

Linear regression: ordinary least squares

Regression can be defined as a statistical technique for investigating and modeling statistical relation between the *mean value* of a dependent variable(s) \mathbf{y} (**response, target**) and corresponding values of independent variable(s) \mathbf{x} – **predictors, covariates**:

$$E\{\mathbf{y}|\mathbf{x}\} = f(\mathbf{x}) \quad (1)$$

The regression equation *is only an approximation* to the true functional relationship between the variables of interest. Thus, **regression does not mean causation**.

Generally, **regression equations are valid only over the region of the regressor variables contained** in the observed data.

Why classical LS linear regression:

- Central topic in machine learning – can be generalized to other popular ML methods – classification problems and logistic regression, neural networks;
- Has a closed-form expression for solution;
- Easy to introduce basic concepts like **bias-variance tradeoff**, **cross-validation**, **resampling**, **regularization** techniques and many other ML topics.

1. Multiple linear regression formulation

Assume the following form of functional relation between the variables:

$$y_i = \mathbf{x}_i^T \boldsymbol{\beta} + \varepsilon_i, \quad i = 1, \dots, n. \quad (2)$$

With $\varepsilon_i \sim N(0, \sigma^2)$ it follows that:

- $E\{y_i|x\} = \mathbf{x}_i^T \boldsymbol{\beta} \Rightarrow$ **regression model is a line (or hyperplane in general)**;
- $\text{Var}\{y|x\} = \sigma_{y|x}^2 = \text{Var}(\beta_0 + \beta_1 x + \varepsilon) = \sigma^2$.

If to concatenate in columns all the observations i , we can re-formulate (2) in the matrix form:

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon} \quad (3)$$

where \mathbf{X} is the **design matrix** n – observations \times p – features:

$$\mathbf{X} = \begin{pmatrix} 1 & x_{12} & \dots & x_{1p} \\ 1 & x_{22} & \dots & x_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n2} & \dots & x_{np} \end{pmatrix} = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_p) \quad (4)$$

Often $\mathbf{x}_1 = \mathbf{1}$ so that β_0 is a **free term**.

$\boldsymbol{\beta}$ is the $p \times 1$ vector of model parameters;

\mathbf{y} is the $1 \times n$ vector of target observations.

In classical linear regression model apart from relation (3) also imposed the following assumptions on stochastic structure of the model – **5 basic assumptions of classical model**:

- Model is sufficient i.e. $E\{y|x\}$ is **linear, no grouping**. (*Solution*: transform data, \ln or x/y^{-1});
- $E\{\varepsilon_i\} = 0$ + random scatter is symmetric (*Solution*: data transformation);
- $E\{\varepsilon_i^2\} = \sigma^2$ – const error variance;
- $E\{\varepsilon_i \varepsilon_j\} = 0, \quad \forall i \neq j$ – uncorrelated errors;
- $\text{rk}(\mathbf{X}) = p < n$ (rank is the maximum number of linearly independent column vectors);

Also:

- Linear regression is **sensitive to outliers** (*Solution*: transform data or remove them);
- $\varepsilon \sim N(0, \sigma^2)$ – explicit assumption on the error form **for hypothesis testing and interval estimation**.

2. Least squares estimation for multiple linear regression

Under the described above assumptions, the main method of model parameters (β) estimation is the **Ordinary Least Squares (OLS)**:

$$\hat{\beta}_{\text{OLS}} = \arg \min_{\beta} \sum_{i=1}^n (\mathbf{y}_i - \mathbf{x}_i^T \beta)^2 \quad (5)$$

Which has a general solution in matrix form as:

$$\hat{\beta}_{\text{OLS}} = \underbrace{(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}}_{\text{MP pseudoinverse}} \quad (6)$$

As a result of this minimization problem, we get **fitted values** $\hat{y}_i = x_i^T \hat{\beta}$ and **residuals** $e_i = y_i - \hat{y}_i$.

2.1 LSE derivation and solution

We wish to find a vector of OLS model parameters that minimizes:

$$S(\beta) = \sum_{i=1}^n \varepsilon_i^2 = \varepsilon^T \varepsilon = (\mathbf{y} - \mathbf{X}\beta)^T (\mathbf{y} - \mathbf{X}\beta) \quad (7)$$

Open the brackets:

$$\begin{aligned} S(\beta) &= \mathbf{y}^T \mathbf{y} - \beta^T \mathbf{X}^T \mathbf{y} - \mathbf{y}^T \mathbf{X} \beta + \beta^T \mathbf{X}^T \mathbf{X} \beta \\ &= \mathbf{y}^T \mathbf{y} - 2\beta^T \mathbf{X}^T \mathbf{y} + \beta^T \mathbf{X}^T \mathbf{X} \beta \end{aligned} \quad (8)$$

$\beta^T \mathbf{X}^T \mathbf{y}$ is a scalar which transpose $(\beta^T \mathbf{X}^T \mathbf{y})^T = \beta \mathbf{X} \mathbf{y}^T$ is the same scalar!

The optimal values of parameters $\hat{\beta}$ should minimize $S(\beta)$, thus:

$$\left. \frac{\partial S}{\partial \beta} \right|_{\hat{\beta}} = (\text{simple lin. alg. from above}) = -2\mathbf{X}^T \mathbf{y} + 2\mathbf{X}^T \mathbf{X} \hat{\beta} = \mathbf{0} \quad (9)$$

Simplify it to the **LS normal equation**:

$$\boxed{\mathbf{X}^T \mathbf{X} \hat{\beta} = \mathbf{X}^T \mathbf{y}} \quad (10)$$

Normal equation in details:

$$\begin{bmatrix} \sum_{i=1}^n x_{i1} & \sum_{i=1}^n x_{i1}^2 & \sum_{i=1}^n x_{i1}x_{i2} & \cdots & \sum_{i=1}^n x_{i1}x_{ip} \\ \sum_{i=1}^n x_{i2} & \sum_{i=1}^n x_{i2}^2 & \sum_{i=1}^n x_{i2}x_{i3} & \cdots & \sum_{i=1}^n x_{i2}x_{ip} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \sum_{i=1}^n x_{ip} & \sum_{i=1}^n x_{ip}x_{i1} & \sum_{i=1}^n x_{ip}x_{i2} & \cdots & \sum_{i=1}^n x_{ip}^2 \end{bmatrix} \begin{bmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \\ \vdots \\ \hat{\beta}_p \end{bmatrix} = \begin{bmatrix} \sum_{i=1}^n y_i \\ \sum_{i=1}^n x_{i1}y_i \\ \vdots \\ \sum_{i=1}^n x_{ip}y_i \end{bmatrix} \quad (11)$$

Solution of the normal equation is (6), provided that $(\mathbf{X}^T \mathbf{X})^{-1}$ exist which is **always true if regressors are linearly independent** – no column of \mathbf{X} is a linear combination of other columns.

2.2 Hat (projection) matrix H

The **fitted values** from the LSE are then:

$$\boxed{\hat{\mathbf{y}} = \mathbf{X}\hat{\beta} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} = \mathbf{H}\mathbf{y}} \quad (12)$$

where

$$\mathbf{H}_{n \times n} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \quad (13)$$

Is $n \times n$ squared **hat matrix**. It maps the vector of observed target values onto a vector of fitted values.

2.3 Residuals e

The difference between the observed and corresponding fitted values is the residuals: $\mathbf{e} = \mathbf{y} - \hat{\mathbf{y}}$.

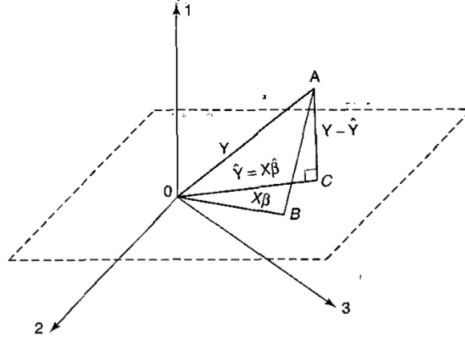
We can express them using (12) also as:

$$\mathbf{e} = \mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}} = (\mathbf{I} - \mathbf{H})\mathbf{y} \quad (14)$$

Residuals play an important role in investigating *model adequacy* and in detecting departures from the underlying assumptions.

2.4 Geometrical interpretation of the least squares

In the following figure the *estimation space* is $p = 2$ which is a subspace of a 3D *sample space* $n = 3$.



We may think of the vector of observations $\mathbf{y} = [y_1, y_2, \dots, y_n]$ as defining the vector in the sample space from the origin to the point A .

The \mathbf{X} matrix consists of p ($n \times 1$) vectors defined from the origin in the sample space. These vectors form a p -dimensional estimation space. In the figure we have only two vectors ($p = 2$). We may represent any point in the estimation space as a linear combination of vectors $\mathbf{1}, \mathbf{x}_1, \dots, \mathbf{x}_p$ so that any point on the estimation space is in the form of $\mathbf{X}\boldsymbol{\beta}$ determining the point B .

What we do in OLS is minimization of the distance between OB and OA . The solution in the estimation space satisfying this is the vector OC ($\hat{\mathbf{y}} = \mathbf{X}\hat{\boldsymbol{\beta}}$) which is the projection of the vector \mathbf{y} onto the sample space $\mathbf{1}, \mathbf{x}_1, \dots, \mathbf{x}_p$.

We see that the vector of residuals $\mathbf{e} = \mathbf{y} - \hat{\mathbf{y}} = \mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}$ is orthogonal to the sample space \mathbf{X} . Therefore, the following is true:

$$\mathbf{X}^T \mathbf{e} = \mathbf{X}^T (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}) = 0, \quad (15)$$

which we recognize as the LS normal equation:

$$\mathbf{X}^T \mathbf{X} \hat{\boldsymbol{\beta}} = \mathbf{X}^T \mathbf{y}. \quad (16)$$

Normal equation in details:

$$\begin{bmatrix} n & \sum_{i=1}^n x_{i1} & \sum_{i=1}^n x_{i2} & \cdots & \sum_{i=1}^n x_{ip} \\ \sum_{i=1}^n x_{i1} & \sum_{i=1}^n x_{i1}^2 & \sum_{i=1}^n x_{i1}x_{i2} & \cdots & \sum_{i=1}^n x_{i1}x_{ip} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \sum_{i=1}^n x_{ip} & \sum_{i=1}^n x_{ip}x_{i1} & \sum_{i=1}^n x_{ip}x_{i2} & \cdots & \sum_{i=1}^n x_{ip}^2 \end{bmatrix} \begin{bmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \\ \vdots \\ \hat{\beta}_p \end{bmatrix} = \begin{bmatrix} \sum_{i=1}^n y_i \\ \sum_{i=1}^n x_{i1}y_i \\ \vdots \\ \sum_{i=1}^n x_{ip}y_i \end{bmatrix} \quad (17)$$

Solution of the normal equation is (6), provided that $(\mathbf{X}^T \mathbf{X})^{-1}$ exist which is *always true if regressors are linearly independent* – no column of \mathbf{X} is a linear combination of other columns.

3. Simple linear regression model

In general, OLS optimal model parameters:

$$\hat{\boldsymbol{\beta}}_{\text{OLS}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \quad (18)$$

Specify this solution for the simple linear regression of a single predictor with a single free term:

$$y_i = x_i\beta_1 + \beta_0 + \varepsilon_i, \quad i = 1, \dots, n. \quad (19)$$

3.1 Derivation of LS estimators $\hat{\beta}_0, \hat{\beta}_1$

As a first step, introduce normalizing factors $1/n$ to the both factors of (18) for the convenience:

$$\hat{\beta} = \left(\frac{1}{n} \mathbf{X}^T \mathbf{X} \right)^{-1} \left(\frac{1}{n} \mathbf{X}^T \mathbf{y} \right) \quad (20)$$

Then let's look at two factors separately:

$$\frac{1}{n} \mathbf{X}^T \mathbf{X} = \frac{1}{n} \begin{bmatrix} n & \sum_i x_i \\ \sum_i x_i & \sum_i x_i^2 \end{bmatrix} = \begin{bmatrix} 1 & \bar{x} \\ \bar{x} & \bar{x}^2 \end{bmatrix} \quad (21)$$

Now, we need to take the inverse:

$$\left(\frac{1}{n} \mathbf{X}^T \mathbf{X} \right)^{-1} = \frac{1}{\bar{x}^2 - \bar{x}^2} \begin{bmatrix} \bar{x}^2 & -\bar{x} \\ -\bar{x} & 1 \end{bmatrix} = \frac{1}{\sigma_x^2} \begin{bmatrix} \bar{x}^2 & -\bar{x} \\ -\bar{x} & 1 \end{bmatrix} \quad (22)$$

Another factor:

$$\frac{1}{n} \mathbf{X}^T \mathbf{y} = \frac{1}{n} \begin{bmatrix} 1 & 1 & \dots & 1 \\ x_1 & x_2 & \dots & x_n \end{bmatrix} \cdot \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \frac{1}{n} \begin{bmatrix} \sum_i y_i \\ \sum_i x_i y_i \end{bmatrix} = \begin{bmatrix} \bar{y} \\ \bar{xy} \end{bmatrix} \quad (23)$$

Now, multiply the two factors:

$$\begin{aligned} \frac{1}{s_X^2} \begin{bmatrix} \bar{x}^2 & -\bar{x} \\ -\bar{x} & 1 \end{bmatrix} \begin{bmatrix} \bar{y} \\ \bar{xy} \end{bmatrix} &= \frac{1}{\sigma_x^2} \begin{bmatrix} \bar{x}^2 \bar{y} - \bar{x} \bar{xy} \\ -\bar{x} \bar{y} + \bar{xy} \end{bmatrix} = \frac{1}{\sigma_x^2} \begin{bmatrix} (s_X^2 + \bar{x}^2) \bar{y} - \bar{x} (c_{XY} + \bar{x} \bar{y}) \\ c_{XY} \end{bmatrix} \\ &= \frac{1}{\sigma_x^2} \begin{bmatrix} \sigma_x^2 \bar{y} + \bar{x}^2 \bar{y} - \bar{x} c_{XY} - \bar{x}^2 \bar{y} \\ c_{XY} \end{bmatrix} = \begin{bmatrix} \bar{y} - \frac{C_{XY}}{\sigma_x^2} \bar{x} \\ \frac{C_{XY}}{\sigma_x^2} \end{bmatrix} \end{aligned} \quad (24)$$

Finally, the model parameters of the simple regression ($i = 1, \dots, n$):

$$\hat{\beta}_1 = \frac{C_{xy}}{SS_x} = \frac{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2} \quad - \text{ slope} \quad (25)$$

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} \quad - \text{ intercept} \quad (26)$$

Conclusion: regression is a linear relationship summary (that's what the correlation measures).

Gauss-Markov theorem. The OLS estimator of β $\hat{\beta}_{OLS} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$ where errors are **uncorrelated, have equal variances and zero expectations** is **BLUE** i.e. $\hat{\beta}$ has the smallest variance among the class of all unbiased estimators that are the linear combination of data. Errors do not need to be normally distributed nor iid.

3.4 Estimation of σ^2

Ideally, we want to have a prior information about σ^2 . If this is not possible, the **Residual Sum of Squares** is used

$$RSS = \sum_{i=1}^n e_i^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2 \quad (27)$$

The residual sum of squares has $n - 2$ degrees of freedom, because *2 degrees of freedom are already associated with estimates of $\hat{\beta}_0$ and $\hat{\beta}_1$ involved in obtaining \hat{y} .*

It can be shown that the expected value of RSS is $E(RSS) = (n - p)\sigma^2$ (the degrees of freedom comes from the rank of $\mathbf{I} - \mathbf{H}$ which is trace $tr(\mathbf{I} - \mathbf{H}) = n - p$). Thus, an unbiased estimator of σ^2 is:

$$\hat{\sigma}^2 = \frac{RSS}{n - p} = MS_{\text{Res}} \quad - \text{Residual Mean Square.} \quad (28)$$

The quantity $\sqrt{\hat{\sigma}^2}$ is sometimes called **Standard Error of Regression**.

Since $\hat{\sigma}^2$ depend on the RSS , any violation of the assumptions on the model errors or any misspecification of the model form may seriously damage the usefulness of its estimate given in (36) \Rightarrow we say that it is **a model dependent estimate** of σ^2 .

4. Coefficient of determination, R^2

Let $SS_R = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2$ and $SS_T = \sum_{i=1}^n (y_i - \bar{y})^2$. Then the quantity

$$R^2 = \frac{SS_R}{SS_T} = 1 - \frac{RSS}{SS_T} \quad (29)$$

is the **coefficient of determination** i.e. *proportion of variation explained by the regressor x* . Since $0 \leq RSS \leq SS_T$, it follows that:

$$0 \leq R^2 \leq 1 \quad (30)$$

Values of R^2 *close to 1 imply that most of the variability in y is explained by the regression model* e.g. if $R^2 = 0.9018$ that is, 90.18% of the variability in strength is accounted for by the regression model. *R^2 should be used with caution*, since it is always possible to make R^2 large by adding enough terms to the model i.e. it is **sensitive to overfitting**.

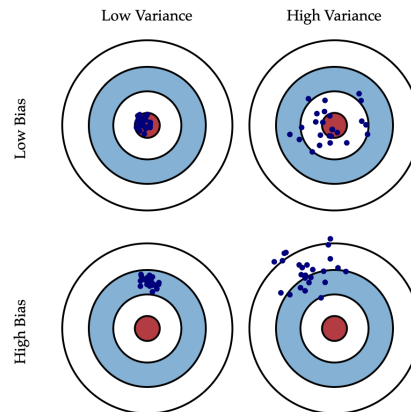
R^2 does not measure the appropriateness of the linear model, for R^2 will often be large even though y and x are *nonlinearly related*. Remember that *large R^2 does not necessarily imply that the regression model will be an accurate predictor*.

5. Bias-variance trade-off

Conceptual definition

- **Error due to bias:** Due to the randomness of the dataset, the underlying model will make predictions covering a certain range. Difference of the average outcome of our model and the correct values we are aiming to predict is a model bias – how far are the model predictions in general.
- **Error due to variance:** For a given data we want to predict, how wide is the spread of the model predictions w.r.t. these data.

Graphical definition



Mathematical definition

Say that exist some true *deterministic* function that we want to approximate with our ML model

$$f = f(\mathbf{x})$$

Although, we cannot observe outcomes of this unknown function directly but with some additive stochastic component, say being centered $E\{\varepsilon\} = 0$:

$$t = f + \varepsilon$$

Consider a dataset:

$$\mathbf{D} = \{(t_i, \mathbf{x}_i), i = 1 \dots N\}$$

Having this dataset \mathbf{D} , we train a parametric model g (e.g. neural network or linear OLS regression) to approximate the function of interest f based on the 'noised' observations t :

$$y_i = g(\mathbf{x}, \mathbf{w})$$

Assume that the measure of quality of our model is MSE:

$$\text{MSE} = E \left\{ (t_i - y_i)^2 \right\}$$

Do some math over this expectation:

$$\begin{aligned} E \left\{ (t_i - y_i)^2 \right\} &= \text{add and subtract hypothetical true values then associate} = \\ &= E \left\{ ((t_i - f_i) + (f_i - y_i))^2 \right\} = E \left\{ (t_i - f_i)^2 \right\} + E \left\{ (f_i - y_i)^2 \right\} + 2E \left\{ (f_i - y_i)(t_i - f_i) \right\} \\ &= E \left\{ \varepsilon^2 \right\} + E \left\{ (f_i - y_i)^2 \right\} + 2(E \{f_i t_i\} - E \{f_i^2\} - E \{y_i t_i\} + E \{y_i f_i\}) = E \left\{ \varepsilon^2 \right\} + E \left\{ (f_i - y_i)^2 \right\} \end{aligned}$$

The last term goes to zero because:

$$\begin{aligned} - E \{f_i t_i\} &= f_i^2 \Leftarrow E \{f_i(f_i + \varepsilon)\} = E \{f_i^2 + f_i \varepsilon\} = f_i \text{ is deterministic} = E \{f_i^2\} + 0 = f_i^2; \\ - E \{y_i t_i\} &= E \{y_i(f_i + \varepsilon)\} = E \{y_i f_i + y_i \varepsilon\} = E \{y_i f_i\} + E \{y_i\} E \{\varepsilon\} = E \{y_i f_i\} + 0. \end{aligned}$$

So now we have

$$E \left\{ (t_i - y_i)^2 \right\} = \sigma_\varepsilon^2 + E \left\{ (f_i - y_i)^2 \right\}$$

where we can apply the same trick as before to the second term but adding and subtracting the mean of the model outcomes:

$$\begin{aligned} E \left\{ (f_i - y_i)^2 \right\} &= E \left\{ (f_i - \bar{y}_i + \bar{y}_i - y_i)^2 \right\} = \\ &= E \left\{ (f_i - \bar{y}_i)^2 \right\} + E \left\{ (\bar{y}_i - y_i)^2 \right\} + 2E \left\{ (\bar{y}_i - y_i)(f_i - \bar{y}_i) \right\} = \\ &= E \left\{ (f_i - \bar{y}_i)^2 \right\} + E \left\{ (\bar{y}_i - y_i)^2 \right\} + 2(E \{f_i \bar{y}_i\} - E \{\bar{y}_i^2\} - E \{y_i f_i\} + E \{y_i \bar{y}_i\}) \end{aligned}$$

The last term also cancels to zero because:

$$\begin{aligned} - E \{f_i \bar{y}_i\} &= f_i \text{ is deterministic} = f_i E \{\bar{y}_i\} = f_i \bar{y}_i; \\ - E \{\bar{y}_i^2\} &= \text{mean is deterministic} = \bar{y}_i^2; \\ - E \{y_i f_i\} &= f_i \bar{y}_i; \\ - E \{y_i \bar{y}_i\} &= \bar{y}_i^2; \end{aligned}$$

Then, if to sum up we have the following:

$$\boxed{\text{MSE} = E \left\{ (t_i - y_i)^2 \right\} = E \left\{ (f_i - \bar{y}_i)^2 \right\} + E \left\{ (\bar{y}_i - y_i)^2 \right\} + \sigma_\varepsilon^2 = \text{Bias} + \text{Var}_y + \sigma_\varepsilon^2}$$

Thus to minimize the MSE of the model, we have to minimize both the bias of the model as well its variance although it is not trivial!

The following is commonly true about the components of the bias-variance equation:

- **A high bias model** has a relatively large error, mostly due to wrong assumptions about the data features;
- **A high variance model** is quite sensitive to small variations in the train data;
- **Irreducible error** is due to the intrinsic *noisiness of the data itself*.

6. Gradient descent optimization methods

TBA