



# Analisis Regresi Linier

---

Kuliah 4 STA1381 – Pengantar Sains  
Data

Septian Rahardiantoro

# Outline

- Pengantar Pemodelan Statistika
- Analisis Regresi
  - Regresi Linier Sederhana
  - Regresi Linier Berganda
- Dummy Variable

# Pengantar Pemodelan Statistika

$$y = f(x_1, x_2, \dots, x_p) + \varepsilon$$

- **Membangun miniatur dari dunia nyata**
  - dinyatakan dalam satu atau beberapa fungsi matematis
- **Menyederhanakan fenomena nyata sehingga mudah memahami pola umum yang ada**
  - memberikan penjelasan terhadap perubahan
  - memberikan penjelasan tentang perbedaan yang terjadi
  - menemukan faktor yang menyebabkan perubahan dan perbedaan

# Pemodelan

- Tujuan/Manfaat:

- Sering digunakan untuk meng-explore dataset yang dimiliki
- Digunakan untuk melakukan prediksi berdasarkan informasi dari variabel prediktor
- Digunakan untuk mengkaji dan memahami bagaimana suatu variabel berhubungan dengan variabel yang lain

- Are not perfect

- “All models are wrong, but some are useful” (GEP Box)

# Beberapa Model Statistika yang Populer

Jenis Variabel Target	Model Statistika
<u>Numerik</u>	<u>Regresi Linier</u>
Kategorik	Regresi Logistik Pohon Klasifikasi (Classification Tree)

# Analisis Regresi

→ Analisis statistika yang memanfaatkan hubungan sebab akibat antara dua atau lebih peubah kuantitatif sehingga salah satu peubah dapat diramalkan dari peubah lainnya.

- Hubungan Antar Peubah:

- Fungsional (deterministik) →  $Y=f(X)$  ; misalnya:  $Y=10X$
- Statistik (stokastik) → amatan tidak jatuh pas pada kurva (terdapat galat)
  - Mis: IQ vs Prestasi, Berat vs Tinggi, Dosis Pupuk vs Produksi

**Analisis Regresi**



- Analisis Regresi digunakan untuk:
  - Menjelaskan dampak perubahan peubah prediktor terhadap peubah respon
    - Memprediksi nilai dari peubah respon berdasarkan nilai dari setidaknya sebuah peubah prediktor

Peubah Respon (peubah tak bebas, peubah terikat, dependent variable):  
peubah yang ingin kita jelaskan

Peubah Prediktor (peubah bebas, independent variable): peubah yang  
digunakan untuk menjelaskan peubah respon

# Regresi Linier

- Syarat Utama: Variabel output (Y) bersifat numerik
- Variabel prediktor (X)
  - numerik OK, kategorik OK
  - satu OK, lebih dari satu OK

- Bentuk model

$$Y = \beta_0 + \beta_1 X_1 + \varepsilon$$

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \varepsilon$$

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_p X_p + \varepsilon$$



# Regresi Linier Sederhana

- Suatu pendekatan untuk memprediksi peubah respon kuantitatif  $Y$  berdasarkan sebuah peubah prediktor  $X$
- Pendekatan ini mengasumsikan bahwa ada hubungan linier antara  $X$  dan  $Y$

The population regression model:

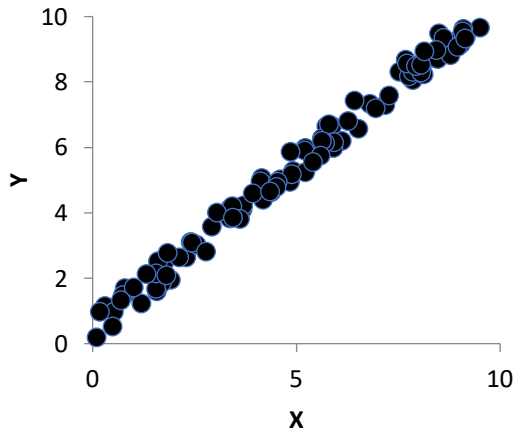
The diagram illustrates the population regression model equation  $y = \beta_0 + \beta_1 x + \epsilon$ . The equation is enclosed in a yellow box. Labels with arrows point to each term: 'Dependent Variable' points to  $y$ ; 'Population y intercept' points to  $\beta_0$ ; 'Population Slope Coefficient' points to  $\beta_1$ ; 'Independent Variable' points to  $x$ ; and 'Random Error term, or residual' points to  $\epsilon$ . Below the equation, two curly braces group the terms: the first brace under  $\beta_0 + \beta_1 x$  is labeled 'Linear component', and the second brace under  $\epsilon$  is labeled 'Random Error component'.

$$y = \beta_0 + \beta_1 x + \epsilon$$

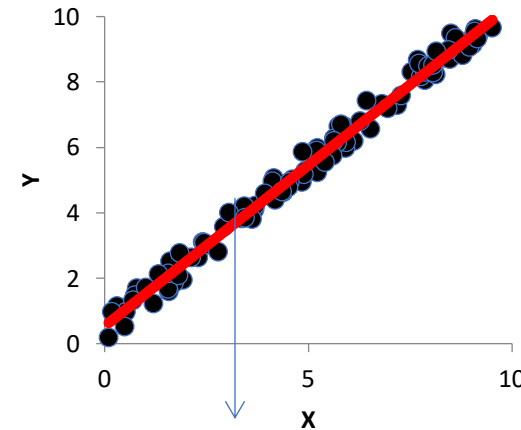
Labels and components:

- Dependent Variable:  $y$
- Population y intercept:  $\beta_0$
- Population Slope Coefficient:  $\beta_1$
- Independent Variable:  $x$
- Random Error term, or residual:  $\epsilon$
- Linear component:  $\beta_0 + \beta_1 x$
- Random Error component:  $\epsilon$

# Regresi Linier Sederhana



Plot data Y dengan X



diperoleh persamaan garis regresi

$$\hat{Y} = b_0 + b_1X$$

Sehingga:

Model regresi

$$Y = \beta_0 + \beta_1X + \varepsilon$$

*b<sub>0</sub>*     *b<sub>1</sub>*

diduga oleh

Persamaan regresi

$$\hat{Y} = b_0 + b_1X$$

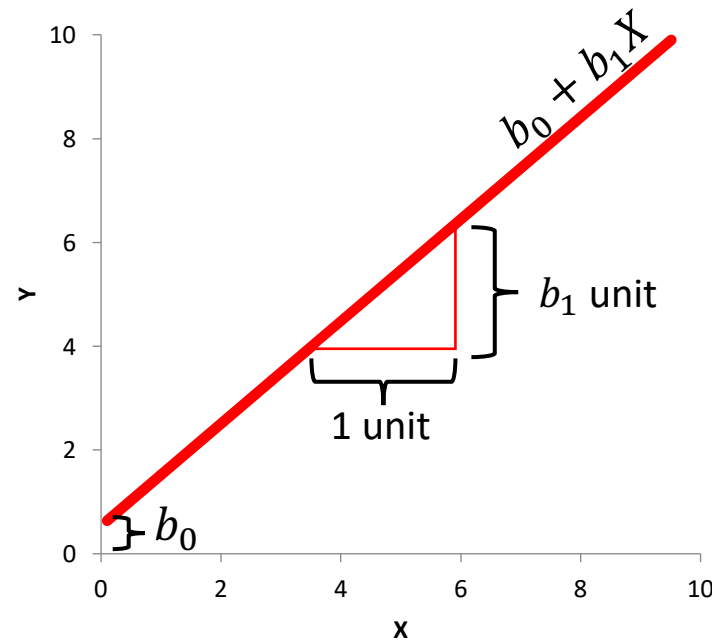
# Regresi Linier Sederhana

- Maka

- $\beta_0$  diduga oleh  $b_0$
- $\beta_1$  diduga oleh  $b_1$



Interpretasi?



$b_0$  adalah nilai rata-ran Y ketika  $X=0$   
(tidak dapat diinterpretasikan oleh X)

$b_1$  adalah perubahan nilai rata-ran Y  
untuk setiap perubahan 1 satuan X.

*numerik*  
*numerik*

# Pendugaan Parameter

- Misalkan  $\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$  adalah prediksi untuk  $Y$  berdasarkan nilai ke- $i$  peubah  $X$  (dengan  $i = 1, 2, 3, \dots, n$ )

$$y = (\hat{\beta}_0 + \hat{\beta}_1 x) + e$$

- Maka residual ke- $i$  didefinisikan oleh:

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$$

$$e_i = y_i - \hat{y}_i \rightarrow e_i = y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i \quad \checkmark$$

$$e = y - \hat{y}$$

- JKG (Jumlah Kuadrat Galat) didefinisikan oleh:

$$= y - (\hat{\beta}_0 + \hat{\beta}_1 x)$$

$$JKG = e_1^2 + e_2^2 + \dots + e_n^2 \rightarrow \sum e_i^2$$

$$e = y - \hat{\beta}_0 - \hat{\beta}_1 x$$

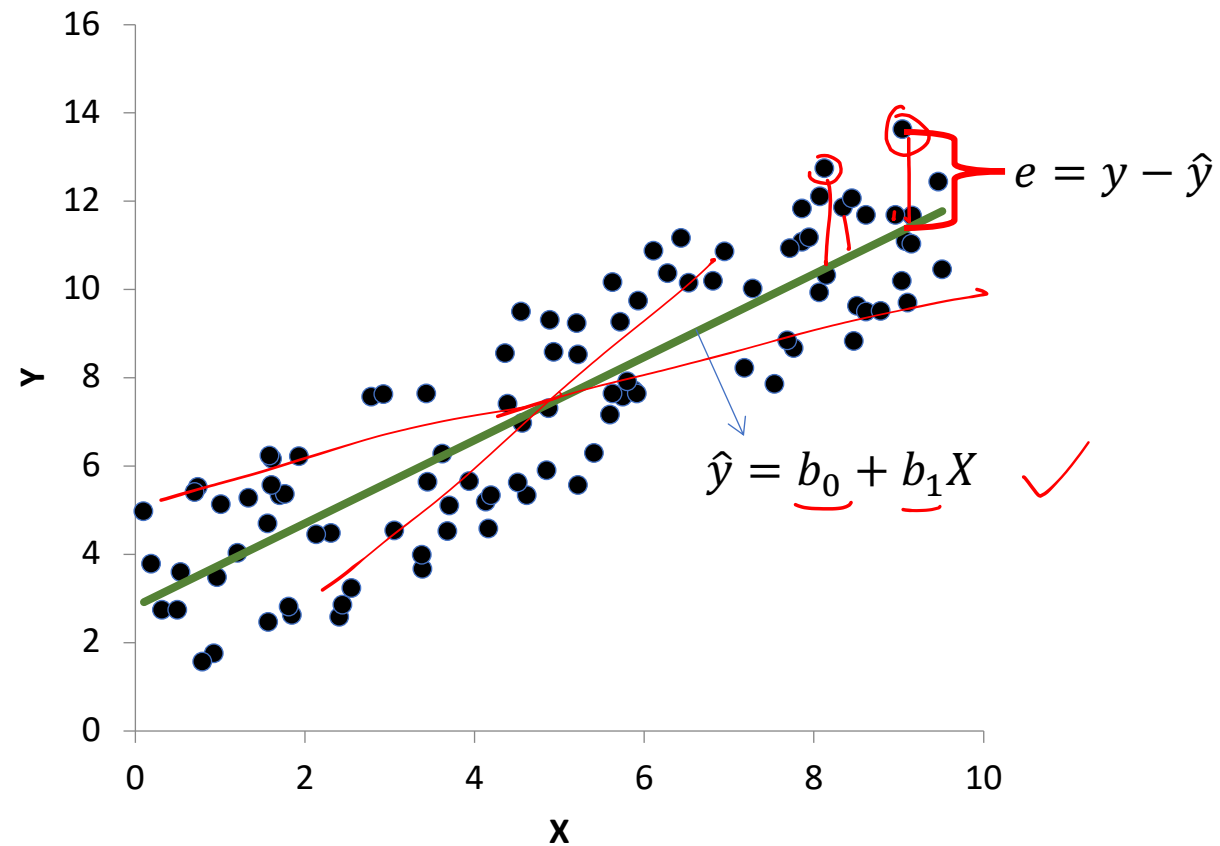
$$JKG = (y_1 - \hat{\beta}_0 - \hat{\beta}_1 x_1)^2 + \dots + (y_n - \hat{\beta}_0 - \hat{\beta}_1 x_n)^2$$

- Penduga MKT (Metode Kuadrat Terkecil), memilih  $\hat{\beta}_0$  dan  $\hat{\beta}_1$  yang meminimumkan  $JKG$ . Dengan perhitungan kalkulus diperoleh:

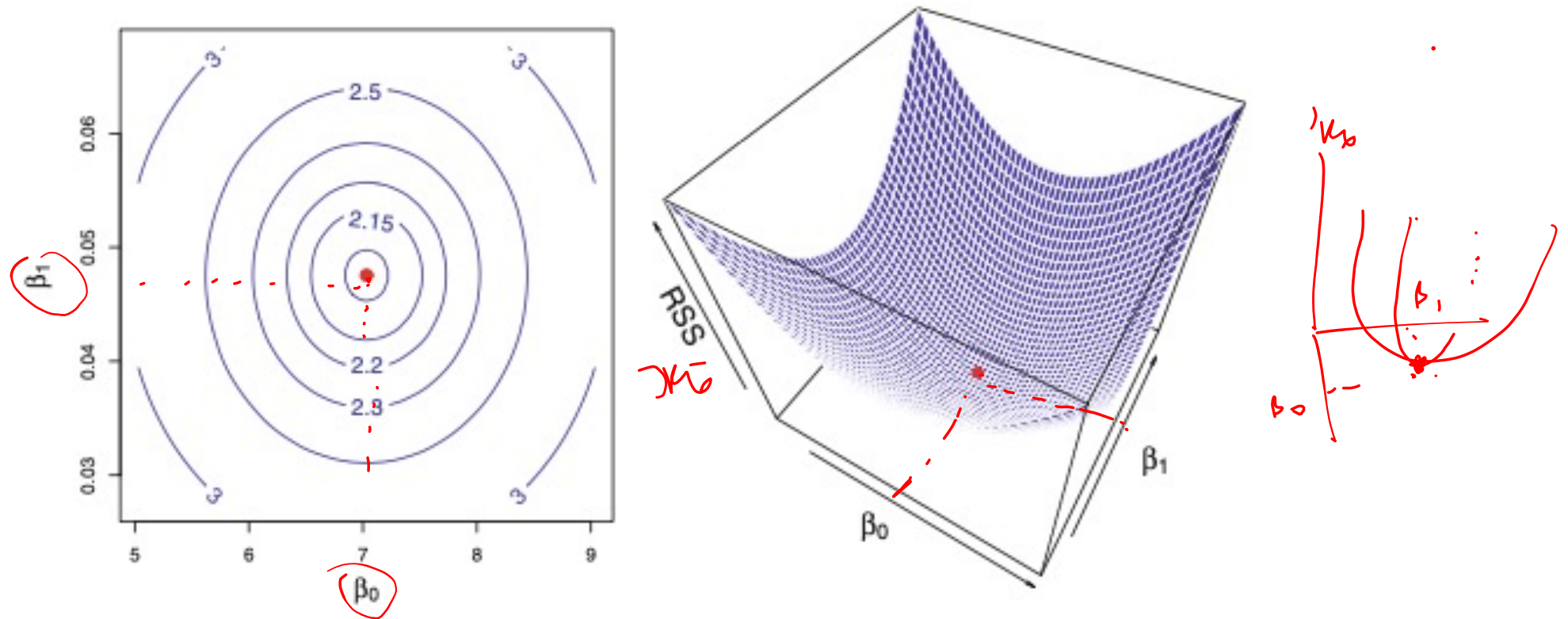
$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}; \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} \quad \checkmark$$

# Pendugaan Parameter

Metode Kuadrat Terkecil → Meminimumkan jumlah kuadrat galat



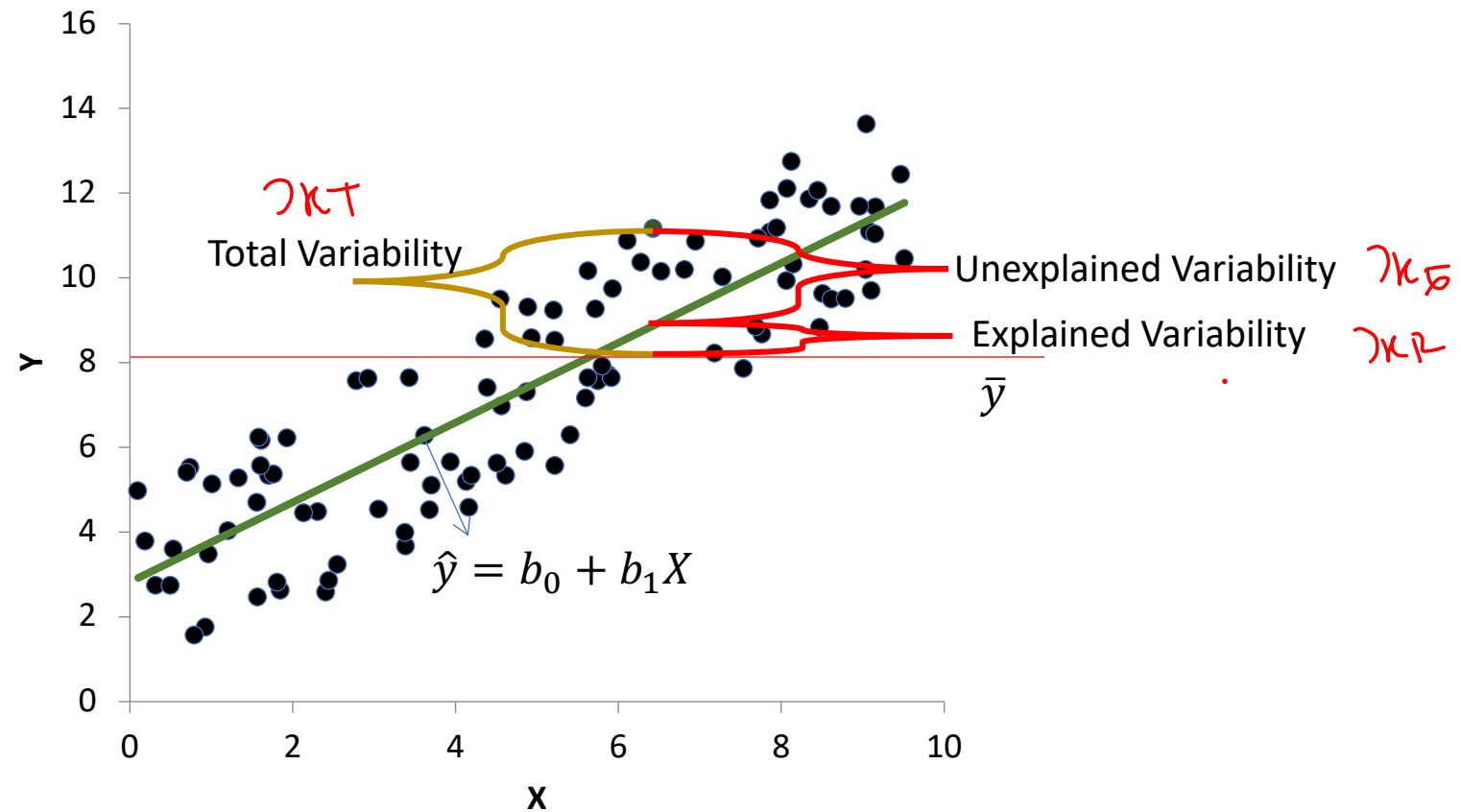
Ilustrasi kontur dan plot 3D pada JKG (RSS) untuk model dengan Y = sales dan X = TV



**FIGURE 3.2.** Contour and three-dimensional plots of the RSS on the Advertising data, using sales as the response and TV as the predictor. The red dots correspond to the least squares estimates  $\hat{\beta}_0$  and  $\hat{\beta}_1$ , given by (3.4).

# Pendugaan Parameter

Keragaman yang dapat dijelaskan dan tidak dapat dijelaskan



# Regresi Linier Berganda

- Analisis regresi linear berganda:
  - Secara umum, kita memodelkan peubah respon  $Y$  sebagai fungsi linier dari  $k$  peubah prediktor ( $X$ ) sebagai:

$$Y = \beta_0 + \beta_1 X_1 + \cdots + \beta_k X_k + \varepsilon$$

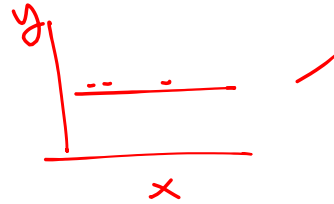
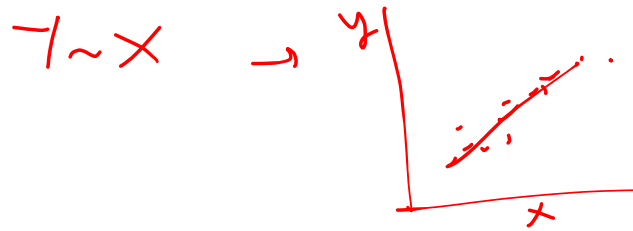
$n = \text{banyak pengamatan}$

- Atau dalam notasi matriks  $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$   $\rightarrow n \times 1$

$$\begin{matrix} \mathbf{y} \\ n \times 1 \end{matrix} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix} \quad \mathbf{X} \begin{matrix} n \times (k+1) \end{matrix} = \begin{pmatrix} 1 & x_{11} & \cdots & x_{k1} \\ 1 & x_{12} & & x_{k2} \\ 1 & x_{13} & & x_{k3} \\ \vdots & \vdots & & \vdots \\ 1 & x_{1n} & & x_{kn} \end{pmatrix} \quad \boldsymbol{\beta} \begin{matrix} (k+1) \times 1 \end{matrix} = \begin{pmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \vdots \\ \beta_k \end{pmatrix} \quad \boldsymbol{\varepsilon} = \begin{pmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_n \end{pmatrix}$$

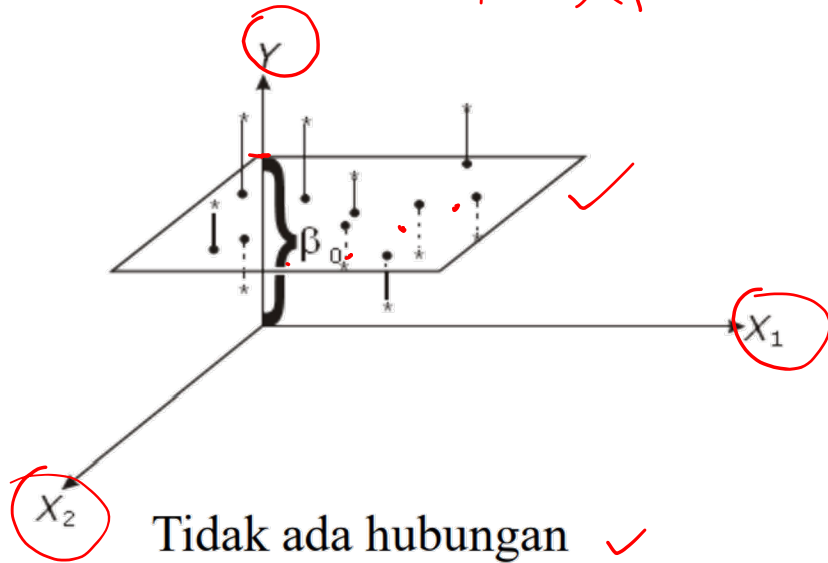


reg linear sederhana

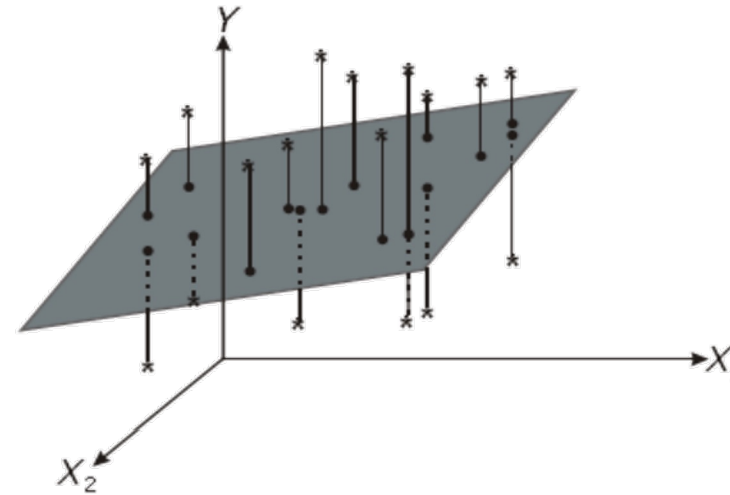


- Jika kita memiliki dua variabel X, model dapat diilustrasikan sebagai berikut

reg linear berganda  $Y \sim X_1 + X_2$



Tidak ada hubungan ✓



Ada hubungan ✓

- Pendugaan koefisien regresi:

- Pendugaan koefisien regresi diperoleh dengan meminimumkan jumlah kuadrat galat (residual)  $\rightarrow$  OLS (Ordinary Least Square) atau MKT (Metode Kuadrat Terkecil)

- Dalam hal ini dicari dugaan dari  $\beta_j, j = 0, 1, 2, \dots, k$  yang meminimumkan  $\sum_i \varepsilon^2$ , dengan  $\varepsilon = Y - \hat{Y}$ , yang dalam notasi matriks diperoleh

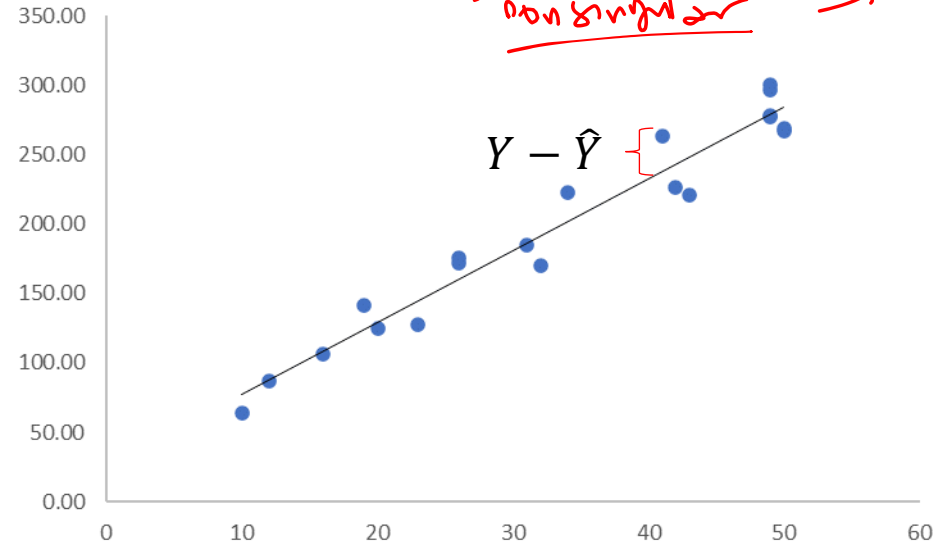
$$\hat{\beta} = (X'X)^{-1}X'y$$

— korelasi tinggi

$X_1$	$X_2$
1	5
2	6
3	7
4	8

$X_1 = X_2 - 4$

$X_1$  bergantung  
dgn  $X_2$



## Asumsi model regresi linear

Nilai mean dari peubah  $Y$  dimodelkan secara akurat oleh fungsi linier dari peubah-peubah  $X$

$$Y = \beta_0 + \beta_1 X_1 + \cdots + \beta_k X_k + \varepsilon$$

Antar peubah  $X$  tidak ada multikolinearitas

Galat acak diasumsikan menyebar normal dengan nilai tengah nol dan memiliki ragam yang konstan  $\sigma^2$  (ragam homogen)

$$\varepsilon \sim N(0, \sigma^2)$$

$$E(\varepsilon) = 0$$

$$\text{Var}(\varepsilon) = \sigma^2$$

$$\text{cov}(\varepsilon_i, \varepsilon_j) = 0$$

Galat bersifat independen/  
saling bebas  
(tidak ada autokorelasi) ✓

# Kesesuaian Model

- Kualitas kecocokan regresi linier biasanya dinilai menggunakan dua besaran terkait: Galat Baku Residual (residual standard error) dan statistik  $R^2$

- **Galat Baku Residual**

- Galat Baku Residual merupakan dugaan simpangan baku dari residual, yakni jumlah rata-rata respon yang akan menyimpang dari garis regresi yang sebenarnya.

$$\text{Galat Baku Residual} = \sqrt{\frac{1}{n-p-1} JKG} = \sqrt{\frac{1}{n-p-1} \sum_{i=1}^n (y_i - \hat{y}_i)^2}$$

- Galat Baku Residual dianggap sebagai ukuran kecocokan model dengan data.
  - Jika prediksi yang diperoleh dengan menggunakan model sangat dekat dengan nilai hasil sebenarnya—yaitu, jika  $\hat{y}_i \approx y_i$  untuk  $i = 1, \dots, n$ —maka Galat Baku Residual akan menjadi kecil, dan kita dapat menyimpulkan bahwa model tersebut sangat cocok dengan data.
  - Di sisi lain, jika  $\hat{y}_i$  sangat jauh dari  $y_i$  untuk satu atau lebih pengamatan, maka Galat Baku Residual mungkin cukup besar, menunjukkan bahwa model tidak sesuai dengan data dengan baik.

- **Statistik  $R^2$**

- Galat Baku Residual memberikan ukuran mutlak ketidaksesuaian model dengan data.
- Tetapi karena diukur dalam satuan  $Y$ , tidak selalu jelas apa yang dimaksud dengan Galat Baku Residual yang baik.
- Statistik  $R^2$  memberikan alternatif ukuran kecocokan model.
- Bentuknya berupa proporsi (proporsi ragam yang dijelaskan) sehingga selalu mengambil nilai antara 0 dan 1, dan tidak bergantung pada skala  $Y$ .

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

- $R^2$  mengukur proporsi keragaman dalam  $Y$  yang dapat dijelaskan dengan menggunakan  $X$ .

# Dummy Variable (Peubah Boneka)

- Dummy variable diterapkan pada peubah prediktor dengan skala kategorik
- Banyaknya dummy variable yang dibentuk dari satu peubah kategorik adalah  $k-1$ , dengan  $k$  adalah banyaknya kategori dalam peubah tersebut.   
  $\times (1, 2, 3) \rightarrow k=3$   
 $\rightarrow dx_2 \quad dx_3$
- Misalkan pada peubah pekerjaan yang terdiri dari 3 kategori: Pegawai BUMN, Pegawai swasta, dan PNS
- Pilih satu kategori sebagai reference, misalkan Pegawai BUMN

Kategori	D1	D2
Pegawai BUMN	0	0
Pegawai swasta	1	0
PNS	0	1

$X \rightarrow \begin{pmatrix} X & B & C & D \end{pmatrix}$

	$D_1$	$D_2$	$D_3$
1. A	0	0	0
2. B	1	0	0
3. C	0	1	0
4. D	0	0	1
5. B	1	0	0
6. C	0	1	0
7. A	0	0	0

# Ilustrasi Kasus

Misalkan telah diperoleh data scraping harga rumah beserta karakteristiknya

Ingin diketahui pengaruh setiap karakteristiknya ke harga rumah tersebut.

# Data: RUMAH.TXT

- ID, Identification number
- **sales\_price**, Sales price of residence (dollars) ✓
- **X1**, Finished area of residence (square feet) ✓
- X2, Total number of bedrooms in residence
- X3, Total number of bathrooms in residence
- X4, Presence or absence of air conditioning: 1 if yes; 0 otherwise
- X5, Number of cars that garage will hold
- X6, Presence or absence of swimming pool: 1 if yes; 0 otherwise
- **X7**, Year property was originally constructed ✓  $\text{usia} = 2023 - X7$
- **X8**, Index for quality of construction (1. Indicates high quality; 2. Indicates medium quality; 3. Indicates low quality) ✓  $\rightarrow$  *ordinal dummy variable* ✓
- X9, Qualitative indicator of architectural style
- X10, Lot size (square feet)
- X11, Presence or absence of adjacency to highway: 1 if yes; 0 otherwise

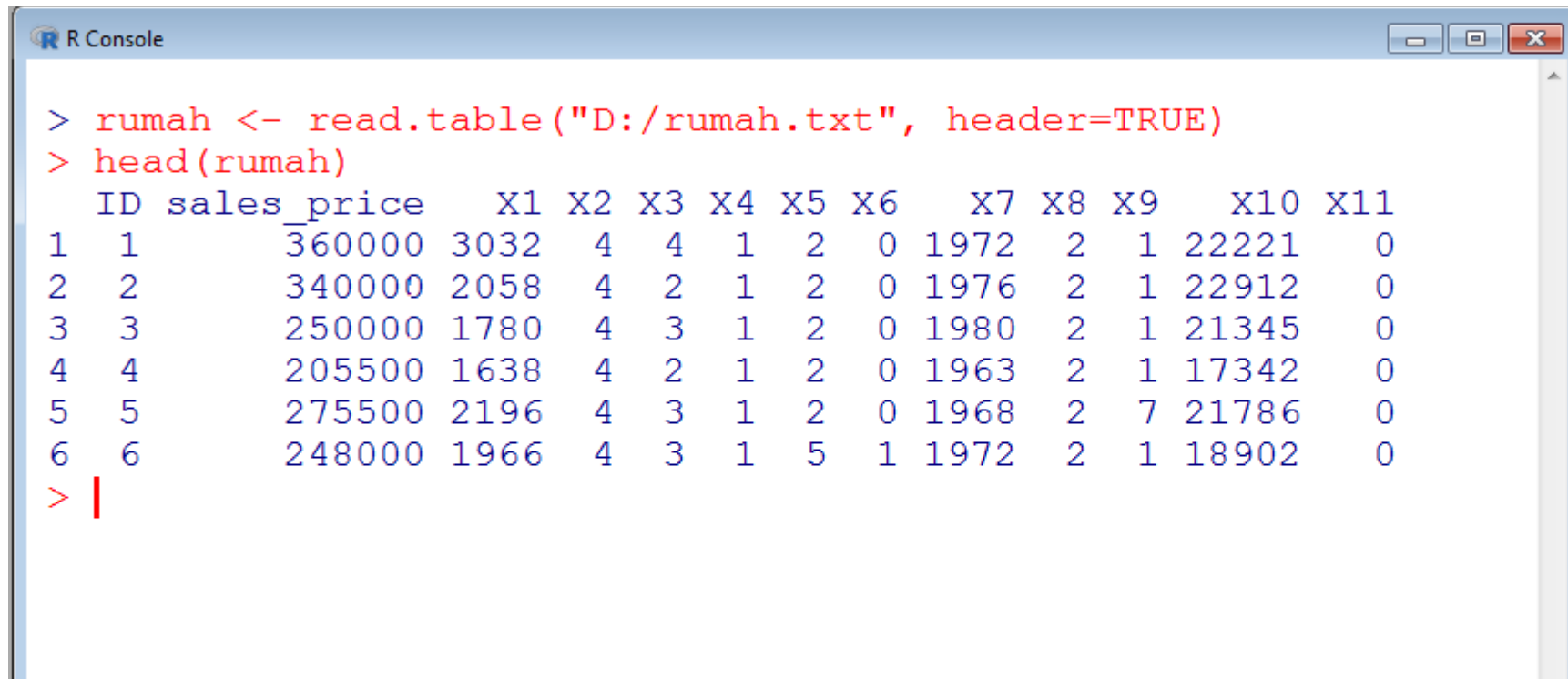


```
rumah <- read.table("D:/rumah.txt", header=TRUE)
head(rumah)
```

Penjelasan:

Baris #1: membaca file TXT dengan nama "rumah.txt" yang tersimpan pada folder D. Opsi "header=TRUE" mengatakan bahwa baris pertama pada file yang dibaca merupakan nama-nama kolom. Hasil pembacaan file tersebut disimpan dalam bentuk data frame di R dengan nama "rumah"

Baris #2: mencetak 6 (enam) baris pertama dari data frame "rumah"

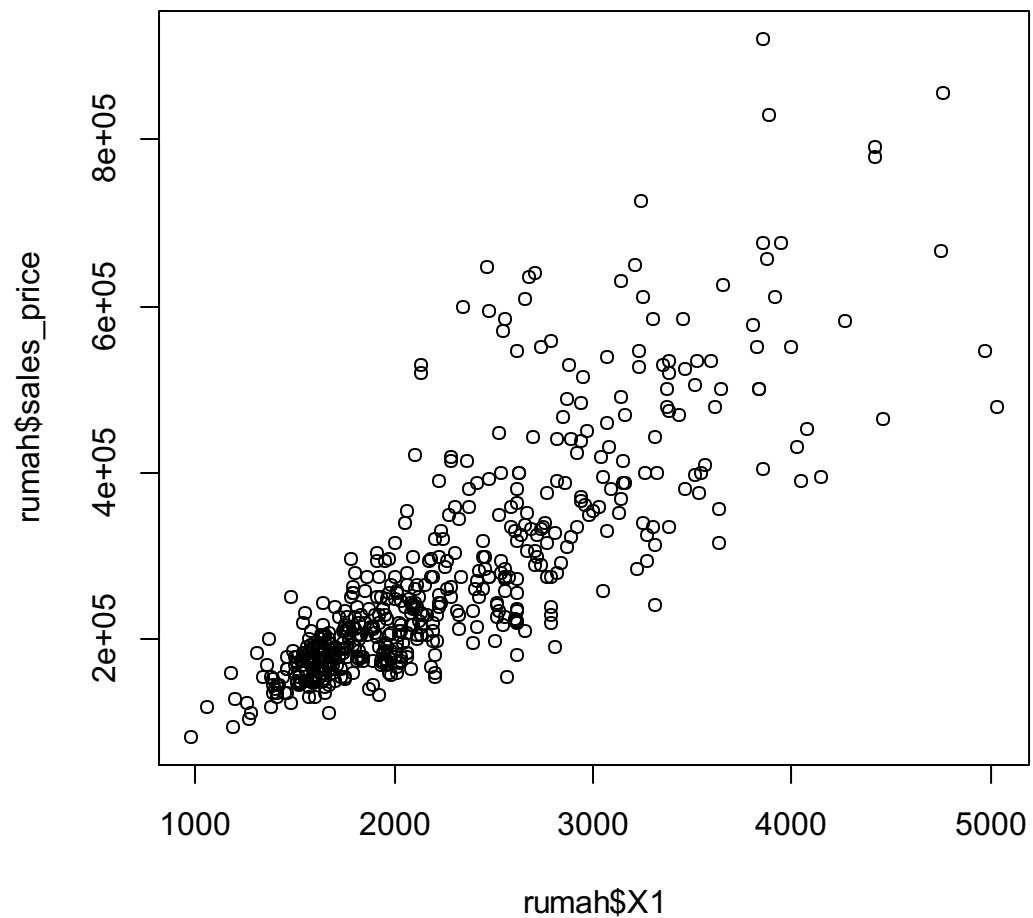


```
> rumah <- read.table("D:/rumah.txt", header=TRUE)
> head(rumah)
```

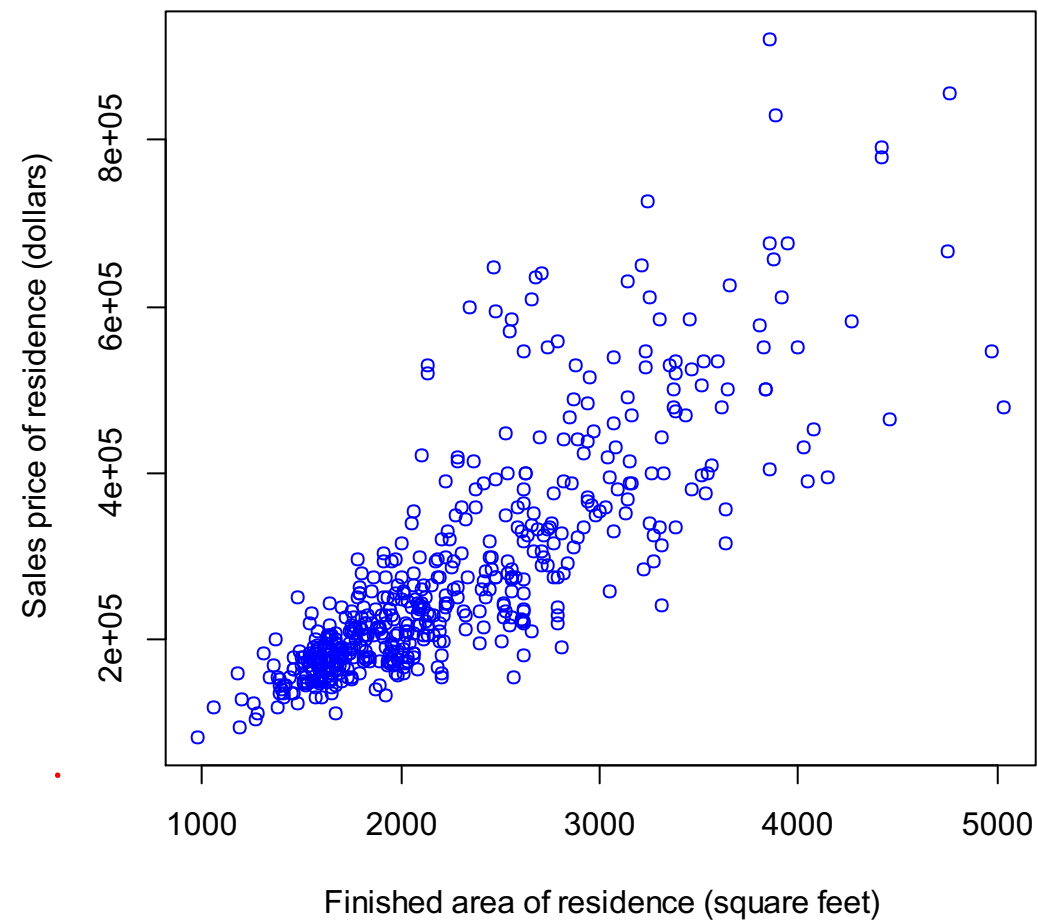
	ID	sales_price	X1	X2	X3	X4	X5	X6	X7	X8	X9	X10	X11
1	1	360000	3032	4	4	1	2	0	1972	2	1	22221	0
2	2	340000	2058	4	2	1	2	0	1976	2	1	22912	0
3	3	250000	1780	4	3	1	2	0	1980	2	1	21345	0
4	4	205500	1638	4	2	1	2	0	1963	2	1	17342	0
5	5	275500	2196	4	3	1	2	0	1968	2	7	21786	0
6	6	248000	1966	4	3	1	5	1	1972	2	1	18902	0

```
> |
```

```
plot(rumah$X1, rumah$sales_price)
```

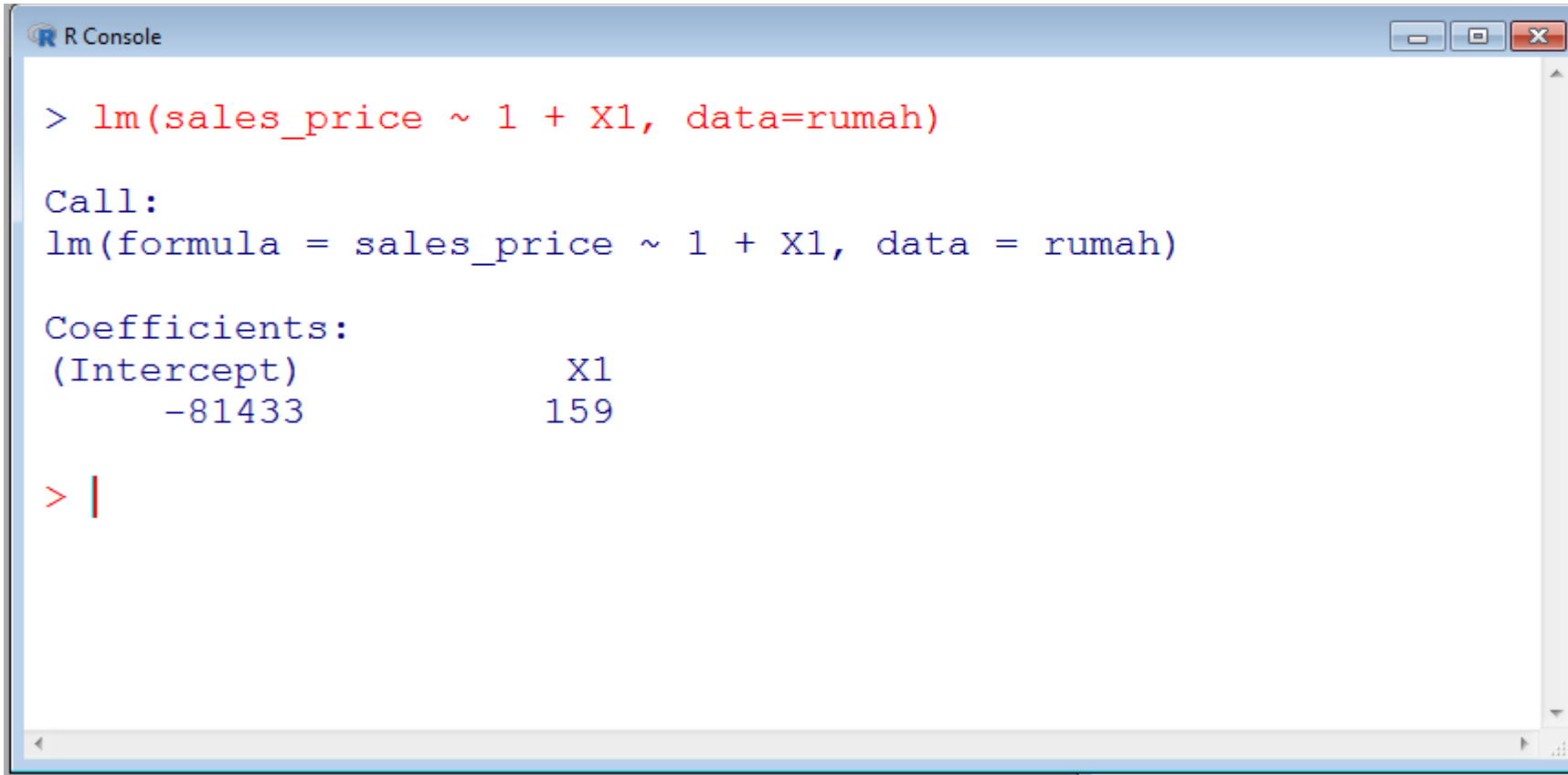


```
plot(rumah$X1, rumah$sales_price,  
     xlab="Finished area of residence (square feet)",  
     ylab="Sales price of residence (dollars)",  
     col="blue")
```



## Model Regresi Linier Sederhana: $Y \sim X_1$

```
lm(sales_price ~ 1 + X1, data=rumah)
```



The image shows a screenshot of an R Console window. The window has a title bar that says "R Console" and standard Windows window controls (minimize, maximize, close). The console displays the following text:

```
> lm(sales_price ~ 1 + X1, data=rumah)
```

Call:

```
lm(formula = sales_price ~ 1 + X1, data = rumah)
```

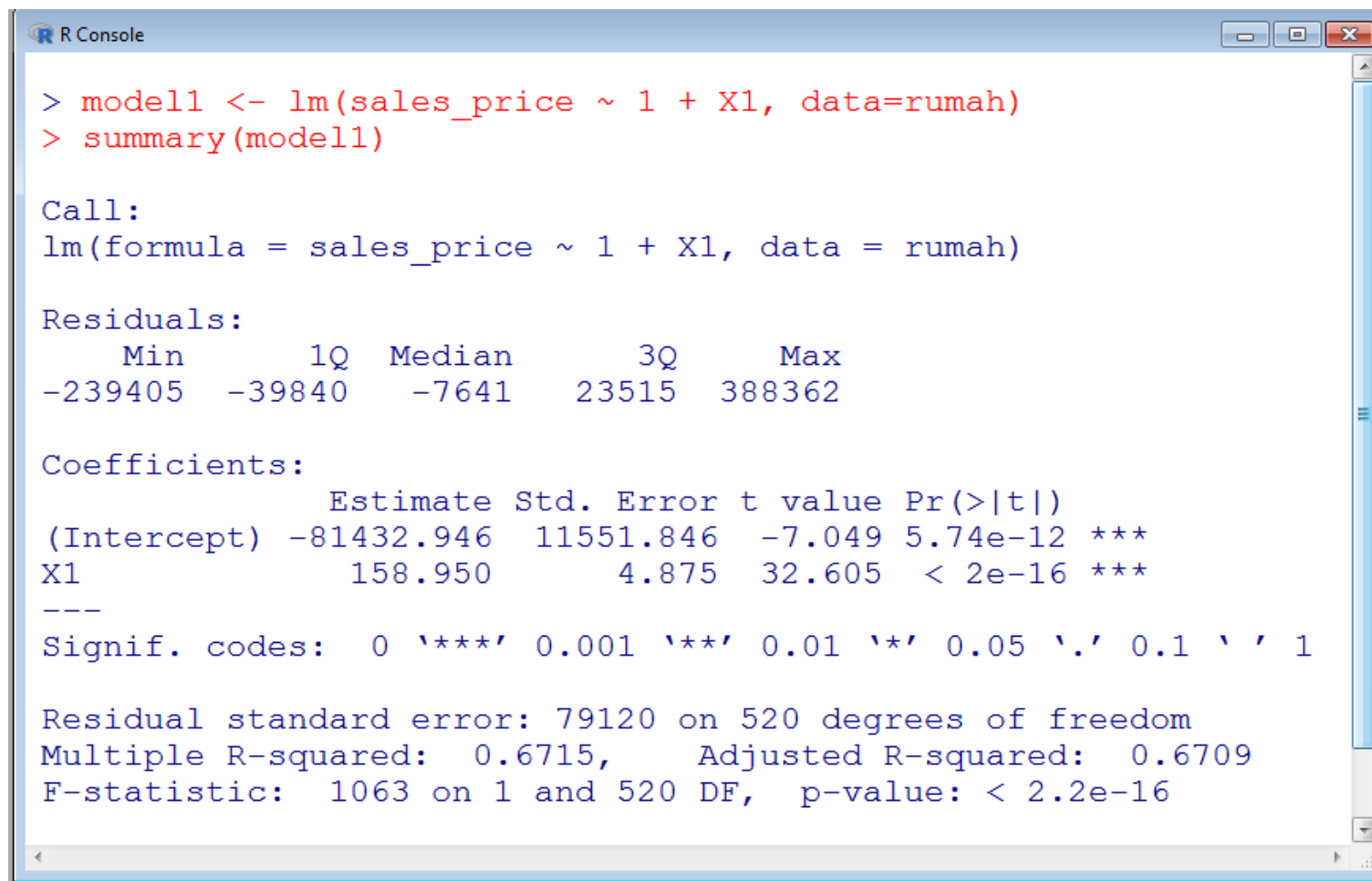
Coefficients:

(Intercept)	X1
-81433	159

```
> |
```

The console output shows the results of the linear regression model. The coefficients are displayed in a table format. The intercept is -81433 and the coefficient for X1 is 159. The prompt "> |" indicates that the user is ready to enter another command.

```
modell1 <- lm(sales_price ~ 1 + X1, data=rumah)
summary(modell1)
```

The image shows a screenshot of an R Console window. The window has a title bar that says "R Console" and standard Windows window controls (minimize, maximize, close). The console displays the following text:

```
> modell1 <- lm(sales_price ~ 1 + X1, data=rumah)
> summary(modell1)

Call:
lm(formula = sales_price ~ 1 + X1, data = rumah)

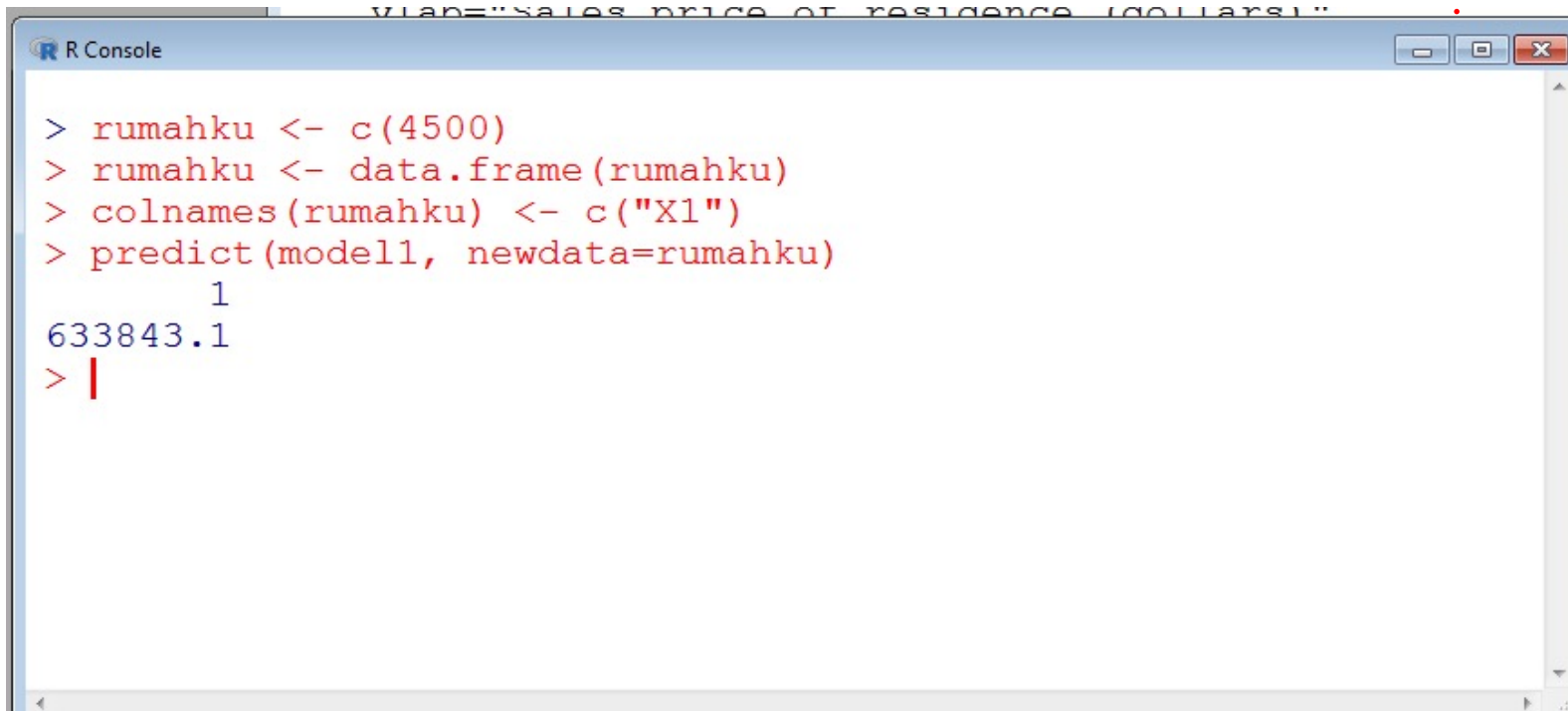
Residuals:
    Min       1Q   Median       3Q      Max
-239405  -39840   -7641   23515  388362

Coefficients:
              Estimate Std. Error t value Pr(>|t|)
(Intercept) -81432.946  11551.846   -7.049 5.74e-12 ***
X1           158.950     4.875   32.605 < 2e-16 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 79120 on 520 degrees of freedom
Multiple R-squared:  0.6715,    Adjusted R-squared:  0.6709
F-statistic: 1063 on 1 and 520 DF,  p-value: < 2.2e-16
```

## Memprediksi Harga Rumah Seluas 4500 ft<sup>2</sup>

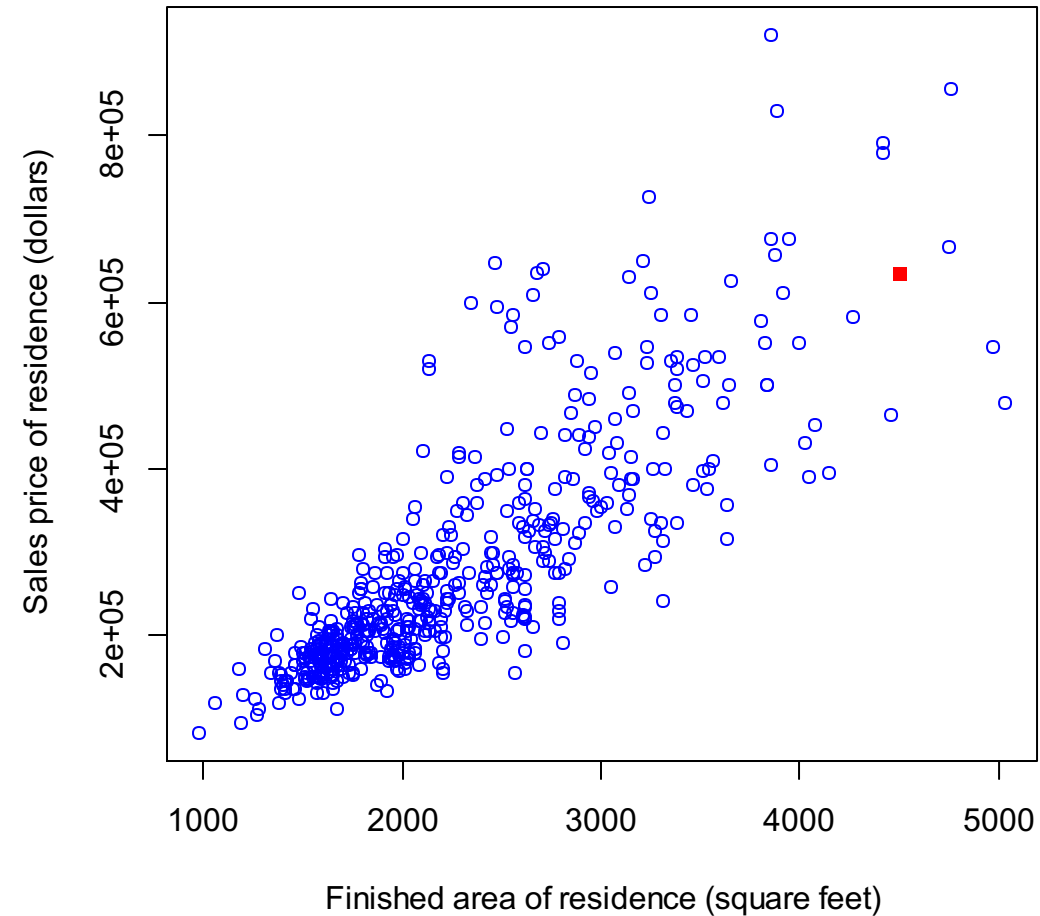
```
rumahku <- c(4500)
rumahku <- data.frame(rumahku)
colnames(rumahku) <- c("X1")
predict(modell, newdata=rumahku)
```



The screenshot shows an R Console window with the following code and output:

```
> rumahku <- c(4500)
> rumahku <- data.frame(rumahku)
> colnames(rumahku) <- c("X1")
> predict(modell, newdata=rumahku)
      1
633843.1
> |
```

```
plot(rumah$X1, rumah$sales_price,  
     xlab="Finished area of residence (square feet)",  
     ylab="Sales price of residence (dollars)",  
     col="blue")  
points(rumahku$X1, predict(model1, newdata=rumahku),  
       col="red", pch=15)
```



# Model Regresi Linier Berganda

```
rumah$umur = 2023 - rumah$X7
rumah$ind.med <- ifelse(rumah$X8==2,1,0)
rumah$ind.low <- ifelse(rumah$X8==3,1,0)
model2 <- lm(sales_price ~ 1 + X1 + umur + ind.med + ind.low, data=rumah)
summary(model2)
```

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )	
(Intercept)	2.585e+05	2.102e+04	12.30	< 2e-16	***
X1	9.756e+01	5.396e+00	18.08	< 2e-16	***
umur	-1.119e+03	1.953e+02	-5.73	1.7e-08	***
ind.med	-1.524e+05	1.040e+04	-14.65	< 2e-16	***
ind.low	-1.709e+05	1.404e+04	-12.17	< 2e-16	***

---

Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 61770 on 517 degrees of freedom

Multiple R-squared: 0.801, Adjusted R-squared: 0.7994

F-statistic: 520.1 on 4 and 517 DF, p-value: < 2.2e-16

```
coef(model2)
> coef(model2)
(Intercept)          X1          umur    ind.med    ind.low
258477.75900    97.55652 -1119.15626 -152382.80986 -170876.66457
```

Rumah indeks tinggi

## • Interpretasi

- dugaan model yang diperoleh

sales\_price = 258478 + 97.56 X1 - 1119.16 Umur - 152382.80 Indeks Medium - 170876.66 Indeks Low

## • UMUR

- Koef = -1119.16, artinya jika rumah bertambah tua satu tahun, harganya turun dengan rata-rata sebesar 1119.16 dollar (signifikan)

## • X1 (luas bangunan)

- Koef = 97.56, artinya jika luas rumah bertambah satu feet<sup>2</sup>, harganya naik dengan rata-rata sebesar 97.56 dollar (signifikan)

## • Indeks Medium

- Rumah dengan indeks medium memiliki rata-rata harga rumah 152382.80 dollar lebih rendah dari rumah dengan indeks tinggi (high) (signifikan)

## • Indeks Low

- Rumah dengan indeks medium memiliki rata-rata harga rumah 170876.66 dollar lebih rendah dari rumah dengan indeks tinggi (high) (signifikan)

→ 3 persamaan

① J. High → Int + X<sub>1</sub> + umur

② J. med → Int + X<sub>1</sub> + umur + Indeks

③ J. low → Int + X<sub>1</sub> + umur + low

Referensi



- atau

```
rumah$ind <- as.factor(rumah$X8)
str(rumah)
model3 <- lm(sales_price ~ 1 + X1 + umur + ind, data=rumah)
summary(model3)
coef(model3)
```

Langsung gunakan factor pada peubah kategorik  
Kategori yang menjadi reference selalu kategori yang pertama

```
> coef(model3)
```

(Intercept)	X1	umur	ind2	ind3
258477.75900	97.55652	-1119.15626	-152382.80986	-170876.66457

Lalu bagaimana jika reference kategori peubah X8 diganti menjadi indeks “Low”?

```
##referensi X8 diganti menjadi "Low"
rumah$ind.high <- ifelse(rumah$X8==1,1,0)
str(rumah)
model4 <- lm(sales_price ~ 1 + X1 + umur + ind.high + ind.med, data=rumah)
summary(model4)
coef(model4)
```

```
> coef(model4)
      (Intercept)           X1          umur      ind.high      ind.med
87601.09443      97.55652 -1119.15626 170876.66457 18493.85470
```

```
##dengan factor dari X8 (ind)
model5 <- lm(sales_price ~ 1 + X1 + umur + relevel(ind,ref=3), data=rumah)
summary(model5)
coef(model5)
```

```
> coef(model5)
      (Intercept)           X1          umur
87601.09443      97.55652 -1119.15626
relevel(ind, ref = 3)1 relevel(ind, ref = 3)2
170876.66457      18493.85470
```

Terima kasih 😊