

Topic VII: Regression Models

Wei You



香港科技大學

THE HONG KONG UNIVERSITY OF
SCIENCE AND TECHNOLOGY

Fall, 2023

Introduction

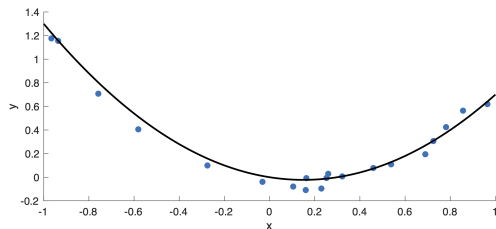
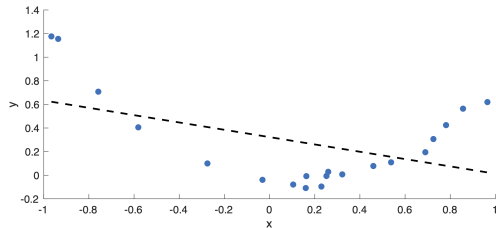
In simple linear regression, we analyze the linear relationship between the response Y and an predictor x .

$$Y = \alpha + \beta x + \varepsilon.$$

- Estimation
- Inference

What if the relationship is not linear?

- $Y = \alpha + \beta x + \varepsilon$? No!
- $Y = \beta_0 + \beta_1 x + \beta_2 x^2 + \varepsilon$.
- Conventional explanation: Fit a parabola to data.



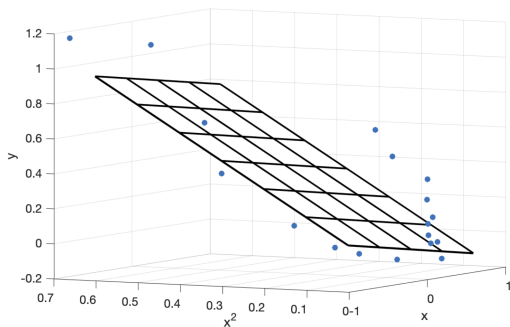
Enlarge the Feature Space

- Alternatively: Let

$$x_1 = x, \quad x_2 = x^2$$

$$Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \varepsilon.$$

- Bottomline: Although Y is quadratic in x , it is linear in the unknown parameters.



More generally, we use non-linear function of the inputs $\phi(x)$ and assume

$$Y = \beta^\top \phi(x) + \varepsilon.$$

If $\phi(x) = (1, x, x^2, x^3, \dots, x^d)$ is the vector of polynomial basis functions, it is called a **polynomial regression**.

Multiple Regression

Multiple linear regression model

For $i = 1, \dots, n$

$$Y_i = \beta_1 x_{i1} + \dots + \beta_p x_{ip} + \varepsilon_i.$$

Or in vector form

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

where $\mathbf{Y} \in \mathbb{R}^{n \times 1}$, $\mathbf{X} \in \mathbb{R}^{n \times p}$, $\boldsymbol{\beta} \in \mathbb{R}^{p \times 1}$.

- \mathbf{X} is called the design matrix, \mathbf{Y} is the response;
- the j -th column of \mathbf{X} , $\mathbf{x}_j = x_{.j} \in \mathbb{R}^n$ is the vector of the j -th predictor;
- $\boldsymbol{\beta}$ is the unknown parameter; ε_i is the error of observation i .
- We usually set $x_{i1} = 1$ to have a non-zero intercept.

Goal: find the parameter values that fits the data the best.

Examples

\mathbf{X} is usually called the **design matrix**.

- The location model: $p = 1$, $\mathbf{X} = (1, \dots, 1)^\top$, $\beta_1 = \mu$.

$$Y = \mu + \varepsilon.$$

- The 2-sample model: $p = 2$, $\mathbf{X} = \begin{pmatrix} 1 & \dots & 1 & 0 & \dots & 0 \\ 0 & \dots & 0 & 1 & \dots & 1 \end{pmatrix}^\top$, $\beta = (\mu_1, \mu_2)^\top$.
- One-way (simple) analysis of variance (ANOVA), compare the mean of k groups:
 $p = k$, $\mathbf{X} = ?$
- Simple linear regression: $Y = \alpha + \beta x + \varepsilon$.

$$\mathbf{X} = \begin{pmatrix} 1 & x_1 \\ \vdots & \vdots \\ 1 & x_n \end{pmatrix}, \quad \beta = (\alpha, \beta)^\top.$$

$$\mathbf{X} = \begin{pmatrix} 1 & 0 \\ \vdots & \vdots \\ 1 & 0 \\ 0 & 1 \\ \vdots & \vdots \\ 0 & 1 \end{pmatrix} \begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix}$$

Handwritten notes: $N(\mu_1, \sigma^2)$ for the first group and $N(\mu_2, \sigma^2)$ for the second group.

Examples

- Higher order regression: $Y = \alpha + \beta_1 x + \beta_2 x^2 + \varepsilon$.

$$\mathbf{X} = \begin{pmatrix} 1 & x_1 & x_1^2 \\ \vdots & \vdots & \vdots \\ 1 & x_n & x_n^2 \end{pmatrix}, \quad \beta = (\alpha, \beta_1, \beta_2)^\top.$$

This seems very powerful combined with the following fact: any continuous function can be approximated arbitrarily well by polynomials. (Weierstrass Theorem)

- Multiplicative model: $Y = \alpha x^\beta \exp(\varepsilon)$. If we take logarithm of both sides, we have $\log(Y) = \log(\alpha) + \beta \log(x) + \varepsilon$.

$$\mathbf{X} = \begin{pmatrix} 1 & \log(x_1) \\ \vdots & \vdots \\ 1 & \log(x_n) \end{pmatrix}, \quad \beta = (\log(\alpha), \beta)^\top.$$

Method of Least Squares

Consider the following linear regression model

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$$

- How to find a good estimator $\boldsymbol{\beta}$?
- For a given $\hat{\boldsymbol{\beta}}$, the fitted values are $\hat{\mathbf{Y}} = \mathbf{X}\hat{\boldsymbol{\beta}}$.
- $\hat{\boldsymbol{\varepsilon}}(\hat{\boldsymbol{\beta}}) \triangleq \mathbf{Y} - \hat{\mathbf{Y}} = \mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}}$ is called the vector of the residuals.
- To estimate $\boldsymbol{\beta}$, minimize the **residual sum of squares (RSS)**:

$$\hat{\boldsymbol{\beta}} = \arg \min_{\boldsymbol{\beta}} \|\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}\|_2^2.$$

RSS is also called sum of squared error (SSE) or sum of squared residuals (SSR).

Least Squares Estimator

$$\hat{\beta} = \arg \min_{\beta} \|Y - X\beta\|_2^2.$$

- Take partial derivative of $\|Y - X\beta\|_2^2 = (Y - X\beta)^\top (Y - X\beta)$ with respect to β , and set them to zero:

$$-2X^\top(Y - X\beta) = 0 \quad \Rightarrow \quad X^\top X\hat{\beta} = X^\top Y$$

Let $f(x) = (f_1(x), \dots, f_m(x))^\top \in \mathbb{R}^m$ for $x \in \mathbb{R}^n$, define $\frac{\partial f}{\partial x} \triangleq [\frac{\partial f_i}{\partial x_j}] \in \mathbb{R}^{m \times n}$.

Matrix differentiation

Let $A \in \mathbb{R}^{m \times n}$, $B \in \mathbb{R}^{n \times n}$, $x \in \mathbb{R}^{n \times 1}$. Let $f(x) = Ax \in \mathbb{R}^m$, $g(x) = x^\top Bx \in \mathbb{R}$, then

$$\frac{\partial f}{\partial x} = A \in \mathbb{R}^{m \times n}, \quad \frac{\partial g}{\partial x} = x^\top (B + B^\top) \in \mathbb{R}^{1 \times n}.$$

Normal equation

$$\mathbf{X}^\top \mathbf{X} \hat{\boldsymbol{\beta}} = \mathbf{X}^\top \mathbf{Y}.$$

- We assumed that $n \geq p$ and \mathbf{X} has full (column) rank, so that $\mathbf{X}^\top \mathbf{X}$ is invertible.
- Given $n \geq p$ and \mathbf{X} is full-rank, we have the **least square estimator (LSE)**

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y},$$

- The **fitted value** of \mathbf{Y} is

$$\hat{\mathbf{Y}} = \mathbf{X} \hat{\boldsymbol{\beta}} = \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}.$$

Collinearity

Collinearity is a phenomenon in which one predictor variable in a multiple regression model can be linearly predicted from the others with a substantial degree of accuracy. In this case, $\mathbf{X}^\top \mathbf{X}$ has extremely small eigenvalues, i.e., it is nearly non-invertible, and the estimation of parameter is sensitive to small changes in the model and data.

Least Squares Estimator – Computation Considerations

$$\mathbf{X}^\top \mathbf{X} \hat{\boldsymbol{\beta}} = \mathbf{X}^\top \mathbf{Y}.$$

The normal equation is usually solved via the Cholesky decomposition of the matrix $\mathbf{X}^\top \mathbf{X}$ or a QR decomposition of \mathbf{X} .

- Cholesky decomposition

$$\mathbf{X}^\top \mathbf{X} = \mathbf{L}\mathbf{L}^\top \Rightarrow \mathbf{L}\mathbf{L}^\top \hat{\boldsymbol{\beta}} = \mathbf{X}^\top \mathbf{Y},$$

where \mathbf{L} is a lower triangular matrix.



- To solve the least square, first let $\mathbf{z} = \mathbf{L}^\top \hat{\boldsymbol{\beta}}$ and solve $\mathbf{L}\mathbf{z} = \mathbf{X}^\top \mathbf{Y}$ for \mathbf{z} using forward substitution, and then solve $\mathbf{L}^\top \hat{\boldsymbol{\beta}} = \mathbf{z}$ for $\boldsymbol{\beta}$ using backward substitution.

Least Squares Estimator – Computation Considerations

$$\mathbf{X}^\top \mathbf{X} \hat{\boldsymbol{\beta}} = \mathbf{X}^\top \mathbf{Y}.$$

- QR decomposition

$$\mathbf{X} = \mathbf{Q}\mathbf{R},$$

where \mathbf{Q} is an orthogonal matrix ($\mathbf{Q}\mathbf{Q}^\top = \mathbf{Q}^\top\mathbf{Q} = \mathbf{I}$) and \mathbf{R} is an upper triangular matrix.

- $\mathbf{X} = \mathbf{Q}\mathbf{R} \Rightarrow \mathbf{X}^\top \mathbf{X} = \mathbf{R}^\top \mathbf{R}$, QR for \mathbf{X} is equivalent to Cholesky for $\mathbf{X}^\top \mathbf{X}$.
- But with n observations and p features, the Cholesky decomposition requires $p^3 + np^2/2$ operations, while the QR decomposition requires np^2 operations.

Depending on the relative size of n and p , the Cholesky can sometimes be faster; on the other hand, it can be less numerically stable (Lawson and Hansen, 1974).

Simple Linear Regression

Example: Simple linear regression. $Y = \alpha + \beta x + \varepsilon$.

$$S_{xx} = \sum_{i=1}^n (x_i - \bar{x})^2 = \sum_{i=1}^n x_i^2 - n\bar{x}^2,$$

$$S_{YY} = \sum_{i=1}^n (Y_i - \bar{Y})^2 = \sum_{i=1}^n Y_i^2 - n\bar{Y}^2$$

$$S_{xY} = \sum_{i=1}^n (x_i - \bar{x})(Y_i - \bar{Y}) = \sum_{i=1}^n x_i Y_i - n\bar{x}\bar{Y}$$

Least Square Estimators

$$\hat{\beta} = \frac{S_{xY}}{S_{xx}}, \quad \hat{\alpha} = \bar{Y} - \hat{\beta}\bar{x}, \quad RSS = \sum_{i=1}^n (Y_i - \hat{\alpha} - \hat{\beta}x_i)^2 = \frac{S_{xx}S_{YY} - S_{xY}^2}{S_{xx}}$$

Geometric Interpretation

Column space of $\mathbf{X} \in \mathbb{R}^{n \times p}$

A dimension p subspace of \mathbb{R}^n , spanned by the column vectors of \mathbf{X} .

$$\mathbf{X}\boldsymbol{\beta} = \beta_1\mathbf{x}_1 + \beta_2\mathbf{x}_2 + \cdots + \beta_p\mathbf{x}_p, \quad \boldsymbol{\beta} \in \mathbb{R}^p.$$

Here $\mathbf{x}_j \in \mathbb{R}^n$ is the j -th column of \mathbf{X} .

- The fitted value

$$\hat{\mathbf{Y}} = \mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}$$

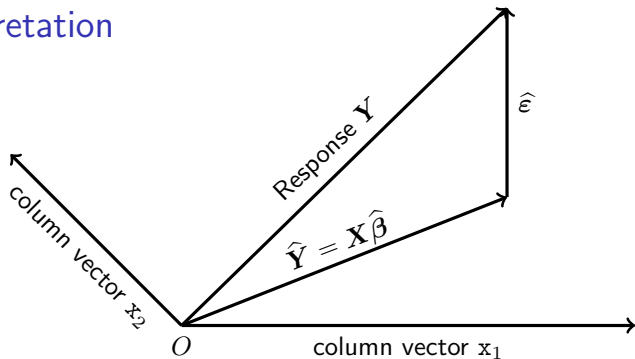
belongs to the column space.

- For the LSE $\hat{\boldsymbol{\beta}}$, we know by normal equation that

$$\hat{\boldsymbol{\varepsilon}}^\top \mathbf{X} = (\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}})^\top \mathbf{X} = 0.$$

It implies that $\hat{\boldsymbol{\varepsilon}}$ is orthogonal to the \mathbf{X} -column space.

Geometric Interpretation



- $\hat{Y} = PY$ is the **projection** of Y on the X -column space, where

$$P \triangleq X(X^T X)^{-1} X^T$$

is the projection matrix.

- $\hat{\epsilon} = QY$ where $Q = I - P$, which is orthogonal to $\hat{Y} = PY$, i.e., $\hat{Y}^T \hat{\epsilon} = 0$.

Note that $PQ = QP = 0$.

Probabilistic Model

Now suppose we have the following assumptions

- Exogeneity: $\mathbb{E}[\varepsilon_i] = 0$.
- Homoscedasticity: $\mathbb{E}[\varepsilon_i^2] = \sigma^2$.
- ε_i 's are independent normal random variables.

The likelihood function of y under parameter β, σ^2 :

$$L(\beta, \sigma^2) \propto \frac{1}{\sigma^n} \exp \left(-\frac{\sum_{i=1}^n (Y_i - \mathbf{x}_i \beta)^2}{2\sigma^2} \right) = \frac{1}{\sigma^n} \exp \left(-\frac{\|\mathbf{Y} - \mathbf{X}\beta\|_2^2}{2\sigma^2} \right).$$

The MLE of β coincides with LSE (does not depend on σ^2).

Best Unbiased Linear Estimator (BLUE)

Assuming exogeneity and Homoscedasticity (but not normality) and under the model $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$, we have the following:

Theorem (Gauss-Markov)

- ① (*Unbiasedness*) $\mathbb{E}[\hat{\boldsymbol{\beta}}] = \mathbb{E}[(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}] = \mathbb{E}[(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top (\mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon})] = \boldsymbol{\beta}$.
- ② (*Conditional standard error*) $\text{Cov}[\hat{\boldsymbol{\beta}}] = \sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1}$.
- ③ *The least-squares estimator $\hat{\boldsymbol{\beta}}$ is the best linear unbiased estimator (BLUE).*

BLUE

For any given vector \mathbf{a} , $\mathbf{a}^\top \hat{\boldsymbol{\beta}}$ is a linear^a unbiased estimator of the parameter $\theta = \mathbf{a}^\top \boldsymbol{\beta}$. Further, for any linear unbiased estimator $\mathbf{b}^\top \mathbf{Y}$ of θ , its variance is at least as large as that of $\mathbf{a}^\top \hat{\boldsymbol{\beta}}$.

^aAn estimator in the form of $\hat{\boldsymbol{\beta}} = \mathbf{A}\mathbf{Y} + \boldsymbol{\mu}$ for some matrix \mathbf{A} and vector $\boldsymbol{\mu}$.

Proof of BLUE. Obviously, $\mathbf{a}^\top \hat{\boldsymbol{\beta}}$ is unbiased and linear. The variance is

$$\text{Var}(\mathbf{a}^\top \hat{\boldsymbol{\beta}}) = \sigma^2 \mathbf{a}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{a}.$$

Consider any linear unbiased estimator $\mathbf{b}^\top \mathbf{Y}$ of $\theta = \mathbf{a}^\top \boldsymbol{\beta}$. Unbiasedness requires that

$$\mathbf{b}^\top \mathbf{X} \boldsymbol{\beta} = \mathbb{E}[\mathbf{b}^\top \mathbf{Y}] = \mathbf{a}^\top \boldsymbol{\beta}, \quad \forall \boldsymbol{\beta}.$$

Hence $\mathbf{b}^\top \mathbf{X} = \mathbf{a}^\top$. Furthermore, the variance is

$$\text{Var}(\mathbf{b}^\top \mathbf{Y}) = \sigma^2 \mathbf{b}^\top \mathbf{b}.$$

Plugging in $\mathbf{a} = \mathbf{X}^\top \mathbf{b}$, we have

$$\text{Var}(\mathbf{b}^\top \mathbf{Y}) - \text{Var}(\mathbf{a}^\top \hat{\boldsymbol{\beta}}) = \sigma^2 \mathbf{b}^\top \mathbf{Q} \mathbf{b} \geq 0.$$

Properties of LSE

- ① $E[\hat{\mathbf{Y}}] = \mathbb{E}[\mathbf{Y}] = \mathbf{X}\boldsymbol{\beta}$.
- ② $\text{Cov}[\hat{\mathbf{Y}}] = \text{Cov}[\mathbf{X}\boldsymbol{\beta} + P\boldsymbol{\varepsilon}] = \sigma^2 P$.
- ③ $\text{Cov}[\hat{\boldsymbol{\varepsilon}}] = \text{Cov}[Q\mathbf{Y}] = \text{Cov}[Q\boldsymbol{\varepsilon}] = \sigma^2 Q$.
- ④ $\text{Cov}[\hat{\mathbf{Y}}, \hat{\boldsymbol{\varepsilon}}] = 0$ as $QP = 0$.

The following is an unbiased estimator of σ^2

$$\hat{\sigma}^2 = \frac{\sum_{i=1}^n (Y_i - \hat{Y}_i)^2}{n - p} = \frac{\|\hat{\boldsymbol{\varepsilon}}\|_2^2}{n - p}.$$

Proof. Note that

$$\sum_{i=1}^n (Y_i - \hat{Y}_i)^2 = \hat{\boldsymbol{\varepsilon}}^\top \hat{\boldsymbol{\varepsilon}} = \mathbf{Y}^\top \mathbf{Q}^\top \mathbf{Q} \mathbf{Y} = \boldsymbol{\varepsilon}^\top \mathbf{Q}^\top \mathbf{Q} \boldsymbol{\varepsilon} = \boldsymbol{\varepsilon}^\top \mathbf{Q} \boldsymbol{\varepsilon}.$$

Use the fact that $\text{tr}(AB) = \text{tr}(BA)$, we have

$$\sum_{i=1}^n (Y_i - \hat{Y}_i)^2 = \text{tr}(\boldsymbol{\varepsilon}^\top \mathbf{Q} \boldsymbol{\varepsilon}) = \text{tr}(\mathbf{Q} \boldsymbol{\varepsilon} \boldsymbol{\varepsilon}^\top)$$

Finally,

$$\mathbb{E} \left[\sum_{i=1}^n (Y_i - \hat{Y}_i)^2 \right] = \mathbb{E} \left[\text{tr}(\mathbf{Q} \boldsymbol{\varepsilon} \boldsymbol{\varepsilon}^\top) \right] = \text{tr} \left(\mathbf{Q} \mathbb{E}[\boldsymbol{\varepsilon} \boldsymbol{\varepsilon}^\top] \right) = \sigma^2 \text{tr}(\mathbf{Q}).$$

To calculate $\text{tr}(\mathbf{Q})$, recall that trace is equal to the summation of the eigenvalues, and note that eigenvalues of \mathbf{Q} is either 0 or 1 and the rank of \mathbf{Q} is $n - p$ (if we assume $\mathbf{X}^\top \mathbf{X}$ is invertible).

Generalized Least Square

What if ϵ is not i.i.d., but $\sim N(\mathbf{0}, \sigma^2 \Sigma)$

- Σ is always positive semi-definite.
- We assume that Σ is **positive definite**, so that $\Sigma = AA^\top$, where A is invertible.
(e.g. Cholesky decomposition)
- We assume that Σ is known.
- We can reduce to the standard model

$$\tilde{Y} \triangleq A^{-1}Y = A^{-1}(X\beta + \epsilon) = A^{-1}X\beta + A^{-1}\epsilon$$

- New independent variables $\tilde{X} = A^{-1}X$, new i.i.d. error $\tilde{\epsilon} = A^{-1}\epsilon$.
- Check:

$$\mathbb{E}[\tilde{\epsilon}] = \mathbb{E}[A^{-1}\epsilon] = A^{-1}\mathbb{E}[\epsilon] = \mathbf{0}$$

$$\text{Cov}[\tilde{\epsilon}] = A^{-1}\text{Cov}[\epsilon](A^{-1})^\top = A^{-1}\sigma^2(AA^\top)(A^{-1})^\top = \sigma^2 I$$

GLS Cont.

- Apply LSE,

$$\begin{aligned}\min \|\tilde{\mathbf{Y}} - \tilde{\mathbf{X}}\boldsymbol{\beta}\|_2^2 &= (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})^\top \Sigma^{-1}(\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}) \\ \Rightarrow \hat{\boldsymbol{\beta}} &= (\mathbf{X}^\top \Sigma^{-1} \mathbf{X})^{-1} \mathbf{X}^\top \Sigma^{-1} \mathbf{Y}.\end{aligned}$$

- This is called the generalized least squares estimate.

$$\hat{\boldsymbol{\beta}} \sim N(\boldsymbol{\beta}, \sigma^2(\mathbf{X}^\top \Sigma^{-1} \mathbf{X})^{-1}).$$

- One can derive test statistics based on the distribution.

GLS Remarks

- A special case is when Σ is a diagonal matrix $diag(v_1, \dots, v_n)$. In this case, we are minimizing the weighted squared error $\sum_i \varepsilon_i^2 / v_i$. This is called the **weighted least square**.
- In practice we usually don't know Σ . So we can run a standard LSE, estimate $\hat{\Sigma}$ (need to impose some structures), and then applied GLS. This is referred to as the Cochrane-Orcutt procedure.
- If Σ is not positive definite (only semi-definite), one can use the Gram-Schmidt process (Schmidt orthonormalization) to find a set of standard normal random variables and the matrix that transform them to the original error terms.
 - Think about the geometric interpretation of covariance and use the fact that uncorrelated means independence for normal random variables.
 - Read the wikipedia page for Gram-Schmidt process.

Properties of LSE – Distributions

Now we impose the additional normal assumption, $\varepsilon \sim N(\mathbf{0}, \sigma^2 I_{n \times n})$.

- $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \varepsilon \sim N(\mathbf{X}\boldsymbol{\beta}, \sigma^2 I_{n \times n})$
- $\hat{\boldsymbol{\beta}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top (\mathbf{X}\boldsymbol{\beta} + \varepsilon) \sim N(\boldsymbol{\beta}, \sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1})$.
- $\hat{\mathbf{Y}} = P\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + P\varepsilon \sim N(\mathbf{X}\boldsymbol{\beta}, \sigma^2 P), P = \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$.
- $\hat{\varepsilon} = Q\mathbf{Y} = Q\varepsilon \sim N(0, \sigma^2 Q), Q = I - P$. Note that $PQ = QP = 0$.
- $\hat{\mathbf{Y}}$ and $\hat{\varepsilon}$ are independent: $\text{Cov}(P\mathbf{Y}, Q\mathbf{Y}) = P\sigma^2 I_{n \times n} Q = 0$.
- $\hat{\boldsymbol{\beta}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{X} \hat{\boldsymbol{\beta}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \hat{\mathbf{Y}} \Rightarrow \hat{\boldsymbol{\beta}}$ and $\hat{\sigma}^2 = \frac{\|\hat{\varepsilon}\|_2^2}{n-p}$ are independent.
- Use the lemma for chi-squared distribution from Lecture 2,

$$(n-p)\hat{\sigma}^2/\sigma^2 = \frac{\|\hat{\varepsilon}\|_2^2}{\sigma^2} = \frac{\sum_{i=1}^n (Y_i - \hat{Y}_i)^2}{\sigma^2} \sim \chi_{n-p}^2.$$

Statistical Inference for β_i

To test whether a particular variable x_i has no effect on \mathbf{Y} :

- Recall that $\hat{\boldsymbol{\beta}}$ and $\hat{\sigma}^2 = \frac{\|\hat{\boldsymbol{\varepsilon}}\|_2^2}{n-p} = RSS/(n-p)$ are independent.
- For each β_i ,

$$\frac{\hat{\beta}_i - \beta_i}{\hat{\sigma} \sqrt{(\mathbf{X}^\top \mathbf{X})_{ii}^{-1}}} = \frac{(\hat{\beta}_i - \beta_i) / \sqrt{\sigma^2 (\mathbf{X}^\top \mathbf{X})_{ii}^{-1}}}{\sqrt{[(n-p)\hat{\sigma}^2/\sigma^2]/(n-p)}} \sim t_{n-p}.$$

Hypothesis Test – β_i

H_0	H_1	TS	Significance level α	p -value
$\beta_i = 0$	$\beta_i \neq 0$	$\frac{\hat{\beta}_i}{\hat{\sigma} \sqrt{(\mathbf{X}^\top \mathbf{X})_{ii}^{-1}}}$	Reject if $ TS > t_{n-p}(1 - \alpha/2)$	$2\mathbb{P}\{T_{n-p} > TS\}$

- The **1 – α confidence interval** for β_i is

$$\left[\hat{\beta}_i - t_{n-p}(1 - \alpha/2) \hat{\sigma} \sqrt{(\mathbf{X}^\top \mathbf{X})_{ii}^{-1}}, \hat{\beta}_i + t_{n-p}(1 - \alpha/2) \hat{\sigma} \sqrt{(\mathbf{X}^\top \mathbf{X})_{ii}^{-1}} \right]$$

Simple Linear Regression

Example: Simple linear regression. $y = \alpha + \beta x + \varepsilon$.

- TS for β , $H_0 : \beta = \beta_0$, implies that $\sqrt{\frac{S_{xx}}{\hat{\sigma}^2}}(\hat{\beta} - \beta_0) \sim t_{n-2}$ under H_0 .
- $(1 - \gamma)$ CI for β :

$$\left(\hat{\beta} - t_{n-2}(1 - \gamma/2) \frac{\hat{\sigma}}{\sqrt{S_{xx}}}, \hat{\beta} + t_{n-2}(1 - \gamma/2) \frac{\hat{\sigma}}{\sqrt{S_{xx}}} \right)$$

- TS for α , $H_0 : \alpha = \alpha_0$, implies that $\sqrt{\frac{S_{xx}n}{\hat{\sigma}^2 \sum_i x_i^2}}(\hat{\alpha} - \alpha_0) \sim t_{n-2}$ under H_0 .
- $(1 - \gamma)$ CI for α :

$$\left(\hat{\alpha} - \sqrt{\frac{\hat{\sigma}^2 \sum_i x_i^2}{S_{xx}n}} t_{n-2}(1 - \gamma/2), \hat{\alpha} + \sqrt{\frac{\hat{\sigma}^2 \sum_i x_i^2}{S_{xx}n}} t_{n-2}(1 - \gamma/2) \right)$$

Statistical Inference for β

To test whether the parameter is β_0 , i.e., $H_0 : \beta = \beta_0$ versus $H_1 : \beta \neq \beta_0$.

$$\frac{(\hat{\beta} - \beta_0)^\top (\mathbf{X}^\top \mathbf{X}) (\hat{\beta} - \beta_0)}{p \hat{\sigma}^2} = \frac{\frac{\varepsilon^\top P \varepsilon}{\sigma^2} / p}{\frac{\varepsilon^\top Q \varepsilon}{\sigma^2} / (n - p)} \sim F_{p, n-p}, \text{ under } H_0.$$

(Check independence!) This is because the numerator equals $\varepsilon^\top P \varepsilon$ and $\text{tr}(P) = p$ and $\hat{\sigma}^2 = \frac{\|\hat{\varepsilon}\|_2^2}{n-p} = \mathbf{Y}^\top Q \mathbf{Y} / (n-p) = \varepsilon^\top Q \varepsilon / (n-p)$ and $\text{tr}(Q) = n-p$.

- Reject if test statistic $\geq F_{p, n-p}(1 - \alpha)$.

$(1 - \alpha)$ confidence set

$$\left\{ \beta : (\beta - \hat{\beta})^\top (\mathbf{X}^\top \mathbf{X}) (\beta - \hat{\beta}) \leq p \hat{\sigma}^2 F_{p, n-p}(1 - \alpha) \right\}$$

An ellipsoid.

Statistical Inference for the True Location of the Hyperplane

- For each experimental condition \mathbf{x}_0 : Let $Y_0 = \mathbf{x}_0^\top \boldsymbol{\beta}$ and $\hat{Y}_0 = \mathbf{x}_0^\top \hat{\boldsymbol{\beta}}$

$$\frac{\hat{Y}_0 - \mathbb{E}[Y_0]}{\hat{\sigma} \sqrt{\mathbf{x}_0^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}_0}} \sim t_{n-p}.$$

⇒ Confidence interval such that the true hyperplane cross **at this \mathbf{x}_0** .

Pointwise confidence band (CB)

$$\mathbb{P}(\mathbf{x}_0^\top \hat{\boldsymbol{\beta}} - w(\mathbf{x}_0) \leq \mathbf{x}_0^\top \boldsymbol{\beta} \leq \mathbf{x}_0^\top \hat{\boldsymbol{\beta}} + w(\mathbf{x}_0)) = 1 - \alpha.$$

- Can we construct a confidence band for the entire true hyperplane ($Y = \mathbf{x}^\top \boldsymbol{\beta}$)?

Simultaneous CB

$$\mathbb{P}(\mathbf{x}^\top \hat{\boldsymbol{\beta}} - w(\mathbf{x}) \leq \mathbf{x}^\top \boldsymbol{\beta} \leq \mathbf{x}^\top \hat{\boldsymbol{\beta}} + w(\mathbf{x}), \text{ for all } \mathbf{x}) = 1 - \alpha.$$

Statistical Inference for the True Location of the Hyperplane

For simultaneous band, by Cauchy-Schwartz for $\langle \alpha, \beta \rangle = \alpha^\top (\mathbf{X}^\top \mathbf{X})^{-1} \beta$

$$\begin{aligned} |\mathbf{x}_0^\top \hat{\beta} - \mathbf{x}_0^\top \beta|^2 &= |\mathbf{x}_0^\top (\mathbf{X}^\top \mathbf{X})^{-1} (\mathbf{X}^\top \mathbf{X}) (\hat{\beta} - \beta)|^2 = \langle \mathbf{x}_0, (\mathbf{X}^\top \mathbf{X}) (\hat{\beta} - \beta) \rangle^2 \\ &\leq (\mathbf{x}_0^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}_0) \left((\hat{\beta} - \beta)^\top (\mathbf{X}^\top \mathbf{X}) (\hat{\beta} - \beta) \right) \end{aligned}$$

From previous slides, we know that

$$\frac{(\hat{\beta} - \beta_0)^\top (\mathbf{X}^\top \mathbf{X}) (\hat{\beta} - \beta_0)}{p \hat{\sigma}^2} \sim F_{p, n-p}$$

Confidence band

$$|\mathbf{x}_0^\top \hat{\beta} - \mathbf{x}_0^\top \beta|^2 \leq (\mathbf{x}_0^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}_0) p \hat{\sigma}^2 F_{p, n-p} (1 - \alpha).$$

Simple Linear Regression

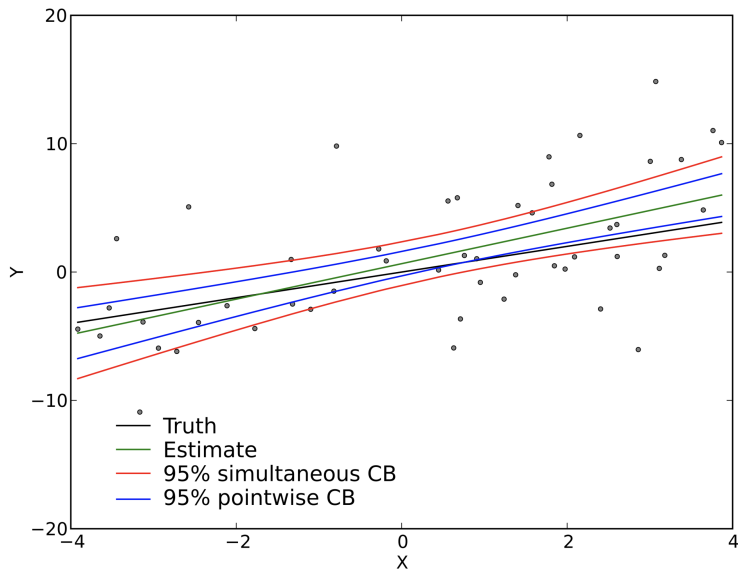
Example: Simple linear regression

Confidence interval for the new observation at x_0

$$\hat{\alpha} + \hat{\beta}x_0 \pm t_{n-2}(1 - \gamma/2) \cdot \sqrt{\frac{RSS}{(n-2)} \left[\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{xx}} \right]}$$

Simultaneous confidence band for the regression line x_0

$$\hat{\alpha} + \hat{\beta}x_0 \pm \sqrt{2F_{n-2}(1 - \gamma/2)} \cdot \sqrt{\frac{RSS}{(n-2)} \left[\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{xx}} \right]}$$



Prediction Interval

In the previous slides, we see confidence band for mean response $\mathbb{E}[Y] = \mathbf{x}^\top \boldsymbol{\beta}$.

What about the response $Y = \mathbf{x}^\top \boldsymbol{\beta} + \varepsilon$?

For each experimental condition \mathbf{x}_0 :

$$\frac{\hat{Y}_0 - Y_0}{\hat{\sigma} \sqrt{\mathbf{1} + \mathbf{x}_0^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}_0}} \sim t_{n-p}.$$

Example: Simple linear regression

$$\text{(For the mean)} \quad \left(\hat{\alpha} + \hat{\beta}x_0 \pm t_{n-2}(1 - \alpha/2) \sqrt{\frac{RSS}{(n-2)} \left[\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{xx}} \right]} \right)$$

$$\text{(For the response)} \quad \left(\hat{\alpha} + \hat{\beta}x_0 \pm t_{n-2}(1 - \alpha/2) \sqrt{\frac{RSS}{(n-2)} \left[\mathbf{1} + \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{xx}} \right]} \right)$$

Nested Model Test – F -statistic

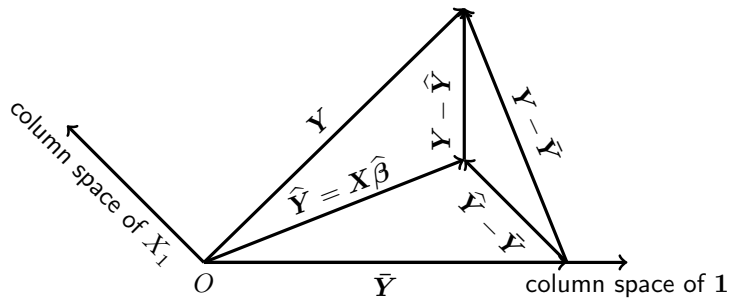
The F statistic checks whether a non-trivial linear model is not rejected.

- $X = \begin{pmatrix} 1 & x_{12} & \cdots & x_{1p} \\ 1 & x_{22} & \cdots & x_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n2} & \cdots & x_{np} \end{pmatrix} = (\mathbf{1}, X_1), \beta = (\beta_1, \beta_2, \dots, \beta_p)^\top.$
- Consider null hypothesis $H_0 : \beta_2 = \beta_3 = \cdots = \beta_p = 0$, the MLE under the null hypothesis is $\bar{\mathbf{Y}}$.

Nested Model Test – F -statistic

The F -statistic

$$F = \frac{\|\bar{\mathbf{Y}} - \hat{\mathbf{Y}}\|_2^2 / (p - 1)}{\|\mathbf{Y} - \hat{\mathbf{Y}}\|_2^2 / (n - p)} \sim F_{p-1, n-p}$$



This can be considered as testing two nested models. $H_0 : \beta_2 = \beta_3 = \dots = \beta_p = 0$ is the simpler model, and if we allow any $\boldsymbol{\beta}$, this is the full model.

Nested Model

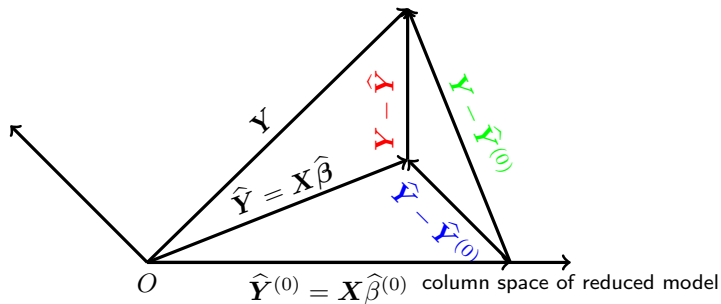
- H_1 : all features play a role $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$.
- Null hypothesis H_0 : a reduced model where some parameters are redundant $\mathbf{B}\boldsymbol{\beta} = \mathbf{b}$.
- For example, $\beta_2 = \beta_3 = \dots = \beta_p = 0$. Or $\beta_1 + \beta_2 = 3$, which means we can reduce $\beta_2 = 3 - \beta_1$ and need only $p - 1$ features.
- For example, \mathbf{B} is $(p - q) \times p$, $\mathbf{b} = \mathbf{0}$. Test if $p - q$ of the coefficients are zero.

$$\mathbf{B} = \begin{pmatrix} 1 & 0 & \dots & 0 & 0 & \dots & 0 \\ 0 & 1 & \dots & 0 & 0 & \dots & 0 \\ \dots & \dots & \dots & \dots & \dots & \dots & \dots \\ 0 & 0 & \dots & 1 & 0 & \dots & 0 \end{pmatrix}$$

- Can we represent the relationship by a simpler model?

Nested Model Test

- If the null hypothesis H_0 is true, say $\beta_2 = 0$. Then the column space of $X^{(0)}$ is a subspace of X .
- Geometric interpretation and Pythagoras theorem.



- Define $RSS = \|\mathbf{Y} - \hat{\mathbf{Y}}\|_2^2$ under H_1 and $RSS_0 = \|\mathbf{Y} - \hat{\mathbf{Y}}^{(0)}\|_2^2$ under H_0 .

Nested Model Test

- Recall the properties of LSE: $RSS/(n-p)$ is an unbiased estimator of σ^2 , under both hypotheses. $(RSS_0 - RSS)/(p-q)$ is an unbiased estimator of σ^2 under H_0 . More importantly, they are **orthogonal**!
- Pythagoras theorem: $\|\mathbf{Y} - \hat{\mathbf{Y}}^{(0)}\|_2^2 = \|\mathbf{Y} - \hat{\mathbf{Y}}\|_2^2 + \|\hat{\mathbf{Y}} - \hat{\mathbf{Y}}^{(0)}\|_2^2$.
- This implies a statistic under H_0 :

$$\frac{(RSS_0 - RSS)/(p-q)}{RSS/(n-p)} = \frac{\|\hat{\mathbf{Y}} - \hat{\mathbf{Y}}^{(0)}\|_2^2/(p-q)}{\|\mathbf{Y} - \hat{\mathbf{Y}}\|_2^2/(n-p)} \sim F_{p-q, n-p}.$$

- Reject the simpler model if the test statistic is large.

Nested Model Test – ANOVA

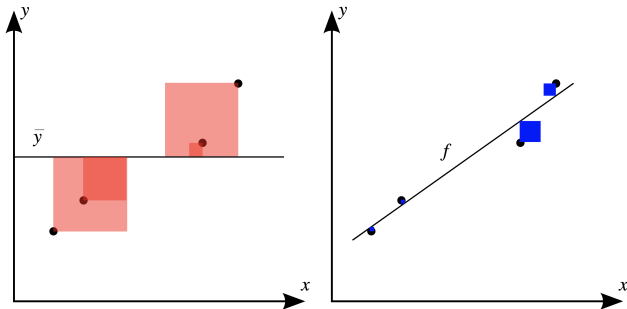
Analysis of variance (ANOVA) is a special case of the nested model test.

- $X = \begin{pmatrix} 1 & \dots & 1 & 0 & \dots & 0 & \dots & 0 & \dots & 0 \\ 0 & \dots & 0 & 1 & \dots & 1 & \dots & 0 & \dots & 0 \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots & & \vdots & \ddots & \vdots \\ 0 & \dots & 0 & 0 & \dots & 0 & \dots & 1 & \dots & 1 \end{pmatrix}^T$, $\beta = (\mu_1, \mu_2, \dots, \mu_k)^T$.
- MLE is $(\bar{Y}_{1\cdot}, \dots, \bar{Y}_{k\cdot})$.
- Consider null hypothesis $H_0 : \mu_1 = \mu_2 = \dots = \mu_k$, the MLE under the null hypothesis is $\hat{\mathbf{Y}}^{(0)} = \bar{\mathbf{Y}}$.
- $\|\mathbf{Y} - \hat{\mathbf{Y}}^{(0)}\|_2^2 = \|\mathbf{Y} - \hat{\mathbf{Y}}\|_2^2 + \|\hat{\mathbf{Y}}^{(0)} - \hat{\mathbf{Y}}\|_2^2 \Rightarrow SST = SSB + SSW$.

Coefficient of Determination

$$R^2 \triangleq \frac{TSS - RSS}{TSS} = 1 - \frac{RSS}{TSS}, \quad TSS = \|\mathbf{Y} - \bar{\mathbf{Y}}\|_2^2.$$

- Total Sum of Squares is the error *not captured by the sample mean*.
- $RSS = \sum (Y_i - \hat{Y}_i)^2$ reflects the variance *not captured by the regression model*.



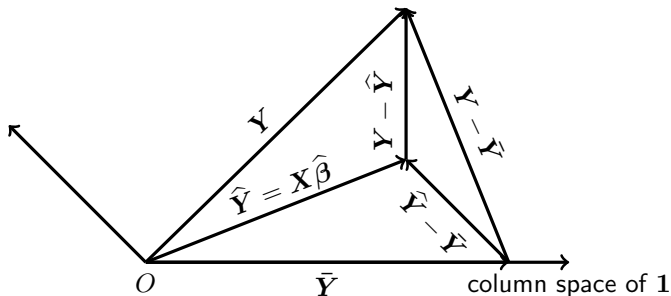
When $R^2 \approx 1$, it means $RSS \approx 0$, the fit is perfect!

When $R^2 \approx 0$, it means the regression model is not much better than \bar{y} .

- One can check that

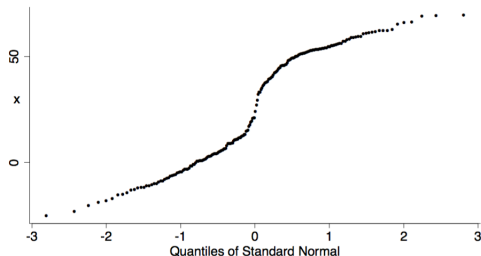
$$R^2 = \frac{\|\hat{\mathbf{Y}} - \bar{\mathbf{Y}}\|_2^2 / (p - 1)}{\|\mathbf{Y} - \bar{\mathbf{Y}}\|_2^2 / (n - p)}.$$

- Also a F distribution? NO!



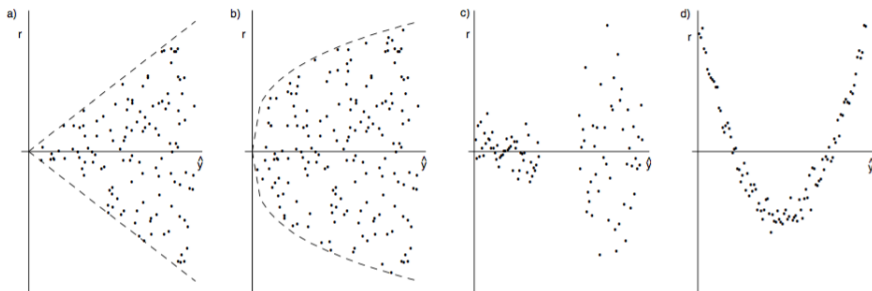
Normality

- Model assumption: ϵ_i has normal distribution.
- **Normal probability plot**: check whether a random sample x_1, \dots, x_n have normal distribution. Plot $\Phi^{-1}((i - 0.5)/n)$ against $(x_{(i)} - \mu)/\sigma$ on the plane.
- If X is indeed normal, the plot should be roughly a straight line.
- Check the **residuals**: $\hat{\epsilon}_i = Y_i - \hat{Y}_i$ should be approximately normal when n is large.



Homoscedasticity

- Model assumption: ϵ_i are i.i.d. random variables.
- Plot \hat{Y}_i against $\hat{\epsilon}_i$. What are the problems of the following figures?



Transformation

In practice we usually take a transformation of the response variable before fitting a regression model.

Example: $Y \propto x_1^{\beta_1} x_2^{\beta_2} \cdots x_p^{\beta_p}$, in which case we take the logarithm

$$\log(Y) = \sum_{i=1}^p \beta_i \log(x_i) + \varepsilon.$$

Box and Cox (1964) proposed a systematic way to find a transformation from the data.

To some extent, Box-Cox transformation can be used to address cases where normality or homoscedasticity assumptions are violated.

Box-Cox Transformation

Consider a parametric family of transformation functions. Let $Y^{(\lambda)}$ denote the transformed response, where λ is a parameter.

$$Y^{(\lambda)} = \begin{cases} \frac{Y^\lambda - 1}{\lambda}, & \text{if } \lambda \neq 0, \\ \log(Y), & \text{if } \lambda = 0. \end{cases}$$

Box-Cox model

$$Y^{(\lambda)} = \mathbf{X}^\top \boldsymbol{\beta} + \varepsilon, \quad \text{where } \varepsilon \stackrel{i.i.d.}{\sim} N(0, \sigma^2)$$

The likelihood function of the Box-Cox model is

$$L(\lambda, \boldsymbol{\beta}, \sigma^2 | \mathbf{X}, \mathbf{Y}) = \frac{1}{(2\pi)^{n/2} \sigma^n} e^{-\frac{1}{2\sigma^2} \|\mathbf{Y}^{(\lambda)} - \mathbf{X}\boldsymbol{\beta}\|^2} \cdot J(\lambda, \mathbf{Y})$$

where the Jacobian is $J(\lambda, \mathbf{Y}) = \prod_{i=1}^n |d\mathbf{Y}^{(\lambda)} / d\mathbf{Y}| = \prod_{i=1}^n |Y_i|^{\lambda-1}$.

MLE for Box-Cox model

- Given a λ , the MLE for $\boldsymbol{\beta}, \sigma^2$ are the usual

$$\hat{\boldsymbol{\beta}}(\lambda) = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}^{(\lambda)}, \quad \hat{\sigma}^2(\lambda) = \frac{1}{n} \|\mathbf{Y}^{(\lambda)} - \hat{\mathbf{Y}}^{(\lambda)}\|^2.$$

- Plugging $\hat{\boldsymbol{\beta}}(\lambda)$ and $\hat{\sigma}^2(\lambda)$ into the likelihood, we have

$$\log L(\lambda) = (\lambda - 1) \sum_{i=1}^n \log(|Y_i|) - \frac{n}{2} \log \hat{\sigma}^2(\lambda) - \frac{n}{2} \Rightarrow \hat{\lambda}_{\text{MLE}} = \arg \max_{\lambda} \log L(\lambda).$$

- The MLE is $(\hat{\lambda}_{\text{MLE}}, \hat{\boldsymbol{\beta}}(\hat{\lambda}_{\text{MLE}}), \hat{\sigma}^2(\hat{\lambda}_{\text{MLE}}))$.

Nonlinear regression

Recall our motivating example – polynomial regression

$$Y = \beta^\top \phi(\mathbf{x}) + \varepsilon,$$

where $\phi(\mathbf{x}) = (1, x, x^2, x^3, \dots, x^d)$.

Weierstrass Theorem

Any continuous $f(X)$ on $[0, 1]$ can be uniformly approximated by a polynomial function.

Spline regression

- One problem of polynomial regression: not suitable for functions with varying degrees of smoothness.
- Solution: use piece-wise polynomial.

Spline

A degree- d spline is a piecewise polynomials with degree d , which is continuous differentiable up to order $d - 1$. The points where discontinuity (in the derivatives) occurs are called the knots.

Example: Linear spline with knots $\tau_1 < \tau_2$.

$$f(x) = \begin{cases} \beta_0 + \beta_1 x, & x \in (-\infty, \tau_1] \\ \beta_0 + \beta_1 x + \beta_2(x - \tau_1)^+, & x \in (\tau_1, \tau_2] \\ \beta_0 + \beta_1 x + \beta_2(x - \tau_1)^+ + \beta_3(x - \tau_2)^+, & x \in (\tau_2, \infty). \end{cases}$$

The basis functions

$$B_0(x) = 1, B_1(x) = x, B_2(x) = (x - \tau_1)^+, B_3(x) = (x - \tau_2)^+.$$

Example: Cubic spline with knots $\tau_1 < \tau_2$.

$$f(x) = \begin{cases} \beta_0 + \beta_1 x + \beta_2 x^2 + \beta_3 x^3, & x \in (-\infty, \tau_1] \\ \beta_0 + \beta_1 x + \beta_2 x^2 + \beta_3 x^3 + \beta_4((x - \tau_1)^+)^3, & x \in (\tau_1, \tau_2] \\ \beta_0 + \beta_1 x + \beta_2 x^2 + \beta_3 x^3 + \beta_4((x - \tau_1)^+)^3 + \beta_5((x - \tau_2)^+)^3, & x \in (\tau_2, \infty). \end{cases}$$

The basis functions

$$B_0(x) = 1, B_1(x) = x, B_2(x) = x^2, B_3(x) = x^3, B_4(x) = ((x - \tau_1)^+)^3, B_5(x) = ((x - \tau_2)^+)^3.$$

- Try to write down cubic splines with k knots.

Spline regression

$$\mathbf{Y} = \boldsymbol{\beta}^\top \mathbf{B}(\mathbf{x}) + \varepsilon.$$

- A multiple linear regression model.
- Versatile and tractable: Cubic splines are widely used.
- Later on, we will look at a fully nonparametric multiple regression called kernel ridge regression.

Influence of Individual Observations on the LSE

Cook's distance for the i th observation

$$D_i = \frac{(\hat{\beta}^{(-i)} - \hat{\beta})^\top (\mathbf{X}^\top \mathbf{X}) (\hat{\beta}^{(-i)} - \hat{\beta})}{p\hat{\sigma}^2} \sim F_{p,n-p}.$$

- Observations with large Cook's distance alters the LSE by a lot, and thus can be considered an outlier.
- It may not be appropriate to simply discard outlier.
- For data that deviates from normal distribution (check by Q-Q plot), especially for those distributions with heavy tails, outliers appear more often.
- So we need robust methods for regression.

Robust Regression

Recall that the multiple linear regression model can be obtained by minimizing the least square error, i.e.,

$$\hat{\beta}_{\text{OLS}} = \arg \min_{\beta} l(\beta) = \arg \min_{\beta} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2 = \arg \min_{\beta} \|\mathbf{Y} - \hat{\mathbf{Y}}\|_2^2.$$

- The loss function is quadratic \Rightarrow loss for outliers have much more impact on the estimation.

Consider the alternative loss function

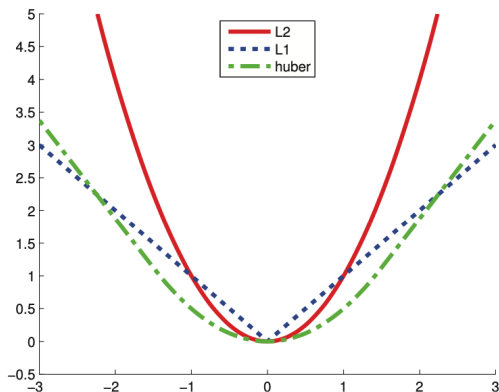
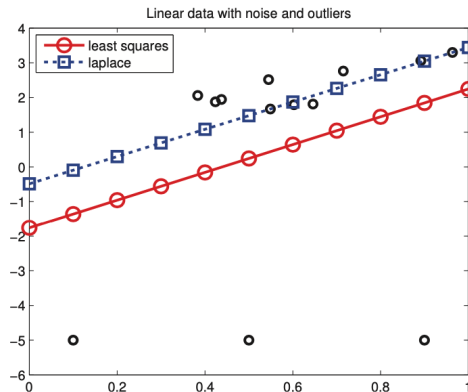
$$\hat{\beta}_{\text{Robust}} = \arg \min_{\beta} \sum_{i=1}^n |Y_i - \hat{Y}_i| = \arg \min_{\beta} \|\mathbf{Y} - \hat{\mathbf{Y}}\|_1.$$

- This is called **robust regression**.

Robust Regression

Robust regression is equivalent to assuming that the response y follows a Laplace distribution.

$$p(Y|\mathbf{X}, \boldsymbol{\beta}, b) = \text{Laplace}(Y|\mathbf{X}, \boldsymbol{\beta}, b) \propto \exp(-|Y - \mathbf{x}^\top \boldsymbol{\beta}|/b).$$



Robust Regression

The loss function in robust regression is not differentiable at 0 and is nonlinear. Fortunately, it can be solved using linear program.

- Consider writing $\varepsilon = \varepsilon^+ - \varepsilon^-$, where $\varepsilon^+ = \max\{0, \varepsilon\}$ and $\varepsilon^- = \max\{0, -\varepsilon\}$.
- One can show that the robust regression estimation is equivalent to

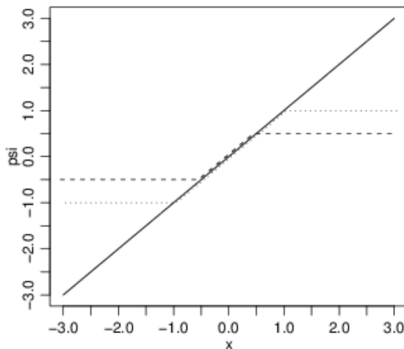
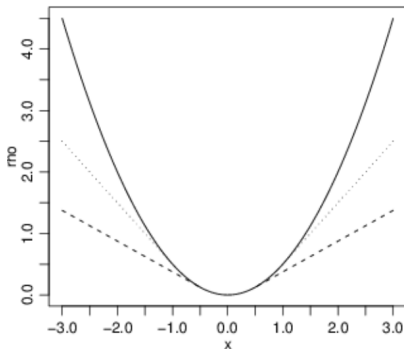
$$\begin{aligned} \hat{\beta}_{\text{Robust}} = \arg \min_{\beta, \varepsilon^+, \varepsilon^-} & \sum_{i=1}^n (\varepsilon_i^+ + \varepsilon_i^-) \\ \text{such that} & \varepsilon_i^+ \geq 0, \varepsilon_i^- \geq 0, \mathbf{X}_i^\top \beta + \varepsilon_i^+ - \varepsilon_i^- = Y_i. \end{aligned}$$

- The resulting linear program can be slow because the dimension is $p + 2n$.

Huber Regression

The loss function in robust regression is not differentiable at 0.

- A trade-off between ordinary least square and robust regression is the **Huber regression**.
- In Huber regression, the error ε is ε^2 for $\varepsilon \leq \delta$ and is $\delta|\varepsilon| - \delta^2/2$.



Robustness of the sample mean

Example: Let X_1, \dots, X_n be i.i.d. $N(\mu, \sigma^2)$. Sample mean can be obtained as the LSE for the location model $p = 1$, $X = (1, \dots, 1)^\top$, $\beta_1 = \mu$.

$$Y = \mu + \varepsilon.$$

- We know \bar{X} is a good estimator for μ under normal distribution and $\text{Var}[\bar{X}_n] = \sigma^2/n$.
- If

$$X_i = \begin{cases} N(\mu, \sigma^2) & w.p. 1 - \delta \\ f(x) & w.p. \delta \end{cases}$$

and $f(x)$ has mean θ and variance τ^2 .

- Then $\text{Var}[\bar{X}] = (1 - \delta)\sigma^2/n + \sigma\tau^2/n + \delta(1 - \delta)(\theta - \mu)^2/n$.
- If $\theta \approx \mu$ and $\sigma \approx \tau$, then we are good.
- If f is Cauchy distribution, then $\text{Var}[\bar{X}] = \infty$.

Median vs Mean

One can check that the median minimizes $\sum |x_i - a|$.

- Median is more robust than the mean. How about its performance?
- For $X \sim F$, we can show that $\sqrt{n}(X_{(n/2)} - \mu)$ is **asymptotically normal** with mean 0 and variance $1/(2f(\mu))^2$.
- The median is the obtained as the estimator of $\beta_1 = \mu$ in the robust regression for the location model $p = 1$, $X = (1, \dots, 1)^\top$, $\beta_1 = \mu$. Robust regression is more robust than least square linear regression.
- There is a **trade-off** between robustness and the performance when the assumed model is correct. The median is not efficient compared to the mean.

Something in-between

Is there an estimator between the mean and the median?

- **Huber estimator:** minimize $\sum \rho(x_i - a)$ where

$$\rho(x) = \begin{cases} x^2/2 & |x| \leq k \\ k|x| - k^2/2 & |x| \geq k \end{cases}$$

Here ρ is differentiable.

- By tuning the parameter k , this estimator can be more “median-like” or “mean-like”.
- We should expect this estimator to be asymptotically normal, because two ends (mean and median) are both so. This is indeed the case.

Reference

Jianqing Fan et al., “Statistical Foundation of Data Science,”

<https://orfe.princeton.edu/~jqfan/fan/classes/525/chapters1-3.pdf>

Kevin Murphy, “Machine Learning: A Probabilistic Prospective,” MIT Press.

John Fox, “Applied Regression Analysis, Linear Models, and Related Methods,” Sage Publications.

Jerome H. Friedman, Robert Tibshirani, and Trevor Hastie, “The Elements of Statistical Learning,” Springer.

Lecture notes by Prof. Kunsch and Meinshausen.