

Linear Regression: Ordinary Least Square

Xiang Zhou

School of Data Science
Department of Mathematics
City University of Hong Kong



Ordinary Linear Regression

Review of linear regression (univariate and multivariate)

- Least-square: is usually credited to Carl Friedrich Gauss (1795), but it was first published by Adrien-Marie Legendre (1805). [history note](#).
The approach was first successfully applied to problems in **astronomy**.
- Loss function: squared error loss $\ell(y, \hat{y}) = |y - \hat{y}|^2$
- Many “fancy” machine learning algorithms *today* in literature are still based on this simple least square method.
- Hypothesis space (model class): linear function (affine function with intercept)

History note : “method of least squares” by Gauss and Legendre

Based on d'Alembert's principle, Gauss derived *Principle of least constraint*:

$$Z = \sum_{i=1}^N \frac{1}{2m_i} (\mathbf{F}_i - m_i \mathbf{A}_i)^2$$

\mathbf{F}_i and \mathbf{A}_i are the forces and accelerations, respectively. For free particles, it recovers the classic Newton's motion $\mathbf{F}_i = m_i \mathbf{A}_i$. If constraints prevent the free choice of the \mathbf{A}_i , we can still minimize Z under the given auxiliary conditions. The solution obtained yields the actual motion of the system realized in nature.

Example

A particle is forced to stay on the surface $z = c(x, y)$ by the action of the force \mathbf{F} . Find the motion of the equation. Hint: $\dot{z} = c_x \dot{x} + c_y \dot{y}$ and $\ddot{z} = c_x \ddot{x} + c_{xx} \dot{x}^2 + c_{yy} \ddot{y} + c_{yy} \dot{x}^2 \approx c_x \ddot{x} + c_y \ddot{y}$. The constraint for $\mathbf{A} = (\ddot{x}, \ddot{y}, \ddot{z})$ is the linear equation $\ddot{z} = c_x \ddot{x} + c_y \ddot{y}$.

Simple linear regression

Data $(x_1, y_1), \dots, (x_n, y_n)$.

The linear regression model assumes a specific linear form for f ,

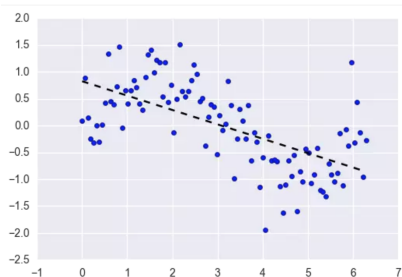
$$f(x) = \beta_0 + \beta x,$$

which is usually thought of as an approximation to the truth.

The loss is also called residual sum of square (RSS)

$$\mathcal{E}(f) = L(\beta_0, \beta) = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

with the prediction $\hat{y}_i = \beta_0 + \beta x_i$



Least squared fitting

Minimize:

$$(\hat{\beta}_0, \hat{\beta}) = \operatorname{argmin}_{\beta_0, \beta} \sum_{i=1}^n (y_i - \beta_0 - \beta x_i)^2.$$

Solution is:

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

- $\hat{y}_i = \hat{\beta}_0 + \hat{\beta}x_i$ are the fitted values
- $r_i = y_i - \hat{y}_i$ are the residuals

Standard errors and confidence intervals

Assume further that

$$y_i = \beta_0 + \beta x_i + \epsilon_i,$$

where $E(\epsilon_i) = 0$ and $\text{Var}(\epsilon_i) = \sigma^2$. Then the standard deviation of $\hat{\beta}$ is

$$se(\hat{\beta}) = \left(\frac{\sigma^2}{\sum (x_i - \bar{x})^2} \right)^{1/2},$$

where σ^2 can be estimated by

$$\hat{\sigma}^2 = \sum (y_i - \hat{y})^2 / (n - 2).$$

Under additional normality assumption of ϵ_i 's, a $(1 - \alpha)100\%$ confidence interval of β is

$$\hat{\beta} \pm z_{\alpha/2} \hat{se}(\hat{\beta}).$$

Ordinary Least Square (OLS)

- The predictor variable $x = (x_0 \equiv 1, x_1, \dots, x_p)$ and **Design Matrix**

$$\mathbf{X} = \begin{bmatrix} 1 & x_{11} & x_{12} & \dots & x_{1p} \\ 1 & x_{21} & x_{22} & \dots & x_{2p} \\ & & \dots & & \\ 1 & x_{n1} & x_{n2} & \dots & x_{np} \end{bmatrix} \in \mathbb{R}^{n \times (p+1)}.$$

n is the number of samples. The first column $x_{i0} \equiv 1$.

- Response vector : $Y = [y_1, y_2, \dots, y_n]^\top$.
- Linear model $\mathcal{H} = \{f : f(x) = \beta^\top x, \beta = (\beta_0, \beta_1, \dots, \beta_p) \in \mathbb{R}^{p+1}\}$.
- Risk minimization:

$$\hat{\beta} = \underset{\beta}{\operatorname{argmin}} \|Y - \mathbf{X}\beta\|_2^2 = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top Y.$$

- Model-based interpretation:

$$Y = \mathbf{X}\beta + \varepsilon, \quad \varepsilon \sim \mathcal{N}(0, \sigma^2).$$

Standardization of Data

The standardization processing is helpful in many cases:

① Centering

- ▶ $x_{ij} \rightarrow x_{ij} - \bar{x}_{.j}$, where $\bar{x}_{.j} = \frac{1}{n} \sum_i x_{ij}$
- ▶ $y_i \rightarrow y_i - \bar{y}$

Then $\sum_i x_{ij} = \sum_i y_i = 0$. the intercept in OLS β_0 vanishes.

- ▶ For centered data: the sample means of the predictor variable x and the response variable y are both zero;
- ▶ The (j, k) -th entry of $\frac{1}{n} \mathbf{X}^T \mathbf{X}$ is $\frac{1}{n} \sum_{i=1}^n (x_{ij} x_{ik}) \approx \text{cov}(X_j, X_k)$. So $\frac{1}{n} \mathbf{X}^T \mathbf{X}$ the (sample) variance-covariance matrix of the predictor variable x .

② Standardization (after centering):

$$x_{ij} \rightarrow \frac{x_{ij}}{\sqrt{\frac{1}{n} \sum_i x_{ij}^2}}.$$

Then $\frac{1}{n} \sum_i x_{ij}^2 \equiv 1, \forall j$.

- ▶ For standardized data, the variance of each factor X_j is unit.
- ▶ It follows that $\text{Trace}(\mathbf{X}^T \mathbf{X}) = \sum_{ij} (x_{ij}^2) = n^2$

Check the linear regression assumption !!!

- The true relationship is linear
- Errors are normally distributed
- Homoscedasticity of errors (or, equal variance around the line).
- Independence of the observations

Read <https://towardsdatascience.com/how-do-you-check-the-quality-of-your-regression-model-in-python>

1

- ① Understanding OLS from the perspective of MLE and Bayes
- ② Understanding uncertainty in $\hat{\beta}$: variance analysis
- ③ Understanding OLS as the minimum variance unbiased estimator of the response : Gauss-Markov theorem
- ④ Understanding OLS from the perspective of linear algebra: orthogonal project, pseudo-inverse, Gram-Schmidt procedure; QR, SVD

Maximize log-likelihood function

$\varepsilon \sim \mathcal{N}(0, \sigma^2)$ leads to the log-likelihood function

$$\begin{aligned}\log \mathcal{L}(\beta; x_i, y_i) &= \log \prod_{i=1}^n p(y_i|x_i)p(x_i) = \sum_{i=1}^n \log p(y_i|x_i) + \sum_{i=1}^n \log p(x_i) \\ &= \sum_{i=1}^n \log \left[\frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y_i - \beta^\top x_i)^2}{2\sigma^2}} \right] + \sum_{i=1}^n \log p(x_i) \\ &= -\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \beta^\top x_i)^2 + \text{terms not depend on } \beta.\end{aligned}$$

Therefore $\hat{\beta}^{\text{MLE}} = \hat{\beta}^{\text{OLS}}$.

Assume the measurement error ε follows other distribution, the other type of loss function ¹ instead of sum of square errors will arise.

¹In statistics, it is called “deviance”. e.g., the Tweedie deviance
Xiang Zhou CityU

- Understanding OLS from the perspective of linear algebra: orthogonal project, pseudo-inverse, Gram-Schmidt procedure; QR, SVD

OLS prediction as the orthogonal projection

- The optimal prediction

$$\hat{Y} = \mathbf{X}\hat{\beta} = \mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T Y =: \text{Proj}_{\mathbf{X}} Y \quad (1)$$

is the orthogonal projection of the vector Y onto the column space of the matrix \mathbf{X} in \mathbb{R}^n

$$\mathbf{X} = \text{span}\{X_0, X_1, \dots, X_p\}$$

- \hat{Y} is the point in \mathbb{R}^n with the shortest Euclidian distance to this subspace \mathbf{X} .
- It would be nice if we have a set of $p+1$ *orthonormal basis vector* of \mathbf{X} . This can be done by Gram-Schmidt procedure (Sec. 3.2.3. in [ESL] under the name “sequential linear regression”).
- In addition, one can use QR, SVD decomposition of $\mathbf{X}^T\mathbf{X}$. To efficiently find the orthogonal projection of the vector Y onto a subspace spanned by X_i in \mathbb{R}^n is a classic topic in numerical linear algebra.

Properties of Projection matrix

Other names used in statistics literature for the projection matrix Proj_X

- influence matrix;
- hat matrix

$$P = \text{Proj}_X = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$$

satisfies

- symmetric: $P = P^T$;
- idempotent: $P^2 = \mathbf{I}_n$ identity matrix;
- $\text{rank} = \dim(X) = p + 1$
- eigenvalues: $p + 1$ ones and $n - (p + 1)$ zeros;
- $\text{trace} = \dim(X)$.

Singular Value Decomposition

- Assume $\mathbf{X} = UDV^T$ is a SVD of the design matrix \mathbf{X} , then $D = \text{diag}\{d_0, \dots, d_p\}$, d_i is the singular value of \mathbf{X} .
- The column vectors of U , $\{U_i, 0 \leq i \leq p\}$, is a set of orthonormal basis of X .
- Then $\mathbf{X}^T \mathbf{X} = VD^2V^T$, and $\text{Proj}_X = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T = (UDV^T)VD^{-2}V^TVDU^T = UU^T$.



$$\hat{Y} = \text{Proj}_X Y = UU^T Y = \sum_{i=0}^p \alpha_i U_i, \quad \text{where } \alpha_i = U_i \cdot Y.$$

Decomposition of Total Sum of Squares

notations

(x_i, y_i) are the data and \hat{y} are the predicted response. For any regression method, define

- SST= total sum of squares for the response variable (proportional to the variance of the response)

$$SST = \sum_i (y_i - \bar{y})^2$$

- SSReg = sum of squares explained by regression

$$SSReg = \sum_i (\hat{y}_i - \bar{y})^2$$

- SSE = sum of squares of errors ¹

$$SSE = \sum_i (y_i - \hat{y}_i)^2$$

¹[ISL] [ESL] name this as RSS= residual sum of squares
Xiang Zhou CityU

coefficient of determination R^2

Definition (coefficient of determination)

$$R^2 = 1 - \frac{SSE}{SST}$$

For OLS with the optimal prediction $\hat{Y} = \mathbf{X}\hat{\beta}$, we have ¹

$$SST = SSReg + SSE$$

For OLS, the coefficient of determination ² is

$$R^2 = \frac{SSReg}{SST} = \frac{\sum_i (\hat{y}_i - \bar{y})^2}{\sum_i (y_i - \bar{y})^2} = \frac{\text{explained sum of squares by regression}}{\text{total sum of square}}$$

¹proof: https://en.wikipedia.org/wiki/Explained_sum_of_squares#Partitioning_in_the_general_ordinary_least_squares_model

²https://scikit-learn.org/stable/modules/model_evaluation.html#r2-score

- Understanding uncertainty in $\hat{\beta}$: unbiasedness, consistence, variance analysis

The distribution of the OLS coefficient $\hat{\beta}$

Since $Y = \mathbf{X}\beta + \varepsilon$, then

$$\begin{aligned}\hat{\beta} &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T Y = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{X}\beta + \varepsilon) \\ &= \beta + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \varepsilon\end{aligned}$$

Note that $\varepsilon \sim N(0, \sigma^2 I_n)$, thus

$$\mathbb{E} \hat{\beta} = \beta \quad (\text{unbiased estimator})$$

$$\begin{aligned}\mathbb{V}(\hat{\beta}) &= \mathbb{V}((\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \varepsilon) \\ &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbb{V}(\varepsilon) (\mathbf{X}^T \mathbf{X}^{-1} \mathbf{X}^T)^T \\ &= \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T I_n \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \\ &= \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1}.\end{aligned}$$

Therefore,

$$\hat{\beta} \sim \mathcal{N}(\beta, \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1}),$$

from which the confidence interval of $\hat{\beta}$ can be calculated.

Consistency of $\hat{\beta}$

Assume that

$$\lim_{n \rightarrow \infty} \left(\frac{\mathbf{X}^T \mathbf{X}}{n} \right) = \Delta$$

exists as a nonstochastic and nonsingular matrix (for example, $|x_{ji}| \leq c$ is bounded). Then

$$\begin{aligned} \lim_{n \rightarrow \infty} \mathbb{E} |\hat{\beta} - \beta|^2 &= \lim_{n \rightarrow \infty} \mathbb{V}(\hat{\beta}) \\ &= \sigma^2 \lim_{n \rightarrow \infty} \frac{1}{n} \left(\frac{X^T X}{n} \right)^{-1} \\ &= \sigma^2 \lim_{n \rightarrow \infty} \frac{1}{n} \Delta^{-1} \\ &= 0 \end{aligned}$$

This implies that OLSE $\hat{\beta}$ converges to true β in quadratic mean. Thus OLSE $\hat{\beta}$ is a consistent estimator of β .

- The distribution of $\hat{Y} = X\hat{\beta}$ is then $\mathcal{N}(\mathbf{X}\beta, \sigma^2 X(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top)$
- When a new data of input x arrives, taking value $x_i = a_i, i = 1, \dots, p$, with $a = (1, a_1, a_2, \dots, a_p)^\top \in \mathbb{R}^{p+1}$, then the prediction from the regression equation is

$$\hat{y} := a^\top \hat{\beta} \sim \mathcal{N}(a^\top \beta, \sigma^2 a^\top (\mathbf{X}^\top \mathbf{X})^{-1} a)$$

which can give the **confidence interval** of $\hat{y} = a^\top \hat{\beta}$.

- But remember that in our model $Y = X\beta + \varepsilon$, it is assumed that the data you *observe* inevitably is contaminated by the measurement error ε . By including this measurement error, the predicted value at this new input $x = a$ is

$$\hat{y} + \varepsilon_a = a^\top \hat{\beta} + \varepsilon_a$$

where ε_a is $\mathcal{N}(0, \sigma_a^2)$ and independent of the training data you used to build the regression equation.

It is clear that the distribution of $\hat{y} + \varepsilon_a$ is

$$\mathcal{N}(a^\top \beta, \sigma^2 a^\top (\mathbf{X}^\top \mathbf{X})^{-1} a + \sigma_a^2),$$

which gives the **prediction interval**.

- Understanding OLS as the best linear unbiased estimator (BLUE) with the smallest MSE.

Gauss-Markov theorem (Rao, 1973)

- Recall that given a training dataset D for supervised learning, the regression function $\hat{f}_D \in \mathcal{H}$. In OLS, we assumed that $\hat{f}_D(x)$ is a linear function of x .
- Now, if we fix a test input $x = a$, $\hat{f}_D(a)$ then is a mapping (statistics) from D to \mathcal{Y} . What if we assume this mapping is linear and consider the **MVU**(minimum variance unbiased) estimator of the ground truth $\beta^T a$ at $x = a$?
- Fix the design matrix \mathbf{X} , then this estimator takes the linear form in the response of training examples Y :

$$Y \rightarrow c^T Y$$

with the coefficient $c \in \mathbb{R}^n$.

Theorem (Gauss-Markov Theorem)

Let u be an unbiased estimate of the ground truth response $a^T\beta$ at the new input $x = a$, and u is in the space of linear transformations from the response training data $Y = \mathbf{X}\beta + \varepsilon$, where $\varepsilon \sim N(0, \sigma^2 I_n)$. This is to say that $u = c^T Y$ for some vector $c \in \mathbb{R}^n$ satisfying $\mathbb{E} u = a^T \beta$ for **any** β in \mathbb{R}^{p+1} . Prove

$$\text{Var}(u) \geq \text{Var}(\hat{y}) = \sigma^2 a^T (\mathbf{X}^T \mathbf{X})^{-1} a$$

where $\hat{y} = a^T \hat{\beta}^{OLS} = a^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T Y$. (see Exercise 3.3 in [ESL].)

Proof.

$\mathbb{E} u = c^T \mathbb{E} Y = c^T \mathbf{X} \beta$ must equal $a^T \beta$ for any β , then

$$\mathbf{X}^T c = a.$$

To minimize $\text{Var}(u) = c^T \mathbb{V}(Y) c = \sigma^2 \|c\|_2^2$, the optimal c is the L_2 -minimal solution of the linear system $\mathbf{X}^T c = a$ (which is exactly the “pseudo-inverse” of \mathbf{X}^T). The remaining is left as an exercise. □

This exercise is optional. If you know Cramer-Rao bound, it is worth trying.

Exercise

Find the Fisher information matrix I , which is the covariance matrix of the parameter-gradient of the log likelihood function $I(\beta) := \mathbb{V}(\partial_\beta \log p(Y; \beta))$ and show that the variance matrix of $\hat{\beta}^{OLS} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T Y$ is the lower bound $I^{-1}(\beta)$