X} \beta + \varepsilon) (\mathbf{X} \beta + \varepsilon)^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-T} \right\} - \beta \beta^T \quad (71)$$

When variables (entries in $\mathbf{x} = [x_1, \dots, x_L]$) are correlated, the matrix $\mathbf{X}^T \mathbf{X}$ will be badly scaled, i.e. $(\mathbf{X}^T \mathbf{X})^{-1}$ will be very large, and the prediction error for this estimator will thus be big. A biased alternative (see my slides for Lecture 3) is to add a regularization (for example ridge regression). In this case, the variance becomes

$$\mathbb{E} \left\{ (\tilde{\beta}_{RR} - \mathbb{E} \tilde{\beta}_{RR}) (\tilde{\beta}_{RR} - \mathbb{E} \tilde{\beta}_{RR})^T \right\} \quad (72)$$

$$= \mathbb{E} \left\{ (\mathbf{X} \mathbf{X}^T + \lambda \mathbf{I})^{-1} \mathbf{X}^T (\mathbf{X} \beta + \varepsilon) (\mathbf{X} \beta + \varepsilon)^T \mathbf{X} (\mathbf{X} \mathbf{X}^T + \lambda \mathbf{I})^{-T} \right\} \quad (73)$$

$$- \mathbb{E} \tilde{\beta}_{RR} (\mathbb{E} \tilde{\beta}_{RR})^T \quad (74)$$

which will be smaller than inverting a badly conditioned matrix because of the identity which is well conditioned.

References

- [1] C. M. Bishop. Pattern recognition and machine learning. *Springer*, 2006.