

# Statistical Learning

Tommaso Proietti

DEF Tor Vergata

Linear Regression

# Specification

Let  $Y$  be a univariate quantitative response variable. We model  $Y$  as follows:

$$Y = f(X) + \varepsilon$$

where  $f(X)$  is the systematic part (that can be predicted using the inputs  $X$ ) and  $\varepsilon$  is a disturbance (error) term, accounting for the all the variation sources different from  $X$ .

Assumptions:

- Linearity. The regression function is linear in the inputs:

$$f(X) = \beta_0 + \sum_{j=1}^p \beta_j X_j$$

- Weak exogeneity:  $E(\varepsilon|X) = 0$ , i.e.  $X$  and  $\varepsilon$  are independent in the mean.
- Homoscedasticity:  $E(\varepsilon^2|X) = \text{Var}(\varepsilon|X) = \sigma^2$ .

The inputs set  $X$  includes quantitative variables, nonlinear transformations and basis expansions (to be defined later), as well as dummy variables coding the levels of qualitative inputs.

Under these assumptions the regression function is interpreted as the conditional mean of  $Y$  given  $X$

$$f(X) = E(Y|X).$$

This is the optimal predictor of  $Y$  under square loss (minimum mean square estimator of  $Y$ ).

# Estimation

Set of training data  $\{(y_i, \mathbf{x}_i), i = 1, \dots, N\}$ , with  $\mathbf{x}_i' = (x_{i1}, x_{i2}, \dots, x_{ip})$ .

We assume that the training sample is a sample of size  $N$ , drawn independently (with replacement) from the underlying population, so that for the  $i$ -th individual we can write

$$y_i = \beta_0 + \sum_{j=1}^p \beta_j x_{ij} + \varepsilon_i, \quad i = 1, \dots, N$$

Writing the  $N$  equations in matrix form:

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$$

with  $E(\boldsymbol{\varepsilon}|\mathbf{X}) = \mathbf{0}$  and  $\text{Var}(\boldsymbol{\varepsilon}|\mathbf{X}) = \sigma^2 \mathbf{I}$ . This means that  $\varepsilon_i$ 's have the same variance and are uncorrelated: the error for the  $i$ -th unit is not correlated to that of any other unit in the sample.

$$\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_i \\ \vdots \\ y_N \end{bmatrix}, \mathbf{X} = \begin{bmatrix} 1 & x_{11} & x_{12} & \dots & x_{1k} & \dots & x_{1p} \\ 1 & x_{21} & x_{22} & \dots & x_{2k} & \dots & x_{2p} \\ \vdots & \vdots & \vdots & \dots & \dots & \dots & \vdots \\ 1 & x_{i1} & x_{i2} & \dots & x_{ik} & \dots & x_{ip} \\ \vdots & \vdots & \vdots & \dots & \dots & \dots & \vdots \\ 1 & x_{N1} & x_{N2} & \dots & x_{Nk} & \dots & x_{Np} \end{bmatrix}, \boldsymbol{\varepsilon} = \begin{bmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_i \\ \vdots \\ \varepsilon_N \end{bmatrix}$$

$$\boldsymbol{\beta} = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \vdots \\ \beta_k \\ \vdots \\ \beta_p \end{bmatrix}$$

We aim at estimating the unknown parameters  $\beta_0, \dots, \beta_p$ , and  $\sigma^2$ .

# Learning by OLS

An important estimation method is Ordinary Least Squares (OLS).

OLS obtains an estimate of  $\beta$  denoted  $\hat{\beta}$ , as the minimizers of the residual sum of squares

$$\text{RSS}(\beta) = (\mathbf{y} - \mathbf{X}\beta)'(\mathbf{y} - \mathbf{X}\beta).$$

There is a closed form solution to the above problem.

Computing the partial derivatives of  $\text{RSS}(\beta)$  w.r.t. each  $\beta_j$  and equating to zero yields a  $p + 1$  linear system of equations in  $p + 1$  unknowns,

$$\mathbf{X}'\mathbf{X}\beta = \mathbf{X}'\mathbf{y}$$

which is solved w.r.t.  $\beta$ . The solution is unique provided that the  $X$ 's are not collinear.

In matrix notation the solution is

$$\hat{\beta} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}$$

Example:  $p = 1$  (simple linear regression)

$$\begin{aligned} N\beta_0 + \beta_1 \sum_i x_i &= \sum_i y_i \\ \beta_0 \sum_i x_i + \beta_1 \sum_i x_i^2 &= \sum_i x_i y_i \end{aligned}$$

We solve this system by substitution: solving wrt  $\beta_0$  the 1st eqn. yields

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

Replacing into the 2nd and solving wrt  $\hat{\beta}_1$  gives:

$$\hat{\beta}_1 = \frac{\frac{1}{N} \sum_i x_i y_i - \bar{x} \bar{y}}{\frac{1}{N} \sum_i x_i^2 - \bar{x}^2} = \frac{\text{Cov}(x, y)}{\text{Var}(x)} = \frac{\sum_i (x_i - \bar{x})(y_i - \bar{y})}{\sum_i (x_i - \bar{x})^2}$$

# Predicted values and residuals

## Predicted (fitted) values

$$\hat{y}_i = \mathbf{x}'_i \hat{\boldsymbol{\beta}} = \hat{\beta}_0 + \hat{\beta}_1 x_{i1} + \dots + \hat{\beta}_p x_{ip}, \quad i = 1, \dots, N.$$

In matrix notation,

$$\hat{\mathbf{y}} = \mathbf{X} \hat{\boldsymbol{\beta}} = \mathbf{X} (\mathbf{X}' \mathbf{X})^{-1} \mathbf{X}' \mathbf{y} = \mathbf{H} \mathbf{y},$$

$\mathbf{H} = \mathbf{X} (\mathbf{X}' \mathbf{X})^{-1} \mathbf{X}'$  is the “hat” matrix, projecting  $\mathbf{y}$  orthogonally onto the vector space generated by  $\mathbf{X}$ .

## Residuals

The OLS residuals are  $e_i = y_i - \hat{y}_i$ ,  $i = 1, \dots, N$ . In matrix notation,

$$\mathbf{e} = \mathbf{y} - \hat{\mathbf{y}} = (\mathbf{I} - \mathbf{H}) \mathbf{y} = \mathbf{M} \mathbf{y},$$

where  $\mathbf{M} = \mathbf{I} - \mathbf{H}$ .



# Properties of $\hat{\mathbf{y}}$ and $\mathbf{e}$

The fitted values and the LS residuals have properties that are simple to show and are good to know.

- Orthogonality.  $\mathbf{X}'\mathbf{e} = \mathbf{0}$ . The residuals are uncorrelated with each of the variable in  $\mathbf{X}$ . In particular,  $\sum_i e_i = 0$  (if the model includes the intercept, the residuals have zero mean) and  $\sum_i x_{ik} e_i = 0, k = 1, \dots, p$ .
- $\hat{\mathbf{y}}'\mathbf{e} = \sum_i \hat{y}_i e_i = 0$
- $\text{RSS}(\hat{\beta}) = \mathbf{e}'\mathbf{e} = \sum_i e_i^2$  (henceforth RSS)
- From  $\mathbf{y} = \hat{\mathbf{y}} + \mathbf{e}$  and the previous properties,  $\|\mathbf{y}\|^2 = \|\hat{\mathbf{y}}\|^2 + \|\mathbf{e}\|^2$
- $\bar{y} = N^{-1}\mathbf{1}'\mathbf{y} = \bar{\hat{y}}$  (average of predicted values)

# Estimation of $\sigma^2$

The OLS residual for the  $i$ -th observation is  $e_i = y_i - \hat{y}_i$ .

For estimating  $\sigma^2$  (the variance of the error term) we could use the variance of the residuals corrected for the number of degrees of freedom:

$$\hat{\sigma}^2 = \frac{\sum_i e_i^2}{N - p - 1}$$

$N - p - 1$  is known as the number of degrees of freedom.

$\hat{\sigma}$  is the *standard error of regression* (SER).

# Goodness of fit

Define

$TSS = \sum_i (y_i - \bar{y})^2$  (Total sum of squares - Deviance of  $y_i$ )

$ESS = \sum_i (\hat{y}_i - \bar{y})^2$  (Explained sum of squares - Deviance of  $\hat{y}_i$ )

$RSS = \sum_i e_i^2$  (Residual sum of squares - Deviance of  $e_i$ )

It is easy to prove that

$$TSS = ESS + RSS$$

A relative measure of g.o.f. is

$$R^2 = \frac{ESS}{TSS} = 1 - \frac{RSS}{TSS}$$

This measure suffers from a serious drawback when used for model selection: the inclusion of (possibly irrelevant) additional regressors always produces an increase in  $R^2$ .

Adjusted  $R^2$

$$\bar{R}^2 = 1 - \frac{RSS/(N - p - 1)}{TSS/(N - 1)} = 1 - \frac{N - 1}{N - p - 1}(1 - R^2)$$

# Properties of the OLS estimators

We are going to look at the statistical properties of the OLS estimators (and other related quantities, such as the fitted values and the residuals), hypothesizing that we can draw repeated training samples from the same population, keeping the  $\mathbf{X}$ 's fixed and drawing different  $\mathbf{Y}$ 's. A distribution of outcomes will arise.

Under the stated assumptions,

$$\mathbb{E}(\hat{\beta}|\mathbf{X}) = \beta, \quad \text{Var}(\hat{\beta}|\mathbf{X}) = \sigma^2(\mathbf{X}'\mathbf{X})^{-1}$$

Also,  $\mathbb{E}(\hat{\sigma}^2) = \sigma^2$ .

Moreover, the OLS estimators are Best Linear Unbiased (*Gauss Markov theorem*), also unconditionally:  $\mathbb{E}(\hat{\beta}_j) = \beta_j, \forall j$  and  $\text{Var}(\hat{\beta}_j) = \text{MSE}(\hat{\beta}_j)$  is the smallest among the class of all linear unbiased estimators.

Example:  $p = 1$  (simple regression model)

$$\text{Var}(\hat{\beta}_0|\mathbf{X}) = \sigma^2 \left[ \frac{1}{N} + \frac{\bar{x}^2}{\sum_i (x_i - \bar{x})^2} \right], \quad \text{Var}(\hat{\beta}_1|\mathbf{X}) = \frac{\sigma^2}{\sum_i (x_i - \bar{x})^2},$$

If we further assume

$$\varepsilon_i | \mathbf{x}_i \sim \text{NID}(0, \sigma^2),$$

- $y_i | \mathbf{x}_i \sim \text{N}(\beta_0 + \beta_1 x_i, \sigma^2)$
- The OLS estimators have a normal distribution:  $\hat{\beta} | \mathbf{X} \sim \text{N}(\beta, \sigma^2 (\mathbf{X}' \mathbf{X})^{-1})$ .
- $\hat{\beta}$  is the maximum likelihood estimator of the coefficients (we will say more about it).
- Letting  $\text{s.e.}(\hat{\beta}_j)$  denote the estimated standard error of the OLS estimator, i.e. the  $j$ -th element of  $\hat{\text{Var}}(\hat{\beta} | \mathbf{X}) = \hat{\sigma}^2 (\mathbf{X}' \mathbf{X})^{-1}$ ,

$$t_j = \frac{\hat{\beta}_j - \beta_j}{\text{s.e.}(\hat{\beta}_j)} \sim t_{N-p-1},$$

a Student's  $t$  r.v. with  $N - p - 1$  d.o.f.

This result is used to test hypotheses on a single coefficient and to produce interval estimates.

- Suppose we wish to test for the significance of subset of the coefficients. The relevant statistic for the null that  $J$  coefficients are all equal to zero is

$$F = \frac{(RSS_0 - RSS)/J}{RSS/(N - p - 1)} \sim F_{J, N-p-1}$$

where  $RSS_0$  is the residual sum of squares of the restricted model (i.e. it has  $J$  less explanatory variables), and  $F_{J, N-K}$  is Fisher's distribution with  $J$  and  $N - p - 1$  d.o.f.

As a special case, the test for  $H_0 : \beta_1 = \beta_2 = \dots = \beta_p = 0$  (all the regression coefficients are zero, except for the intercept) is:

$$F = \frac{ESS/p}{RSS/(N - p - 1)} = \frac{R^2/p}{(1 - R^2)/(N - p - 1)}$$

Example: clothing data. Regression of total sales (logs) on hours worked and store size (logs).

Call:

```
lm(formula = y ~ x1 + x2, data = clothing)
```

Residuals:

	Min	1Q	Median	3Q	Max
	-1.89464	-0.21835	0.01462	0.29581	1.63896

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	7.73152	0.21850	35.385	< 2e-16 ***
x1	0.88949	0.05791	15.359	< 2e-16 ***
x2	0.31488	0.04308	7.309	1.49e-12 ***

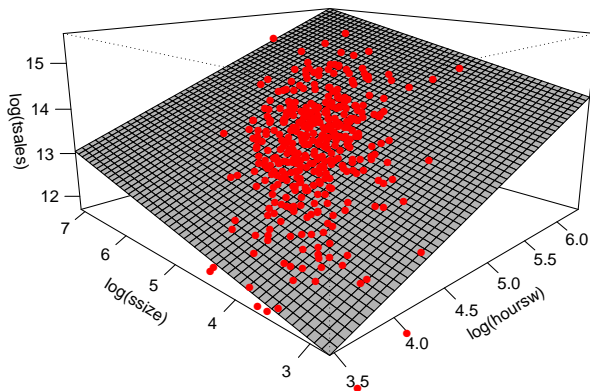
---

Residual standard error: 0.4311 on 397 degrees of freedom

Multiple R-squared: 0.6342, Adjusted R-squared: 0.6324

F-statistic: 344.2 on 2 and 397 DF, p-value: < 2.2e-16



**Figure: Fitted values**

# Example: house price regression in R

```
lm(formula = price ~ sqft + Age + Pool + Bedrooms + Pool + Fireplace
    Waterfront + DOM)
```

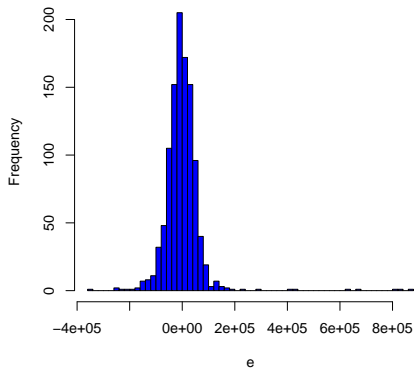
Coefficients:

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	4336.851	12084.856	0.359	0.720
sqft	97.862	3.485	28.082	< 2e-16 ***
Age	-694.671	141.292	-4.917	1.02e-06 ***
Pool	815.028	8926.466	0.091	0.927
Bedrooms	-20923.878	4567.168	-4.581	5.16e-06 ***
Fireplace	-97.130	5152.184	-0.019	0.985
Waterfront	63376.186	9389.874	6.749	2.43e-11 ***
DOM	-20.988	24.841	-0.845	0.398

Residual standard error: 76390 on 1072 degrees of freedom

Multiple R-squared: 0.6162, Adjusted R-squared: 0.6137

F-statistic: 245.9 on 7 and 1072 DF, p-value: < 2.2e-16

**Figure:** Histogram of residuals

# Properties of the OLS residuals

- Does  $e_i$  reflect the properties of  $\epsilon_i$ , and how?
- The  $i$ -th residual  $e_i$  has zero expectation (under the model's assumptions) and

$$\text{Var}(e_i|\mathbf{X}) = \sigma^2(1 - h_i),$$

where

$$h_i = \mathbf{x}_i'(\mathbf{X}'\mathbf{X})^{-1}\mathbf{x}_i$$

is the **leverage** of observation  $i$  (the  $i$ -th diagonal element of  $\mathbf{H}$ ).

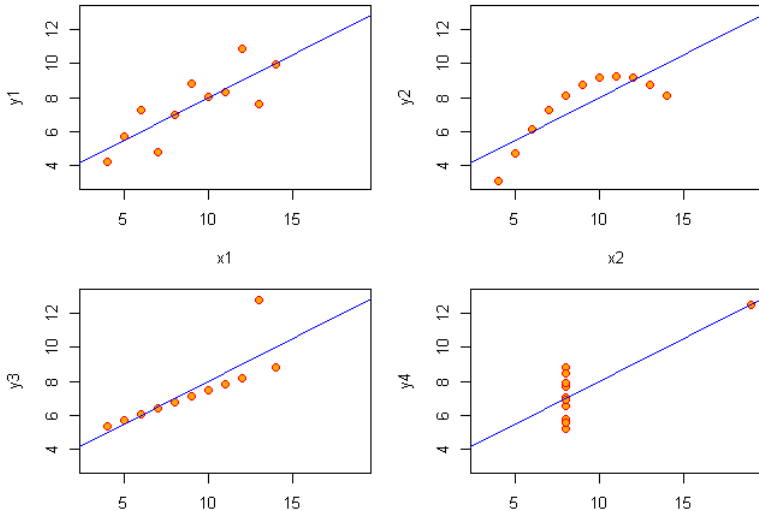
- The leverage,  $h_i$ , is an important diagnostic quantity. It is a measure of remoteness of the inputs for the  $i$ -th individual in the space of the inputs.
- The barplot of  $h_i$  versus  $i$  illustrates the influence of the  $i$ -th observation on the fit.
- It is possible to show that  $\frac{1}{N} \leq h_i \leq 1$ ,  $\sum_i h_i = p + 1$  and that the average of the  $h_i$ s is  $(p + 1)/N$ .
- The standardized residual is defined as

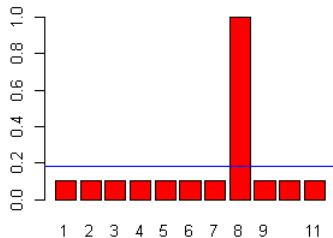
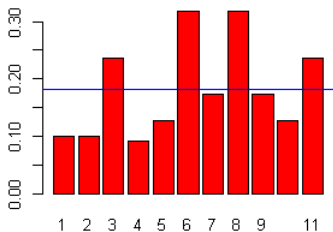
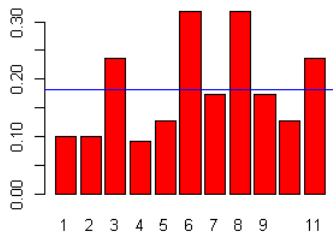
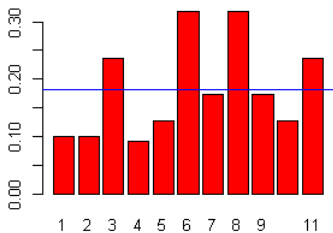
$$r_i = \frac{e_i}{\hat{\sigma}\sqrt{1 - h_i}},$$

and is used in regression diagnostics (normal probability plot, check for homoscedasticity).

**Figure:** Anscombe Quartet: 4 datasets giving the same set fitted values

## Anscombe's 4 Regression data sets



**Figure:** Anscombe Quartet: Leverage plot ( $h_i$  vs  $i$ )

# Specification issues: omitted variable bias, collinearity, etc.

Consider the regression model  $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \gamma\mathbf{w} + \varepsilon$ , where  $\mathbf{w}$  is a vector of  $N$  measurements on an added variable  $W$ . We can show that the OLS estimator of  $\gamma$  is

$$\hat{\gamma} = \frac{\mathbf{w}'\mathbf{M}\mathbf{y}}{\mathbf{w}'\mathbf{M}\mathbf{w}}, \quad \mathbf{M} = \mathbf{I} - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}',$$

and

$$\text{Var}(\hat{\gamma}) = \frac{\sigma^2}{\mathbf{w}'\mathbf{M}\mathbf{w}}.$$

These expressions have a nice interpretation:  $\hat{\gamma}$  results from

- Regress  $Y$  on  $X$  and obtain the residuals  $\mathbf{e}_{Y|X} = \mathbf{M}\mathbf{y}$ ;
- Regress  $W$  on  $X$  and obtain the residuals  $\mathbf{e}_{W|X} = \mathbf{M}\mathbf{w}$ ;
- Regress  $\mathbf{e}_{Y|X}$  on  $\mathbf{e}_{W|X}$  to obtain the LS estimate  $\hat{\gamma}$  (partial regr. coeff.).

# Multicollinearity

- $M\mathbf{w}$  are the residuals of the regression of  $W$  on  $X$ .
- $\mathbf{w}'M\mathbf{w}$  is the residual sum of squares from that regression.
- If  $X$  explains most of the variation in  $W$ , this is very small, and  $\text{Var}(\hat{\gamma})$  will be very high.
- We can conclude that the effect that  $W$  exerts on  $Y$  is poorly (i.e. very imprecisely) estimated.
- The variance inflation factor provides a measure of the increase in the variance of the OLS estimator due to multicollinearity: it is defined as  $1/(1 - R_w^2)$ , where  $R_w^2$  is the R-squared of the regression of  $W$  on  $X$  (write  $\mathbf{w}'M\mathbf{w} = \mathbf{w}'\mathbf{w}(1 - R_w^2)$ ).



# Omission of relevant variables

If  $W$  is omitted from the regression, the OLS estimator  $\hat{\beta}$  is biased,

$$E(\hat{\beta}|\mathbf{X}) = \beta + (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{w}\gamma,$$

but also more precise:

$$\text{Var}(\hat{\beta}|\mathbf{X}) = \sigma^2(\mathbf{X}'\mathbf{X})^{-1} \leq \sigma^2(\mathbf{X}'\mathbf{M}_W\mathbf{X})^{-1}; \quad \mathbf{M}_W = \mathbf{I} - \mathbf{w}(\mathbf{w}'\mathbf{w})^{-1}\mathbf{w}',$$

where  $\sigma^2(\mathbf{X}'\mathbf{M}_W\mathbf{X})^{-1}$  is the variance of the estimator of  $\beta$  when  $W$  is included in the model.

If  $W$  is irrelevant (true  $\gamma = 0$ ) there is only a cost associated to including it in the model.

In any case, if  $W$  is orthogonal to  $X$  the distribution of  $\hat{\beta}$  is not affected by  $W$ .

# Prediction of a new sample observation

We are interested in predicting  $Y$  for a new unit with values of  $X$  in the vector  $\mathbf{x}_o = (1, x_{o1}, \dots, x_{op})'$ . The optimal predictor under square loss is

$$\hat{y}_o = \hat{f}(\mathbf{x}_o) = \hat{\beta}_0 + \hat{\beta}_1 x_{o1} + \dots + \hat{\beta}_p x_{op} = \mathbf{x}_o' \hat{\beta}$$

Under the stated assumptions, the predictor is unbiased, i.e. the prediction error has zero mean,  $E(Y_o - \hat{y}_o) = 0$ , and

$$\text{Var}(Y_o - \hat{y}_o) = \sigma^2(1 + \mathbf{x}_o'(\mathbf{X}'\mathbf{X})^{-1}\mathbf{x}_o)$$

Note: the stated assumptions concern the selection of the set of predictors in  $X$  and the properties of  $\varepsilon$ .

Predictive accuracy depends on

- ① Ability to select all the relevant variable. Failure to do so will result in an inflated  $\sigma^2$  and  $E(\varepsilon|\mathbf{X}) \neq 0$ , so that  $\hat{\beta}$  is no longer optimal.
- ② How removed is the unit profile  $\mathbf{x}_o$  from the rest:  $\mathbf{x}_o'(\mathbf{X}'\mathbf{X})^{-1}\mathbf{x}_o$  measures the (Mahalanobis) distance of the values of the  $X_1, \dots, X_p$  for the  $o$ -th unit in the space of the  $X$ 's.

If the true model is

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \gamma\mathbf{w} + \varepsilon, \mathbb{E}(\varepsilon|\mathbf{X}, \mathbf{w}) = \mathbf{0}, \text{Var}(\varepsilon|\mathbf{X}, \mathbf{w}) = \sigma^2\mathbf{I},$$

the above predictor is biased and its mean square error is

$$\mathbb{E}[(Y_0 - \hat{f}(\mathbf{x}_o))^2|\mathbf{X}, \mathbf{w}, \mathbf{x}_o, w_o] = \sigma^2(1 + \mathbf{x}_o'(\mathbf{X}'\mathbf{X})^{-1}\mathbf{x}_o) + \text{Bias}^2,$$

with

$$\text{Bias} = (\boldsymbol{\beta} - \mathbb{E}(\hat{\boldsymbol{\beta}}))'\mathbf{x}_o + \gamma w_o = (w_o - \mathbf{w}'\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{x}_o)\gamma$$

The term in parenthesis is the error of prediction of  $w_o$  using the information in  $\mathbf{x}_o$ .