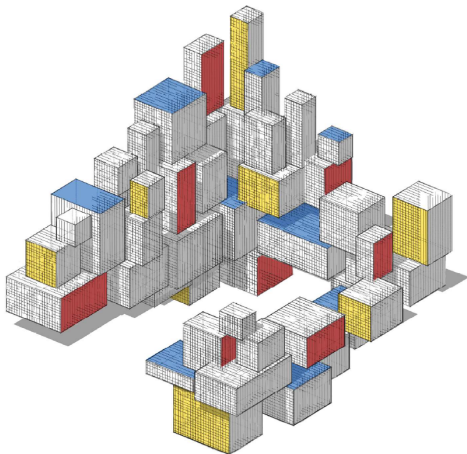


Inference for Normal Linear Models



Purpose

In this lecture we discuss:

Inference for Normal Linear Models

So far we have not assumed any distribution for the random vector of errors $\boldsymbol{\varepsilon} = [\varepsilon_1, \dots, \varepsilon_n]^\top$ in a linear model $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$.

When the error terms $\{\varepsilon_i\}$ are assumed to be normally distributed (that is, $\{\varepsilon_i\} \sim_{\text{iid}} \mathcal{N}(0, \sigma^2)$), whole new avenues open up for inference on linear models.

We already saw that for such **normal linear models**, estimation of $\boldsymbol{\beta}$ and σ^2 can be carried out via maximum likelihood methods, yielding the estimators

- $\hat{\boldsymbol{\beta}} = \mathbf{X}^+ \mathbf{y}$
- $\hat{\sigma}^2 = \|\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}\|^2/n$

The next theorem shows that $\hat{\boldsymbol{\beta}}$ and $\hat{\sigma}^2 n/(n-p)$ are independent and unbiased estimators of $\boldsymbol{\beta}$ and σ^2 , respectively.

Theorem: Properties of the Estimators for a Normal Linear Model

Consider the linear model $Y = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$, with $\boldsymbol{\varepsilon} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}_n)$, where $\boldsymbol{\beta}$ is a p -dimensional vector of parameters and σ^2 a dispersion parameter. The following results hold.

1. The maximum likelihood estimators $\hat{\boldsymbol{\beta}}$ and $\hat{\sigma}^2$ are independent.
2. $\hat{\boldsymbol{\beta}} \sim \mathcal{N}(\boldsymbol{\beta}, \sigma^2(\mathbf{X}^\top \mathbf{X})^+)$.
3. $n\hat{\sigma}^2/\sigma^2 \sim \chi_{n-p}^2$, where $p = \text{rank}(\mathbf{X})$.

Proof

Since $\widehat{\boldsymbol{\beta}} = \mathbf{X}^+ \mathbf{Y}$ is a linear transformation of a normal random vector, it has a multivariate normal distribution.

The mean vector and covariance matrix follow from:

$$\mathbb{E} \widehat{\boldsymbol{\beta}} = \mathbf{X}^+ \mathbb{E} \mathbf{Y} = \mathbf{X}^+ \mathbf{X} \boldsymbol{\beta} = \boldsymbol{\beta}$$

and

$$\mathbb{C}\text{ov}(\widehat{\boldsymbol{\beta}}) = \mathbf{X}^+ \sigma^2 \mathbf{I}_n (\mathbf{X}^+)^{\top} = \sigma^2 (\mathbf{X}^{\top} \mathbf{X})^+.$$

Theorem: Normal and Noncentral χ^2 Distributions

Let $\mathbf{X} \sim \mathcal{N}(\boldsymbol{\mu}, \mathbf{I}_n)$ be an n -dimensional normal random vector and let $\mathcal{V}_k \subset \mathcal{V}_m$ be linear subspaces of dimensions k and m , respectively, with $k < m \leq n$. Let \mathbf{X}_k and \mathbf{X}_m be orthogonal projections of \mathbf{X} onto \mathcal{V}_k and \mathcal{V}_m , and let $\boldsymbol{\mu}_k$ and $\boldsymbol{\mu}_m$ be the corresponding projections of $\boldsymbol{\mu}$. Then, the following holds.

1. The random vectors \mathbf{X}_k , $\mathbf{X}_m - \mathbf{X}_k$, and $\mathbf{X} - \mathbf{X}_m$ are independent.
2. $\|\mathbf{X}_k\|^2 \sim \chi_k^2(\|\boldsymbol{\mu}_k\|)$, $\|\mathbf{X}_m - \mathbf{X}_k\|^2 \sim \chi_{m-k}^2(\|\boldsymbol{\mu}_m - \boldsymbol{\mu}_k\|)$, and $\|\mathbf{X} - \mathbf{X}_m\|^2 \sim \chi_{n-m}^2(\|\boldsymbol{\mu} - \boldsymbol{\mu}_m\|)$.

Proof

Define $\mathbf{Y}^{(2)} = \mathbf{X}\hat{\boldsymbol{\beta}}$. Note that \mathbf{Y}/σ has a $\mathcal{N}(\boldsymbol{\mu}, \mathbf{I}_n)$ distribution, with expectation vector $\boldsymbol{\mu} = \mathbf{X}\boldsymbol{\beta}/\sigma$.

A direct application of the projection theorem shows that $(\mathbf{Y} - \mathbf{Y}^{(2)})/\sigma$ is independent of $\mathbf{Y}^{(2)}/\sigma$.

Since $\hat{\boldsymbol{\beta}} = \mathbf{X}^+\mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{X}^+\mathbf{Y}^{(2)}$ and $\widehat{\sigma^2} = \|\mathbf{Y} - \mathbf{Y}^{(2)}\|^2/n$, it follows that $\widehat{\sigma^2}$ is independent of $\hat{\boldsymbol{\beta}}$.

Finally, by the same theorem, the random variable $\|\mathbf{Y} - \mathbf{Y}^{(2)}\|^2/\sigma^2$ has a χ_{n-p}^2 distribution, as $\mathbf{Y}^{(2)}$ has the same expectation vector as \mathbf{Y} . \square

As a corollary, we see that each estimator $\hat{\beta}_i$ of β_i has a normal distribution with expectation β_i and variance $\sigma^2 \mathbf{u}_i^\top \mathbf{X}^+ (\mathbf{X}^+)^\top \mathbf{u}_i = \sigma^2 \|\mathbf{u}_i^\top \mathbf{X}^+\|^2$, where $\mathbf{u}_i = [0, \dots, 0, 1, 0, \dots, 0]^\top$ is the i -th unit vector; in other words, the variance is $\sigma^2 [(\mathbf{X}^\top \mathbf{X})^+]_{ii}$.

Hypothesis Testing

It is of interest to test the hypothesis $H_0 : \beta_i = 0$ against $H_1 : \beta_i \neq 0$, using the test statistic

$$T = \frac{\widehat{\beta}_i / \|\mathbf{u}_i^\top \mathbf{X}^+\|}{\sqrt{\text{RSE}}}, \quad (1)$$

where RSE is the residual squared error; that is $\text{RSE} = \text{RSS}/(n - p)$.

This test statistic has a t_{n-p} distribution under H_0 .

To see this, write $T = Z/\sqrt{V/(n - p)}$, with

$$Z = \frac{\widehat{\beta}_i}{\sigma \|\mathbf{u}_i^\top \mathbf{X}^+\|} \quad \text{and} \quad V = n \widehat{\sigma^2} / \sigma^2.$$

By the previous theorem, $Z \sim \mathcal{N}(0, 1)$ under H_0 , $V \sim \chi_{n-p}^2$, and Z and V are independent. The result follows from the properties of the t distribution.

Comparing Two Normal Linear Models

Suppose $\mathbf{Y} = [Y_1, \dots, Y_n]^\top$ is of the form

$$\mathbf{Y} = \underbrace{\mathbf{X}_1\boldsymbol{\beta}_1 + \mathbf{X}_2\boldsymbol{\beta}_2}_{\mathbf{X}\boldsymbol{\beta}} + \boldsymbol{\varepsilon}, \quad \boldsymbol{\varepsilon} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}_n), \quad (2)$$

where $\boldsymbol{\beta}_1$ and $\boldsymbol{\beta}_2$ are unknown vectors of dimension k and $p - k$, respectively; and \mathbf{X}_1 and \mathbf{X}_2 are full-rank model matrices of dimensions $n \times k$ and $n \times (p - k)$, respectively.

We wish to test the hypothesis $H_0 : \boldsymbol{\beta}_2 = \mathbf{0}$ against $H_1 : \boldsymbol{\beta}_2 \neq \mathbf{0}$.

Idea: compare the residual sum of squares for both models, expressed as $\|\mathbf{Y} - \mathbf{Y}^{(2)}\|^2$ and $\|\mathbf{Y} - \mathbf{Y}^{(1)}\|^2$.

Using Pythagoras' theorem we saw that

$$\|\mathbf{Y} - \mathbf{Y}^{(2)}\|^2 - \|\mathbf{Y} - \mathbf{Y}^{(1)}\|^2 = \|\mathbf{Y}^{(2)} - \mathbf{Y}^{(1)}\|^2$$

So we can base the decision whether to retain or reject H_0 on the basis of the quotient of $\|Y^{(2)} - Y^{(1)}\|^2$ and $\|Y - Y^{(2)}\|^2$.

Theorem: Test Statistic for Comp. Two Normal Linear Models

For the model (2), let $Y^{(2)}$ and $Y^{(1)}$ be the projections of Y onto the space spanned by the p columns of X and the k columns of X_1 , respectively. Then under $H_0 : \beta_2 = \mathbf{0}$ the test statistic

$$T = \frac{\|Y^{(2)} - Y^{(1)}\|^2 / (p - k)}{\|Y - Y^{(2)}\|^2 / (n - p)} \quad (3)$$

has an $F(p - k, n - p)$ distribution.

Proof

Define $\mathbf{X} := \mathbf{Y}/\sigma$ with expectation $\boldsymbol{\mu} := \mathbf{X}\boldsymbol{\beta}/\sigma$, and $X_j := Y^{(j)}/\sigma$ with expectation μ_j , $j = k, p$.

Note that $\mu_p = \boldsymbol{\mu}$ and, under H_0 , $\mu_k = \mu_p$. We can directly apply the projection theorem to find that

- $\|\mathbf{Y} - \mathbf{Y}^{(2)}\|^2/\sigma^2 = \|\mathbf{X} - \mathbf{X}_p\|^2 \sim \chi_{n-p}^2$ and,
- under H_0 , $\|\mathbf{Y}^{(2)} - \mathbf{Y}^{(1)}\|^2/\sigma^2 = \|\mathbf{X}_p - \mathbf{X}_k\|^2 \sim \chi_{p-k}^2$.

Moreover, these random variables are independent of each other.

Finally, apply the following result from probability:

Theorem: Relationship Between χ^2 and F Distributions

Let $U \sim \chi_m^2$ and $V \sim \chi_n^2$ be independent. Then,

$$\frac{U/m}{V/n} \sim F(m, n).$$

Hypothesis Testing for Nested Models

For nested models $[\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_k]$, $k = 1, 2, \dots, d$, the F test can now be used to test whether certain \mathbf{X}_i are needed or not. In particular, software packages will report the outcomes of

$$F_i = \frac{\|\mathbf{Y}^{(i)} - \mathbf{Y}^{(i-1)}\|/p_i}{\|\mathbf{Y} - \mathbf{Y}^{(d)}\|/(n - p)},$$

in the order $i = 2, 3, \dots, d$.

Under the null hypothesis that $\mathbf{Y}^{(i)}$ and $\mathbf{Y}^{(i-1)}$ have the same expectation (that is, adding \mathbf{X}_i to \mathbf{X}_{i-1} has no additional effect on reducing the approximation error), the test statistic F_i has an $F(p_i, n - p)$ distribution, and the corresponding P-values quantify the strength of the decision to include an additional variable in the model or not. This procedure is called **analysis of variance** (ANOVA).

Crop Yield (cont.)

We continue the crop yield example. Decompose the linear model as

$$Y = \underbrace{\begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix}}_{\mathbf{X}_1} \underbrace{\beta_0}_{\beta_1} + \underbrace{\begin{bmatrix} 0 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}}_{\mathbf{X}_2} \underbrace{\begin{bmatrix} \beta_{12} \\ \beta_{13} \\ \beta_{14} \end{bmatrix}}_{\beta_2} + \underbrace{\begin{bmatrix} C \\ C \\ C \\ C \end{bmatrix}}_{\mathbf{X}_3} \underbrace{\begin{bmatrix} \beta_{22} \\ \beta_{23} \\ \beta_{24} \end{bmatrix}}_{\beta_3} + \varepsilon.$$

Is the crop yield dependent on treatment levels as well as blocks?

We first test $\beta_3 = \mathbf{0}$ versus $\beta_3 \neq \mathbf{0}$.

$Y^{(2)}$ is the projection of Y onto the $(p = 7)$ -dimensional space spanned by the columns of $\mathbf{X} = [\mathbf{X}_1, \mathbf{X}_2, \mathbf{X}_3]$.

$Y^{(1)}$ is the projection of Y onto the $(k = 4)$ -dimensional space spanned by the columns of $\mathbf{X}_{12} := [\mathbf{X}_1, \mathbf{X}_2]$

The test statistic, T_{12} say, under H_0 has an $F(3, 9)$ distribution.

Crop Yield (cont.)

The Python code below calculates the outcome of the test statistic T_{12} and the corresponding P-value. We find $t_{12} = 34.9998$, which gives a P-value 2.73×10^{-5} . This shows that the block effects are extremely important for explaining the data.

Using the extended model (including the block effects), we can test whether $\beta_2 = \mathbf{0}$ or not; that is, whether the treatments have a significant effect on the crop yield in the presence of the Block factor.

This is done in the last six lines of the code below. The outcome of the test statistic is 4.4878, with a P-value of 0.0346.

By including the block effects, we effectively reduce the uncertainty in the model and are able to more accurately assess the effects of the treatments, to conclude that the treatment seems to have an effect on the crop yield.

```
import numpy as np
from scipy.stats import f
from numpy.linalg import lstsq, norm

yy = np.array([9.2988, 9.4978, 9.7604, 10.1025,
               8.2111, 8.3387, 8.5018, 8.1942,
               9.0688, 9.1284, 9.3484, 9.5086,
               8.2552, 7.8999, 8.4859, 8.9485]).reshape(4,4).T

nrow, ncol = yy.shape[0], yy.shape[1]
n = nrow * ncol
y = yy.reshape(16,)
X_1 = np.ones((n,1))

KM = np.kron(np.eye(ncol),np.ones((nrow,1)))
KM[:,0]
X_2 = KM[:,1:ncol]
IM = np.eye(nrow)
C = IM[:,1:nrow]

X_3 = np.vstack((C, C))
X_3 = np.vstack((X_3, C))
X_3 = np.vstack((X_3, C))
```

```

X = np.hstack((X_1,X_2))
X = np.hstack((X,X_3))

p = X.shape[1] #number of parameters in full model
betahat = lstsq(X, y,rcond=None)[0] #estimate under the full model

ym = X @ betahat

X_12 = np.hstack((X_1, X_2)) #omitting the block effect
k = X_12.shape[1] #number of parameters in reduced model
betahat_12 = lstsq(X_12, y,rcond=None)[0]
y_12 = X_12 @ betahat_12
T_12=(n-p)/(p-k)*(norm(y-y_12)**2 - norm(y-ym)**2)/norm(y-ym)**2
pval_12 = 1 - f.cdf(T_12,p-k,n-p)

X_13 = np.hstack((X_1, X_3)) #omitting the treatment effect
k = X_13.shape[1] #number of parameters in reduced model
betahat_13 = lstsq(X_13, y,rcond=None)[0]
y_13 = X_13 @ betahat_13
T_13=(n-p)/(p-k)*(norm(y-y_13)**2 - norm(y-ym)**2)/norm(y-ym)**2
pval_13 = 1 - f.cdf(T_13,p-k,n-p)

```

Confidence and Prediction Intervals

Suppose we wish to predict how a new response variable will behave on the basis of a new explanatory vector \mathbf{x} .

Thus, consider a new \mathbf{x} and let $Y \sim \mathcal{N}(\mathbf{x}^\top \boldsymbol{\beta}, \sigma^2)$, with $\boldsymbol{\beta}$ and σ^2 unknown.

First we are going to look at the *expected* value $\mathbb{E}Y = \mathbf{x}^\top \boldsymbol{\beta}$.

This can estimate it via $\hat{Y} = \mathbf{x}^\top \hat{\boldsymbol{\beta}}$, where $\hat{\boldsymbol{\beta}} \sim \mathcal{N}(\boldsymbol{\beta}, \sigma^2(\mathbf{X}^\top \mathbf{X})^+)$.

It follows that $\hat{Y} \sim \mathcal{N}(\mathbf{x}^\top \boldsymbol{\beta}, \sigma^2 \|\mathbf{x}^\top \mathbf{X}^+\|^2)$.

Let $Z \sim \mathcal{N}(0, 1)$ be the standardized version of \hat{Y} and $V = \|\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}}\|^2 / \sigma^2 \sim \chi_{n-p}^2$. Then the random variable

$$T := \frac{(\mathbf{x}^\top \hat{\boldsymbol{\beta}} - \mathbf{x}^\top \boldsymbol{\beta}) / \|\mathbf{x}^\top \mathbf{X}^+\|}{\|\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}}\| / \sqrt{(n-p)}} = \frac{Z}{\sqrt{V/(n-p)}}$$

has a t_{n-p} distribution.

Confidence Interval for the Expected Response

After rearranging the identity $\mathbb{P}(|T| \leq t_{n-p;1-\alpha/2}) = 1 - \alpha$, where $t_{n-p;1-\alpha/2}$ is the $(1 - \alpha/2)$ quantile of the t_{n-p} distribution, we arrive at the stochastic **confidence interval**

$$\mathbf{x}^\top \hat{\boldsymbol{\beta}} \pm t_{n-p;1-\alpha/2} \sqrt{\text{RSE}} \|\mathbf{x}^\top \mathbf{X}^+\|,$$

where we have identified $\|Y - \mathbf{X}\hat{\boldsymbol{\beta}}\|^2/(n-p)$ with RSE.

This confidence interval quantifies the uncertainty in the learner (regression surface).

Prediction Interval

Construct an interval such that Y lies in this interval with a certain guaranteed probability. We have *two* sources of variation:

1. $Y \sim \mathcal{N}(\mathbf{x}^\top \boldsymbol{\beta}, \sigma^2)$ itself is a random variable.
2. Estimating $\mathbf{x}^\top \boldsymbol{\beta}$ via \widehat{Y} brings another source of variation.

Since $Y \sim \mathcal{N}(\mathbf{x}^\top \boldsymbol{\beta}, \sigma^2)$ and $\widehat{Y} \sim \mathcal{N}(\mathbf{x}^\top \boldsymbol{\beta}, \sigma^2 \|\mathbf{x}^\top \mathbf{X}^+\|^2)$ are independent, it follows that $Y - \widehat{Y} \sim \mathcal{N}(0, \sigma^2(1 + \|\mathbf{x}^\top \mathbf{X}^+\|^2))$.

Letting $Z \sim \mathcal{N}(0, 1)$ be the standardized version of $Y - \widehat{Y}$, and repeating the steps used for the construction of the confidence interval, we arrive at the [prediction interval](#)

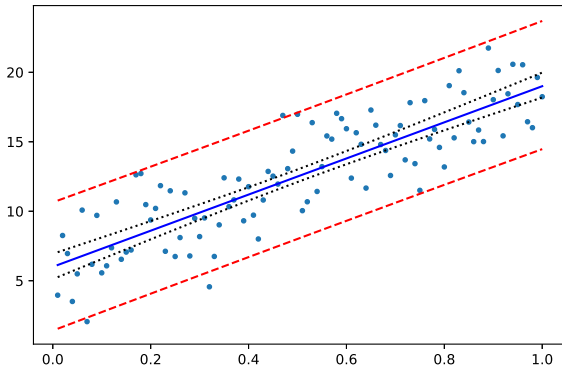
$$\mathbf{x}^\top \widehat{\boldsymbol{\beta}} \pm t_{n-p; 1-\alpha/2} \sqrt{\text{RSE}} \sqrt{1 + \|\mathbf{x}^\top \mathbf{X}^+\|^2}.$$

This prediction interval captures the uncertainty from an as-yet-unobserved response as well as the uncertainty in the parameters of the regression model itself.

Example: Confidence Limits in Linear Regression

The data ($n = 100$ samples) are drawn from a simple linear regression model with parameters $\beta = [6, 13]^\top$ and $\sigma = 2$, where the x -coordinates are evenly spaced on the interval $[0, 1]$.

The figure shows 95% numeric confidence and prediction curves as well as the true regression line.



```
import numpy as np
import matplotlib.pyplot as plt
from scipy.stats import t
from numpy.linalg import inv, lstsq, norm
np.random.seed(123)

n = 100
x = np.linspace(0.01, 1, 100).reshape(n, 1)
# parameters
beta = np.array([6, 13])
sigma = 2
Xmat = np.hstack((np.ones((n, 1)), x)) #design matrix
y = Xmat @ beta + sigma*np.random.randn(n)

# solve the normal equations
betahat = lstsq(Xmat, y, rcond=None)[0]
# estimate for sigma
sqMSE = norm(y - Xmat @ betahat)/np.sqrt(n-2)
tquant = t.ppf(0.975, n-2) # 0.975 quantile
ucl = np.zeros(n) #upper conf. limits
lcl = np.zeros(n) #lower conf. limits
upl = np.zeros(n)
lpl = np.zeros(n)
```

```

rl = np.zeros(n) # (true) regression line
u = 0

for i in range(n):
    u = u + 1/n;
    xvec = np.array([1,u])
    sqc = np.sqrt(xvec.T @ inv(Xmat.T @ Xmat) @ xvec)
    sqp = np.sqrt(1 + xvec.T @ inv(Xmat.T @ Xmat) @ xvec)
    rl[i] = xvec.T @ beta;
    ucl[i] = xvec.T @ betahat + tquant*sqMSE*sqc;
    lcl[i] = xvec.T @ betahat - tquant*sqMSE*sqc;
    upl[i] = xvec.T @ betahat + tquant*sqMSE*sqp;
    lpl[i] = xvec.T @ betahat - tquant*sqMSE*sqp;

plt.plot(x,y, '.')
plt.plot(x,rl, 'b')
plt.plot(x,ucl, 'k:')
plt.plot(x,lcl, 'k:')
plt.plot(x,upl, 'r--')
plt.plot(x,lpl, 'r--')

```