

# Computational Statistics

---

Daniel Balle 2018

# Multiple Linear Regression

In multiple regression we are given data  $(Y_i, \mathbf{x}_i)$  where we assume the response variable is a linear function of the predictors:

$$Y_i = \boldsymbol{\beta}^\top \mathbf{x}_i + \epsilon_i$$

or  $\mathbf{Y} = X\boldsymbol{\beta} + \boldsymbol{\epsilon}$  with  $\mathbf{Y} \in \mathbb{R}^n$  and  $X \in \mathbb{R}^{n \times p}$ . Errors  $\epsilon_i$  are usually assumed *iid* with  $\mathbb{E}[\epsilon_i] = 0$  and  $\text{Var}(\epsilon_i) = \sigma^2$ .

Our goal is to estimate the unknown parameters  $\boldsymbol{\beta} \in \mathbb{R}^p$ .

# Least Squares Estimator

The least squares estimator for a *linear* model  $\mathbf{Y} = X\boldsymbol{\beta} + \boldsymbol{\epsilon}$  is given by minimizing the least squares error:

$$\hat{\boldsymbol{\beta}} = \arg \min_{\boldsymbol{\beta}} \|\mathbf{Y} - X\boldsymbol{\beta}\|^2$$

*Remark.* This is solved theoretically by  $\hat{\boldsymbol{\beta}} = (X^\top X)^{-1} X^\top \mathbf{Y}$ . This estimator is *unbiased*, i.e.  $\mathbb{E}[\hat{\boldsymbol{\beta}}] = \boldsymbol{\beta}$  and  $\mathbb{E}[\hat{\mathbf{Y}}] = X\boldsymbol{\beta}$ .

# Least Squares Residuals

With the *least squares* estimator  $\hat{\beta}$  the residuals  $r_i = Y_i - \hat{\beta}^\top \mathbf{x}_i$  give an *unbiased* estimate of the errors  $\epsilon_i$ ,  $\mathbb{E}[\mathbf{r}] = \mathbf{0}$ . And

$$\hat{\sigma}^2 = (n - p)^{-1} \sum r_i^2$$

provides an *unbiased* estimate of the variance, i.e.  $\mathbb{E}[\hat{\sigma}^2] = \sigma^2$ .

# Least Squares Projection

Geometrically the *least squares* method performs an orthogonal projection  $\mathbf{Y} \mapsto \hat{\mathbf{Y}} = X\hat{\boldsymbol{\beta}}$  with projection matrix:

$$\hat{\mathbf{Y}} = P\mathbf{Y} \implies P = X(X^\top X)^{-1}X^\top$$

*Remark.* The residuals can then be expressed as  $\mathbf{r} = (I - P)\mathbf{Y}$ .

# Linear Regression in R

In `R` a *linear regression* model can be fitted as follows:

```
fitted <- lm(formula = LOGRUT ~ ., data = asphalt1)
```

# Bias-Variance Tradeoff

The bias-variance decomposition for any *supervised* learning task of  $y = f(x) + \epsilon$  is the expected generalization error:

$$\underbrace{\mathbb{E}[(f(x) - \hat{f}(x))^2]}_{\text{MSE}(x)} = \underbrace{(\mathbb{E}[\hat{f}(x)] - f(x))^2}_{\text{Bias}} + \underbrace{\mathbb{E}[\hat{f}(x)^2] - E[\hat{f}(x)]^2}_{\text{Var}(\hat{f}(x))}$$

*Remark.* Optimizing this trade-off is called *regularization*, and avoids the problem of *overfitting*.

# Kernel Density Estimation

Given realizations  $X_i \in \mathbb{R} \sim F$ , the nonparametric kernel density estimator  $\hat{f}$  of the unknown density function  $f = F'$  is

$$\hat{f}(x) = \frac{1}{nh} \sum K \left( \frac{x - X_i}{h} \right)$$

where  $K(\cdot)$  is a *kernel* function, usually symmetric around 0, and the tuning parameter  $h$  the *bandwidth*.

$$K(x) = K(-x) \quad K(x) \geq 0 \quad \int_{-\infty}^{\infty} K(x) dx = 1$$