

# Linear regression

- We use a **linear model** to approximate an unknown function which depends on several variables, to predict a continuous target variable.
- Objective: to find the optimal parameters of such a model.

$$\hat{y}_i = \beta_0 + \sum_{j=1}^p \beta_j x_{ij}$$

predicted value of the target variable for the  $i^{\text{th}}$  data point

intercept

weight of the  $j^{\text{th}}$  feature in the model

total number of features in the model.

value of the  $j^{\text{th}}$  feature in the  $i^{\text{th}}$  data point

$$\hat{y}_i = \beta_0 + \sum_{j=1}^p \beta_j x_{ij}$$

$$\begin{pmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \hat{y}_3 \\ \vdots \\ \hat{y}_n \end{pmatrix} = \begin{pmatrix} 1 & x_{11} & x_{12} & \dots & x_{1p} \\ 1 & x_{21} & x_{22} & \dots & x_{2p} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & x_{n2} & \dots & x_{np} \end{pmatrix} \begin{pmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \vdots \\ \beta_p \end{pmatrix}$$

$\hat{Y} = X \beta$

→  $\hat{Y} = X\beta$  in matrix notation.

Mean-squared error (MSE), also called the residual sum of squares:

$$\text{MSE}(\beta) = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

→ true value of the target variable

→ predicted value of the target

$$= \frac{1}{n} \sum_{i=1}^n \left( y_i - \left\{ \beta_0 + \sum_{j=1}^p \beta_j x_{ij} \right\} \right)^2$$

To find the optimal set of parameters,

$$\frac{\partial \text{MSE}(\beta)}{\partial \beta_j} = 0.$$

for any extremum

To ensure a minimum:  $\frac{\partial^2 \text{MSE}}{\partial \beta_j^2} > 0.$

$$SSE = n \times MSE$$

$$SSE(\beta) = (Y - X\beta)^T (Y - X\beta)$$

$$= (Y^T - (X\beta)^T)(Y - X\beta)$$

$$= (Y^T - \beta^T X^T)(Y - X\beta)$$

$$SSE(\beta) = Y^T Y - Y^T X \beta - \beta^T X^T Y + \beta^T X^T X \beta$$

$$\frac{\partial (SSE)}{\partial \beta} = 0 \Rightarrow \frac{\partial (SSE)}{\partial \beta} = -2X^T Y + 2X^T X \beta$$

$$\frac{\partial^2 (SSE)}{\partial \beta^2} = 2X^T X$$

positive semi-definite  
matrix

ensures a minimum in most cases

For the optimal  $\beta$ , denoted as  $\hat{\beta}$

$$-2X^T Y + 2X^T X \hat{\beta} = 0$$

$$2X^T X \hat{\beta} = 2X^T Y$$

$$\hat{\beta} = (X^T X)^{-1} X^T Y$$

expression for the optimal set of parameters  
"best-fit linear model".

If  $X$  or  $Y$  have complex entries,

$$\hat{\beta} = (X^H X)^{-1} X^H Y$$

where

$$X^H = \overline{X^T}$$

complex conjugate

Hermitian transpose: complex conjugate transpose.

The matrix  $X^+ = (X^H X)^{-1} X^H$  is referred to as a pseudo-inverse or a Moore-Penrose inverse. of the matrix  $X$  and is the generalization of the concept of an inverse to a rectangular matrix.

$$AX = b \rightarrow X = A \backslash b$$
$$X = (A^H A)^{-1} A^H b$$

Coefficient of determination ( $R^2$ ):

It is a measure of the goodness of fit of a model.

$$R^2 = 1 - \frac{SS_{\text{res}}}{SS_{\text{tot}}}$$

residual sum of squares

total sum of squares

$$SS_{res} = \sum_{i=1}^n (y_i - \hat{y}_i)^2 \equiv SSE$$

$\swarrow$  true value       $\searrow$  predicted value

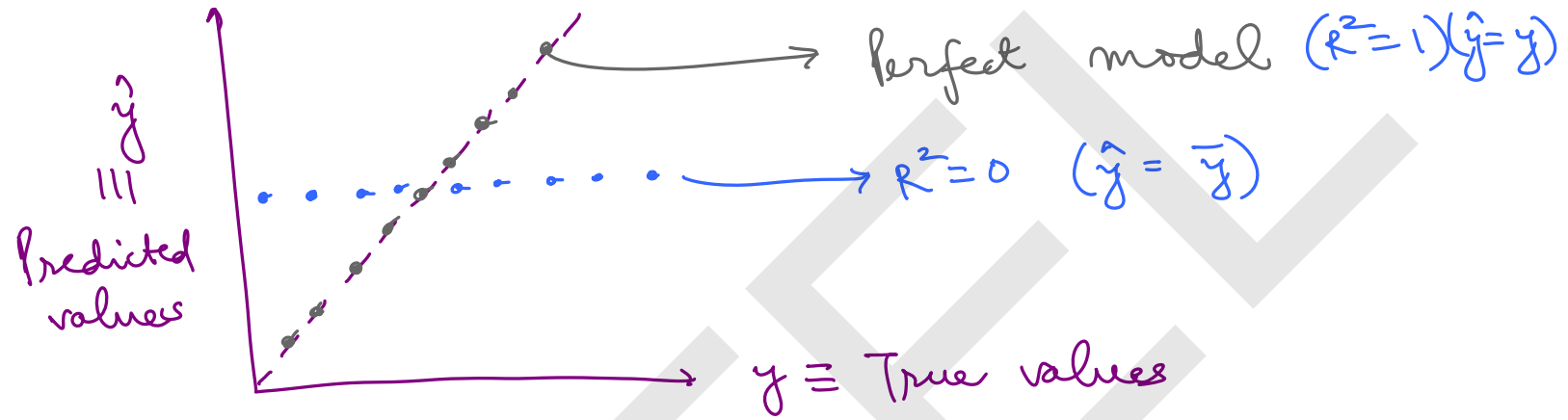
$$SS_{tot} = \sum_{i=1}^n (y_i - \bar{y})^2$$

$\searrow$  mean of the true values

$$R^2 = 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2}$$

\* For a perfect model:  $\hat{y}_i = y_i \Rightarrow R^2 = 1 - 0 = \boxed{1}$

Parity plot



\* For imperfect models:  $R^2 < 1$

eg:  $R^2 = 0$

When  $\hat{y}_i = \bar{y}$ , then  $R^2 = 1 - 1 = \boxed{0}$

→ the model always predicts the mean value.

Note that  $R^2$  can be less than zero for very inferior model.