

# MACHINE LEARNING

## LINEAR REGRESSION

Sebastian Engelke

MASTER IN BUSINESS ANALYTICS



**UNIVERSITÉ  
DE GENÈVE**

# Simple linear regression

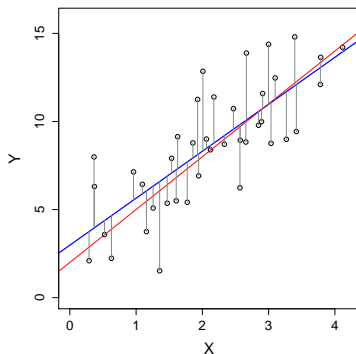
- ▶ Recall:  $Y = f(X) + \epsilon$
- ▶ **Simple linear model:**

$$f(X) = \beta_0 + \beta_1 X,$$

where  $\beta_0$  and  $\beta_1$  are model coefficients or parameters.

- ▶ For training data  $\{(x_i, y_i)\}_{i=1}^n$ , suppose we have an estimate  $(\hat{\beta}_0, \hat{\beta}_1)$  of  $(\beta_0, \beta_1)$ .
- ▶ **Fitted value:**  $\hat{y}_i = \hat{f}(x_i) = \hat{\beta}_0 + \hat{\beta}_1 x_i$
- ▶ **Residual:**  $e_i = y_i - \hat{y}_i$
- ▶ **Residual sum of squares (RSS)**

$$\begin{aligned}\text{RSS}(\hat{\beta}_0, \hat{\beta}_1) &= e_1^2 + e_2^2 + \cdots + e_n^2 \\ &= \sum_{i=1}^n (y_i - \hat{f}(x_i))^2 \\ &= \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)^2\end{aligned}$$



## Estimation of parameters by least squares

- ▶ To estimate  $(\beta_0, \beta_1)$ , we minimize the RSS:

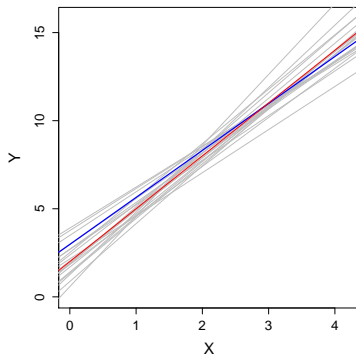
$$\begin{aligned}(\hat{\beta}_0, \hat{\beta}_1) &= \operatorname{argmin}_{(\beta_0, \beta_1) \in \mathbb{R}^2} \operatorname{RSS}(\beta_0, \beta_1) \\ &= \operatorname{argmin}_{(\beta_0, \beta_1) \in \mathbb{R}^2} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2\end{aligned}$$

- ▶ How do we solve this optimization problem?
- ▶ Numerically with computer algorithm (**gradient descent**, etc.)
- ▶ Analytically by differentiation (if possible):

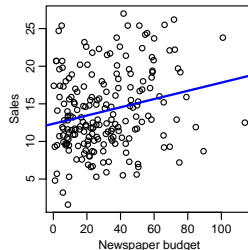
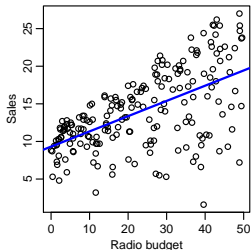
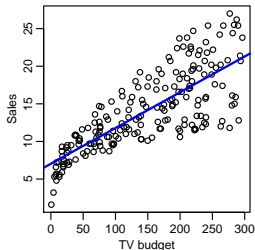
$$\begin{aligned}\hat{\beta}_1 &= \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} \\ \hat{\beta}_0 &= \bar{y} - \hat{\beta}_1 \bar{x},\end{aligned}$$

where  $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$  and  $\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$  are **sample means**.

- ▶ Note:  $(\hat{\beta}_0, \hat{\beta}_1)$  are **random**, since for different training data we obtain different estimates!



# Advertising data



Coefficient	Estimate	Std. Error	<i>t</i> value	<i>p</i> -value
Intercept ( $\beta_0$ )	7.032594	0.457843	15.36	$< 2e-16$
TV ( $\beta_1$ )	0.047537	0.002691	17.67	$< 2e-16$

Coefficient	Estimate	Std. Error	<i>t</i> value	<i>p</i> -value
Intercept ( $\beta_0$ )	9.31164	0.56290	16.542	$< 2e-16$
radio ( $\beta_1$ )	0.20250	0.02041	9.921	$< 2e-16$

Coefficient	Estimate	Std. Error	<i>t</i> value	<i>p</i> -value
Intercept ( $\beta_0$ )	12.35141	0.62142	19.88	$< 2e-16$
newspaper ( $\beta_1$ )	0.05469	0.01658	3.30	0.00115

## Multiple linear regression

- ▶ We consider  $p$  predictors  $X = (X_1, \dots, X_p)$  for the response  $Y$ . The model is

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_p X_p + \epsilon.$$

- ▶ For the **Advertising** data set:

$$\text{sales} = \beta_0 + \beta_1 \times \text{TV} + \beta_2 \times \text{radio} + \beta_3 \times \text{newspaper} + \epsilon.$$

- ▶ The training data  $\{(x_i, y_i)\}_{i=1}^n$  with  $x_i = (x_{i1}, \dots, x_{ip})^\top$  can be written as an  $n \times p$  matrix

$$\mathbf{X} = \begin{bmatrix} x_{11} & x_{12} & \dots & x_{1p} \\ x_{21} & x_{22} & \dots & x_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \dots & x_{np} \end{bmatrix}$$

- ▶ For notational convenience, we sometimes use the **dot product** between  $x, y \in \mathbb{R}^p$

$$x^\top y = \sum_{i=1}^p x_i y_i = x_1 y_1 + \dots + x_p y_p.$$

## Estimation and prediction in multiple linear regression

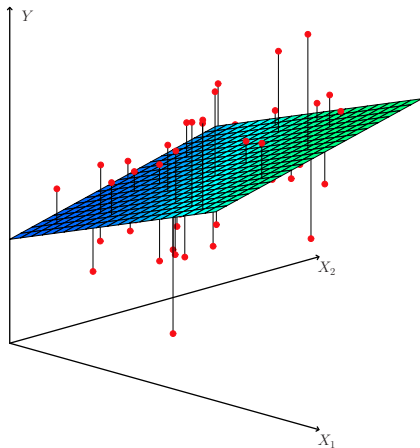
- As before, we obtain parameter estimates  $\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_p$  as those parameters that minimize the **RSS**

$$\begin{aligned}\text{RSS}(\beta_0, \dots, \beta_p) &= \sum_{i=1}^n (y_i - \beta_0 - \mathbf{x}_i^\top \boldsymbol{\beta})^2 \\ &= \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_{i1} - \dots - \beta_p x_{ip})^2\end{aligned}$$

- We make **predictions** at a new point  $\mathbf{x}_0 = (x_{01}, \dots, x_{0p})$  by

$$\begin{aligned}\hat{y}_0 &= \hat{\beta}_0 + \mathbf{x}_0^\top \hat{\boldsymbol{\beta}} \\ &= \hat{\beta}_0 + \hat{\beta}_1 x_{01} + \dots + \hat{\beta}_p x_{0p},\end{aligned}$$

where  $\hat{\boldsymbol{\beta}} = (\hat{\beta}_1, \dots, \hat{\beta}_p)$ .



## Advertising data

Coefficient	Estimate	Std. Error	<i>t</i> value	<i>p</i> -value
Intercept	2.939	0.3119	9.42	$< 2e-16$
TV	0.046	0.0014	32.81	$< 2e-16$
radio	0.189	0.0086	21.89	$< 2e-16$
newspaper	-0.001	0.0059	-0.18	0.8599

Correlations:

	TV	radio	newspaper
TV	1	0.0548	0.0567
radio		1	0.354
newspaper			1

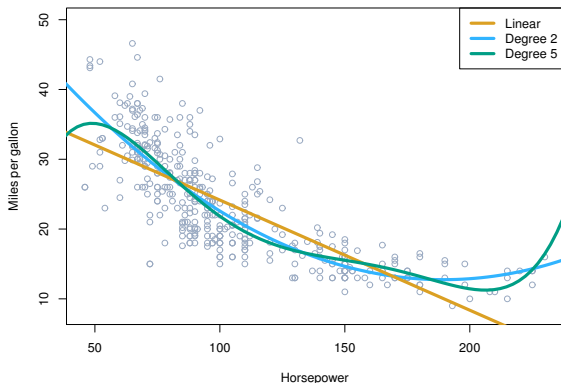
## Example: the Auto data set

- ▶ The **Auto** data set contains measurements for  $n = 397$  cars, such as **mpg** (miles per gallon), **horsepower**, **year**, etc.
- ▶ We can fit the model

$$\text{mpg} = \beta_0 + \beta_1 \times \text{horsepower} + \beta_2 \times \text{horsepower}^2 + \epsilon,$$

this is still a linear model (in  $X_1 = \text{horsepower}$  and  $X_2 = \text{horsepower}^2$ ).

- ▶ In this case, both  $\beta_1$  and  $\beta_2$  are significant.
- ▶ Even higher polynomials can be fitted.





## Beyond linearity: linear basis functions

- ▶ Most often the true function  $f(X)$  is not linear in the predictors  $X_1, \dots, X_p$ ! But we can adjust the linear model.
- ▶ Instead of the  $X = (X_1, \dots, X_p)$  we consider transformation of them. Let  $h_m(X) : \mathbb{R}^p \rightarrow \mathbb{R}$  be the  $m$ th transformation,  $m = 1, \dots, M$ . The model is

$$f(X) = \sum_{m=1}^M \beta_m h_m(X).$$

- ▶ This model is **linear** in the new predictors, the **basis functions**  $h_1(X), \dots, h_M(X)$ ; estimation and prediction can be applied as before!

Possible choices for  $h_m$  are:

- ▶  $h_m(X) = X_m$ ,  $m = 1, \dots, p$ , essentially recovers the original linear model.
- ▶  $h_m(X) = X_j^2$  or  $h_m(X) = X_j X_k$  allows to achieve higher order polynomials.
- ▶  $h_m(X) = \log(X_j)$ ,  $\sqrt{X_j}$  covers other non-linear transformation.
- ▶ etc.