

STA 141A – Fundamentals of Statistical Data Science

Department of Statistics; University of California, Davis

Instructor: Dr. Akira Horiguchi (ahoriguchi@ucdavis.edu)

Ao1 TA: Zhentao Li (ztlli@ucdavis.edu)

Ao2 TA: Zijie Tian (zjztian@ucdavis.edu)

Ao3 TA: Lingyou Pang (lyopang@ucdavis.edu)

Section 6: Regression Analysis with R

Spring 2025 (Mar 31 – Jun 05), MWF, 01:10 PM – 02:00 PM, Young 198

Based on Chapter 3 of ISL book James et al. (2021).

- For more R code examples, see R Markdown files in <https://www.statlearning.com/resources-second-edition>

Section 6: Regression Analysis with R

- Linear Regression
- Idea of polynomial regression

SECTION 6: REGRESSION ANALYSIS WITH R

LINEAR REGRESSION

AN EXAMPLE – 1

Consider the data set in `Advertising.csv` consisting of the sales of a product in 200 different markets, with advertising budgets for the product in each of those markets for three different media: TV, Radio, Newspaper.

```
1 adv <- read.csv("Advertising.csv", row.names="X")  
2 str(adv)
```

```
> str(adv)  
'data.frame':  200 obs. of  4 variables:  
 $ TV      : num  230.1 44.5 17.2 151.5 180.8 ...  
 $ Radio   : num  37.8 39.3 45.9 41.3 10.8 48.9 32.8 19.6 2.1 2.6 ...  
 $ Newspaper: num  69.2 45.1 69.3 58.5 58.4 75 23.5 11.6 1 21.2 ...  
 $ Sales   : num  22.1 10.4 9.3 18.5 12.9 7.2 11.8 13.2 4.8 10.6 ...
```

AN EXAMPLE – 2

We want to investigate the relationship between Sales and the total budget spent for advertisement on TV, Radio, and Newspaper.

- Then, we sum row-wise, but exclude the last column (which is the 4th column after we deleted the 1st column).

```
1 adv$Budget <- rowSums(adv[, -4])  
2 str(adv)
```

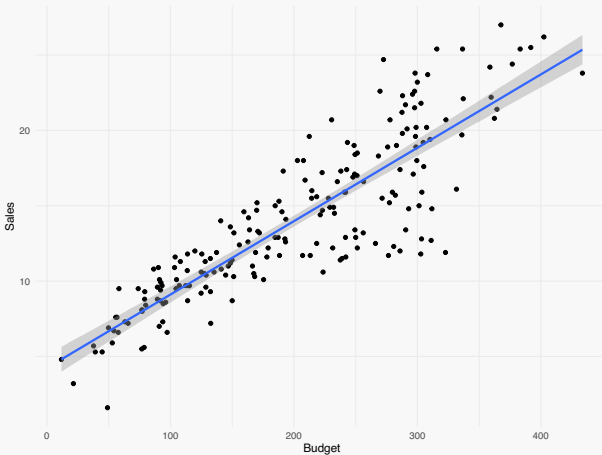
```
> str(adv)  
'data.frame': 200 obs. of 5 variables:  
 $ TV : num 230.1 44.5 17.2 151.5 180.8 ...  
 $ Radio : num 37.8 39.3 45.9 41.3 10.8 48.9 32.8 19.6 2.1 2.6 ...  
 $ Newspaper: num 69.2 45.1 69.3 58.5 58.4 75 23.5 11.6 1 21.2 ...  
 $ Sales : num 22.1 10.4 9.3 18.5 12.9 7.2 11.8 13.2 4.8 10.6 ...  
 $ Budget : num 337 129 132 251 250 ...
```

Reasonable research questions for this data set:

- Is there a relationship between Budget and Sales?
- If there is a relationship, is it linear?
- How strong is the relationship between Budget and Sales?
- Which media contribute to sales?
- How accurately can we estimate the effect of each medium on sales?
- How accurately can we predict future sales?

AN EXAMPLE – 4

```
1 ggplot(adv, aes(Budget, Sales)) +  
2   geom_point() +  
3   geom_smooth(method="lm") +  
4   theme_minimal()
```



LINEAR REGRESSION – MATRIX REPRESENTATION

Recall: want to estimate the relationship between Y and X_1, \dots, X_p .

- Linear regression simplifies the task of estimating the “true” relationship to the task of estimating $p + 1$ parameters $\beta_0, \beta_1, \dots, \beta_p$:

$$\underline{Y} = \underline{\beta}_0 + \underline{\beta}_1 X_1 + \dots + \underline{\beta}_p X_p + \underline{\varepsilon} \quad \text{-- “epsilon”} \quad (1)$$

where the error term ε is a catch-all for what is missed by this model.

- For n observations

$$(\underline{x}_1, y_1) = (X_{11}, X_{12}, \dots, X_{1p}, y_1),$$

$$(\underline{x}_2, y_2) = (X_{21}, X_{22}, \dots, X_{2p}, y_2),$$

$$\vdots$$

$$(\underline{x}_n, y_n) = (X_{n1}, X_{n2}, \dots, X_{np}, y_n),$$

there is the following matrix representation

$$y_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip} + \varepsilon_i, \quad \text{for } i = 1, 2, \dots, n \quad (2)$$

$$\Leftrightarrow \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix} = \begin{pmatrix} 1 & X_{11} & X_{12} & \cdots & X_{1p} \\ 1 & X_{21} & X_{22} & \cdots & X_{2p} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & X_{n1} & X_{n2} & \cdots & X_{np} \end{pmatrix} \begin{pmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_p \end{pmatrix} + \begin{pmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_n \end{pmatrix} \quad (3)$$

$$\Leftrightarrow \mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon} \quad (4)$$

ORDINARY LEAST SQUARES (OLS)

A “good” estimator for the regression parameters $\beta = (\beta_0, \dots, \beta_p)^T \in \mathbb{R}^{p+1}$ produces small residuals.

- A *residual* is the difference between a response and its predicted value.
- The *i*th *residual* is $y_i - \hat{y}_i$ for $i = 1, \dots, n$.
- β can be estimated by using the *Ordinary Least Squares* (OLS) method.
- The *OLS estimator* for β is defined to be the vector that minimizes the residual sum of squares:

$$\hat{\beta}_{OLS} := \operatorname{argmin}_{\beta \in \mathbb{R}^{p+1}} \underbrace{\|\mathbf{y} - \mathbf{X}\beta\|^2}_{\text{residual sum of squares}} \quad (5)$$

where for any $\mathbf{z} = (z_1, \dots, z_{p+1})^T \in \mathbb{R}^{p+1}$ holds $\|\mathbf{z}\|^2 = z_1^2 + \dots + z_{p+1}^2$.

- It can be shown that $\hat{\beta}_{OLS} = \underbrace{(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}}$ (if the inverse $(\mathbf{X}^T \mathbf{X})^{-1}$ exists).

The OLS estimator $\hat{\beta}_{OLS} \in \mathbb{R}^{p+1}$ produces:

- $p = 0$: sample mean of the y_1, \dots, y_n (best estimate w/o predictor info)
- $p = 1$: “line of best fit”
- $p = 2$: “plane of best fit”
- $p \geq 3$: “hyperplane of best fit”

Given a line/plane/hyperplane of best fit:

- Residual is vertical displacement between point and line/plane/hyperplane

ORDINARY LEAST SQUARES (OLS) – EXAMPLE LINE OF BEST FIT

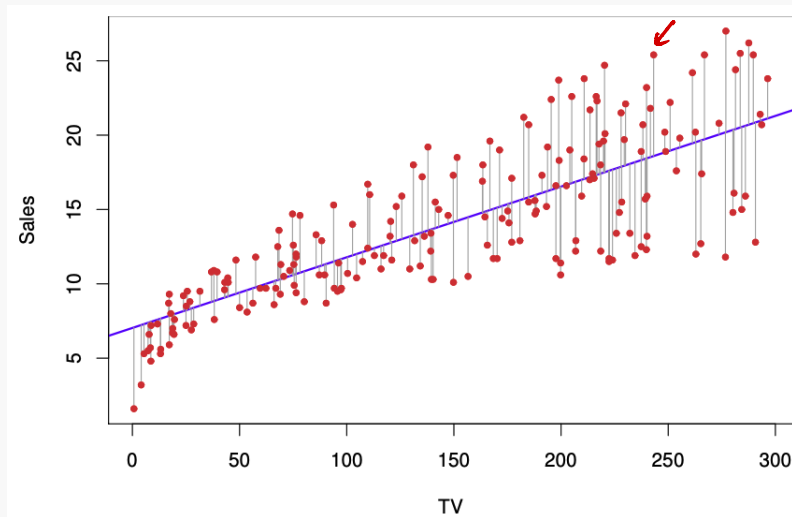


Figure 1: Image by James et al. (2021). “For the Advertising data, the least squares fit for the regression of sales onto TV is shown.” The line is chosen to minimize the sum of the squared vertical distances between each observation (shown in red) and the line.

ORDINARY LEAST SQUARES (OLS) – EXAMPLE PLANE OF BEST FIT

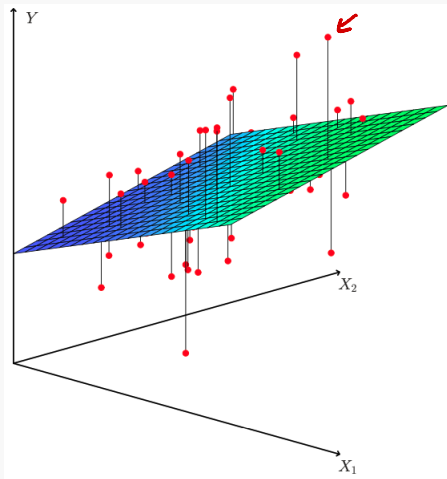


Figure 2: Image by James et al. (2021). “In a three-dimensional setting, with two predictors and one response, the least squares regression line becomes a plane. The plane is chosen to minimize the sum of the squared vertical distances between each observation (shown in red) and the plane.”

End of May 2 lecture

Questions of interest:

1. Is at least one of the predictors X_1, X_2, \dots, X_p useful in predicting the response Y ?
2. Do all the predictors help to explain Y , or is only a subset of the predictors useful?
3. How well does the model fit the data?
4. Given a set of predictor values, what response value should we predict, and how accurate is our prediction?

Questions apply to regression generally, but answers might be specific to linear regression.

1. IS THERE A RELATIONSHIP BETWEEN THE RESPONSE AND PREDICTORS?

Are all regression coefficients equal to zero?

- What would this imply about the relationship?
- One can use the hypothesis test

$$H_0: \beta_1 = \beta_2 = \cdots \beta_p = 0 \quad \text{vs.} \quad H_a: \beta_j \neq 0 \text{ for at least one } j, \quad (6)$$

for which the following F -statistic is needed:

$$F = \frac{(TSS - RSS)/p}{RSS/(n - p - 1)} \geq 1. \quad (7)$$

- If H_0 is true, we might expect $F \approx 1$.
- If H_a is true, we might expect $F > 1$.

1. IS THERE A RELATIONSHIP BETWEEN THE RESPONSE AND PREDICTORS?

We might test whether one specific regression coefficient is zero or not.

- For a specific $j = 0, \dots, p$, one can use the hypothesis test

$$H_0: \beta_j = 0 \quad \text{vs.} \quad H_a: \beta_j \neq 0 \quad (8)$$

for which the following t -statistic is needed:

$$t = \frac{\hat{\beta}_j}{SE(\hat{\beta}_j)} . \quad (9)$$

- $SE(\hat{\beta}_j)$ is the *standard error of $\hat{\beta}_j$* .
- Under H_0 holds $t \sim t_{n-p-1}$, where t_{n-p-1} is the *Student's t -distribution*.

2. DECIDING ON IMPORTANT VARIABLES

Variable selection: the task of determining which predictors are associated with the response.

1. A model contains a subset of the predictors.
2. For p predictors, there are 2^p possible models.
3. Many ways to choose a model in linear regression.
4. Outside linear regression, this is still an active research area!

Potential problems: i) Non-linearity

- Linear regression assumes a linear relationship between the predictors and the response.
- Residual plots can be used to detect non-linearity: If no pattern is visible, linearity is a reasonable assumption, otherwise not.
- If there are non-linear associations, a simple approach is to check whether non-linear transformations of the predictors help, such as $\log(X)$, \sqrt{X} or X^2 .

Potential problems: ii) Correlation of the error terms

- The errors are assumed to be uncorrelated.
- If the errors are correlated, the estimated standard errors tend to underestimate the true standard errors.
- Correlations frequently occur in the context of time series, where observations are analyzed over time, e.g. daily temperatures.

Potential problems: iii) Non-constant variances of the error terms (Heteroskedasticity)

- The errors are assumed to be homoskedastic (meaning their variances are constant across observations).
- The standard errors, confidence intervals, and hypothesis tests associated with the linear model rely on this assumption.
- One can identify non-constant variances in the errors (*heteroskedasticity*) from the presence of a funnel shape in the residual plot.
- One possible solution is to transform the response Y by using a concave function such as $\log(Y)$ or \sqrt{Y} .

Potential problems: iv) Outliers

- An *outlier* is a point which is far from the value predicted by the model.
- Outliers can arise for a variety of reasons, e.g., incorrect recording of an observation during data collection.
- Outliers might inflate the variance, the RSE, and other measures.

3. MODEL FIT

Potential problems: v) Collinearity

- *Collinearity* refers to the situation in which two or more predictors are closely related, resulting in uncertainty in the coefficient estimates, and thus in the standard error for $\hat{\beta}_j$ to grow.
- Collinearity can also be present between more than two variables (*multicollinearity*), even if no pair of variables are closely related.
- The *variance inflation factor* (VIF) quantifies the severity of (multi)collinearity: It measures how much the variance of an estimated regression coefficient is increased due to collinearity.
- VIF for each variable can be computed by

$$VIF(\hat{\beta}_j) = \frac{1}{1 - R_j^2},$$

where R_j^2 is the R^2 from a regression of X_j onto all other predictors.

- As a rule of thumb, VIF exceeding 5 indicates a large amount of collinearity, then we should ...
 1. ... drop the j th predictor;
 2. ... or combine the j th with other collinear predictors together into a single predictor.

Consider the linear model with $n > p + 1$, where the vectors of the predictors $X_1, \dots, X_n \in \mathbb{R}^p$ are linearly independent.

- In the linear model the OLS estimator $\hat{\beta}_{OLS}$ for β is the *BLUE* (Best Linear Unbiased Estimator) if $E(\varepsilon) = \mathbf{0}_n$ and $\text{Var}(\varepsilon) = \sigma^2 \mathbf{I}_n$, where $\sigma^2 > 0$ is the variance, and \mathbf{I}_n is the $n \times n$ identity matrix.
- Usually, the variance σ^2 is unknown and has to be estimated. An unbiased estimator for σ^2 is given by

$$\hat{\sigma}^2 := \frac{1}{n-p-1} \|\mathbf{y} - \mathbf{X}\hat{\beta}\|^2 \quad (10)$$

- For $p = 0$, we can use the sample variance s^2 in R (command `(var)`), which is for an independent sample Y_1, \dots, Y_n with $\bar{Y} := \frac{1}{n} \sum_{i=1}^n Y_i$ defined by

$$\hat{\sigma}^2 := \frac{1}{n-1} \sum_{i=1}^n (Y_i - \bar{Y})^2 \quad (11)$$

Exercise: Show that $\hat{\sigma}^2$ is an unbiased estimator for the variance σ^2 in our model, i.e. show that $E(\hat{\sigma}^2) = \sigma^2$.

- If $\hat{\sigma}$ is almost surely (i.e. with prob. 1) non-constant ($\hat{\sigma}$ is almost surely constant if all X_i equal to \bar{X} with prob. 1), Jensen's inequality gives

$$E(\hat{\sigma}) < \sqrt{E(\hat{\sigma}^2)} = \sqrt{\sigma^2} = \sigma. \quad (12)$$

This means that $\hat{\sigma}$ is not an unbiased estimator for σ , although $\hat{\sigma}^2$ is an unbiased estimator for σ^2 .

Exercise: Illustrate the bias of $\hat{\sigma}$ for σ , that is $\hat{\sigma} - \sigma$:

1. Generate a random vector x of length $n = 10$ with i.i.d. $N(0, 4)$ -distributed entries
2. Calculate the bias $\hat{\sigma} - \sigma$
3. Repeat the steps $nsim = 10,000$ times and report the average bias.

We consider the linear model with $p = 1$ and analyze the advertising data set.

```
1 adv <- read.csv("Advertising.csv")
2 fit <- lm(Sales ~ Newspaper, data = adv)
3 summary(fit)
```

```
Call:
lm(formula = Sales ~ Newspaper, data = adv)

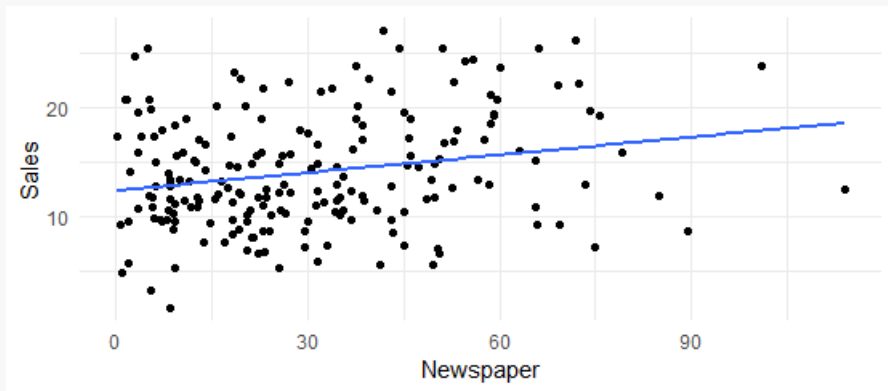
Residuals:
    Min       1Q   Median       3Q      Max
-11.2272  -3.3873  -0.8392   3.5059  12.7751

Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept)  12.35141    0.62142   19.88 < 2e-16 ***
Newspaper     0.05469    0.01658    3.30  0.00115 **
---
Signif. codes:
  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 5.092 on 198 degrees of freedom
Multiple R-squared:  0.05212,    Adjusted R-squared:  0.04733
F-statistic: 10.89 on 1 and 198 DF,  p-value: 0.001148
```

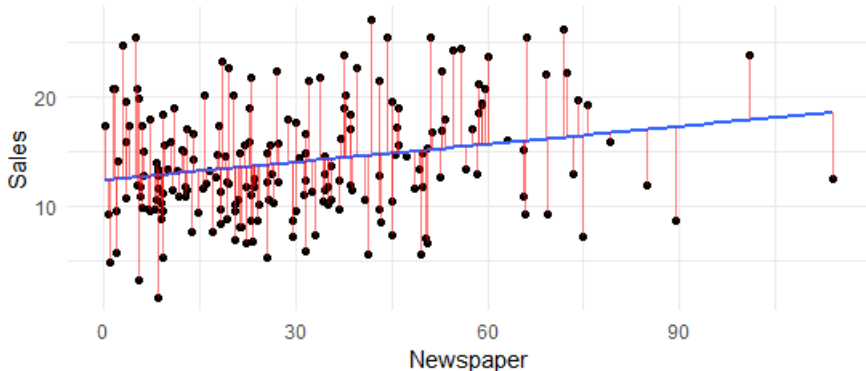

R^2 – EXAMPLE: REGRESSION LINE

```
1 ggplot(adv, aes(Newspaper, Sales)) +  
2   geom_point() +  
3   geom_smooth(method = "lm", se = F) +  
4   theme_minimal()
```



R^2 – EXAMPLE: REGRESSION LINE (ERRORS VISIBLE)

```
1 fit <- lm(Sales ~ Newspaper, data = adv)
2 ggplot(adv, aes(Newspaper, Sales)) +
3   geom_point() +
4   geom_smooth(method = 'lm', se = F) +
5   geom_segment(aes(xend=Newspaper, yend=fit$fitted), color='red', alpha
6     =0.5) +
7   theme_minimal()
```

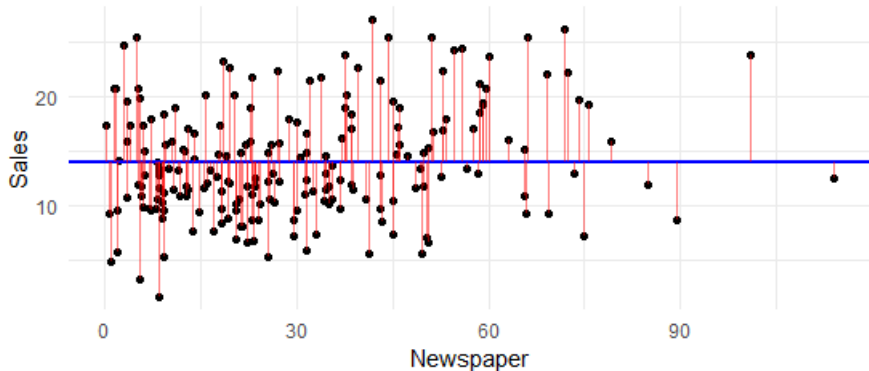


In the plot above, the blue line has the property that the sum of the squared red lines is minimal. Let's compare the sum of squared errors for this optimal fit to the fit with the line having the mean value.

```
1 summary(fit)$df[2] # (n - p - 1)
2 summary(fit)$sigma * summary(fit)$df[2] # 1008.311112571
3 m <- mean(adv$Sales) # 14.0225
4 sum((adv$Sales - m)^2) # 5417.149
```

R^2 – EXAMPLE: SUBOPTIMAL FIT

```
1 ggplot(adv, aes(Newspaper, Sales)) +  
2   theme_minimal() +  
3   geom_point() +  
4   geom_abline(slope=0, intercept=m, color='blue', size=1) +  
5   geom_segment(aes(xend=Newspaper, yend=m), color='red', alpha=0.5)
```



- σ^2 is the error variance, the variability in Y which is not explained by $X\beta$.
- $X\hat{\beta}$ explains most variability in Y .
- A measure for goodness of the fit with the linear model is the *coefficient of determination* R^2 which is defined by

$$R^2 := \frac{\text{TSS} - \text{RSS}}{\text{TSS}} = 1 - \frac{\text{RSS}}{\text{TSS}}. \quad (13)$$

- TSS is the *total sum of squares* which measures the total variance before the regression (see previous slide), defined by

$$\text{TSS} := \sum_{i=1}^n (Y_i - \bar{Y})^2. \quad (14)$$

- RSS is the *residual sum of squares* which measures the variability after performing the regression, defined by

$$\text{RSS} := \|\mathbf{y} - \mathbf{X}\hat{\beta}\|^2. \quad (15)$$

- By definition, R^2 is the proportion of the total variability minus the variability after the regression, in relation to the total variability.
- R^2 has values between 0 and 1.
- A value close to 1 indicates that a large proportion of the variability in the response has been explained by the regression.
- A value close to 0 indicates that the regression did not explain much of the variability in the response. Maybe, because the linear model is wrong, or the inherent error variance σ^2 is high, or both.
- In our example, we have

$$R^2 \approx 1 - \frac{1008.311}{5417.149} \approx 0.814,$$

meaning that approximately 81.4% of the variability have been explained by the regression.

- By including more (not perfectly collinear) predictors into the model will always increase explained variation.
- The *adjusted R^2* , denoted as \bar{R}^2 , measures as R^2 also how much variability have been explained by the regression, but also takes into accounts the number of predictors. It is defined as

$$\bar{R}^2 = 1 - (1 - R^2) \frac{n - 1}{n - p - 1} = 1 - \frac{\text{RSS}/(n - p - 1)}{\text{TSS}/(n - 1)}. \quad (16)$$

- \bar{R}^2 is smaller than R^2 if $\text{RSS} \neq 0$ and $p > 0$.
- As smaller p make inference easier, one should choose p such that \bar{R}^2 is the largest.

Residual plots show the fitted values \hat{y}_i against the observed values y_i , or the predictor values x_i against the residuals $e_i := y_i - \hat{y}_i$.

■ Residual plots are mainly useful for two things:

1. To validate/reject the suggested model.
2. To extract further information about the data.

■ Residual plots can have the following properties, among others:

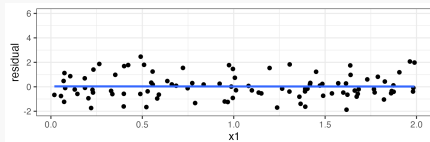
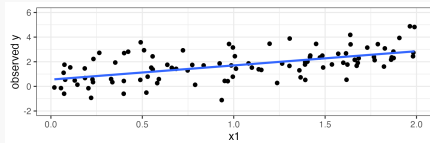
1. The values in the residual plot are scattered around zero without a visible trend \Rightarrow model assumption is reasonable.
2. The values in the residual plot exhibit a visible trend/pattern \Rightarrow model assumption is NOT reasonable.
3. The scatter plot or residual plot exhibits unusual values being far away from most of data \Rightarrow Outliers!
4. The magnitudes of the measurement errors are not roughly constant across observations \Rightarrow Heteroskedasticity (variance heterogeneity).

We create the following data consisting of 100 rows.

```
1 # Create data
2 df1 <- data.frame(x1=runif(n=100, min=0, max=2))
3 df1$yobs <- 0.1 + df1$x1 * 1.5 + rnorm(n=100, sd=1)
4
5 # Fit linear model to data and compute fitted values
6 fit1 <- lm(yobs ~ x1, data=df1)
7 df1$ypred <- fit1$coefficients[1] + df1$x1 * fit1$coefficients[2]
8
9 # Subtract fitted values from observed values
10 df1$residual <- df1$yobs - df1$ypred
```

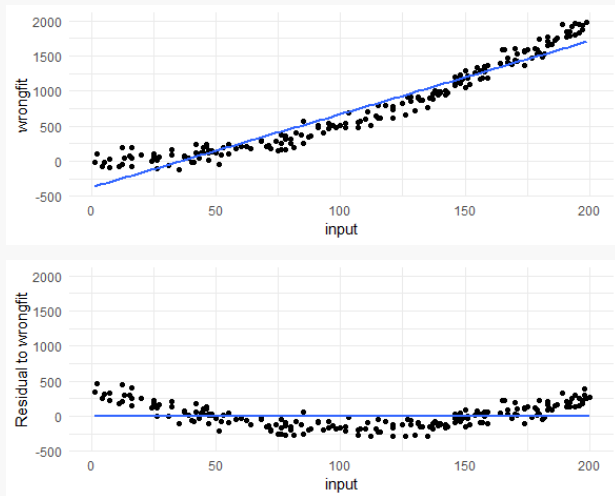
RESIDUAL PLOTS: DATA SET 1

```
1 library(ggplot2)
2 # Plot line of best fit
3 ggplot(df1, aes(x1, yobs)) +
4   geom_point() +
5   geom_smooth(method='lm', se=F) +
6   scale_y_continuous(limits=c(-2,6))+
7   labs(y='observed y') +
8   theme_minimal()
9
10 # Plot residuals
11 ggplot(df1, aes(x1, residual)) +
12   geom_point() +
13   geom_smooth(method='lm', se=F) +
14   scale_y_continuous(limits=c(-2,6))+
15   theme_minimal()
```



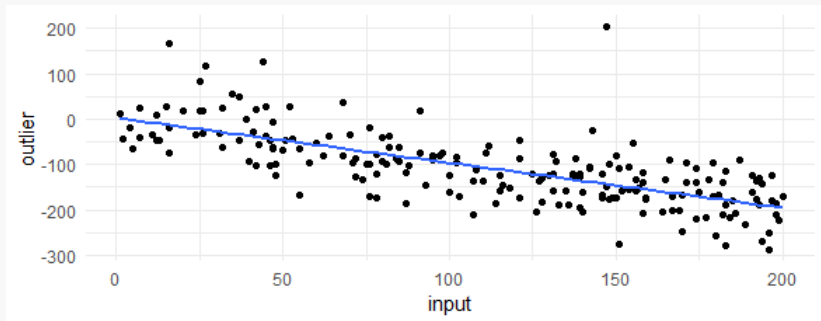
No visible pattern in residual plot \Rightarrow Proper fit of the data using a linear model

RESIDUAL PLOTS: DATA SET 2



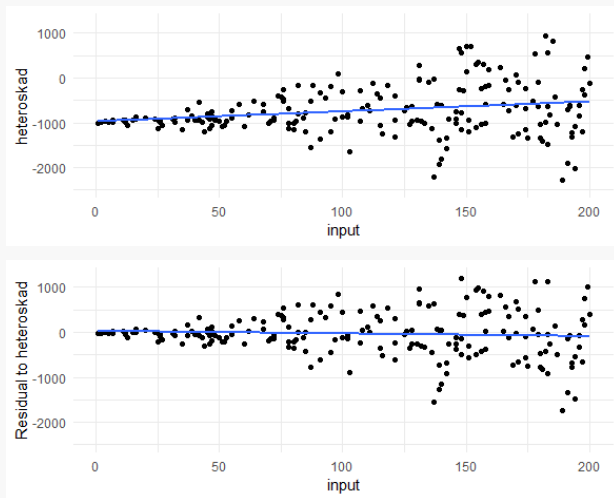
Visible pattern in the residual plot \Rightarrow The data are not properly fitted by the linear model. Maybe a quadratic relationship is reasonable.

RESIDUAL PLOTS: DATA SET 3



A very "unusual" value around $x = 150 \Rightarrow$ interpretable as an outlier.

RESIDUAL PLOTS: DATA SET 4



Error amplitudes increase as input increases \Rightarrow Signal seems to be well-modeled as a linear function, but errors are heteroskedastic.

4. PREDICTIONS

With coefficient estimates $\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_p$, it is straightforward to predict the response Y_{n+1} at a set of predictor values $x_{n+1,1}, \dots, x_{n+1,p}$. (How?)

Three types of uncertainty associated with this prediction:

1. Inaccuracy in the coefficient estimates $\hat{\beta}$ — quantify uncertainty using *confidence intervals*.
2. How well can the true model be captured by even the best linear model?
3. Inaccuracy in the prediction \hat{Y}_{n+1} — quantify uncertainty using *prediction intervals*.
 - ▶ Width of prediction interval incorporates both model uncertainty and observation variance.

SECTION 6: REGRESSION ANALYSIS WITH R

IDEA OF POLYNOMIAL REGRESSION

Polynomial regression extends the simple linear model by also allowing sums of predictors raised by powers, thus "polynomial".

- In polynomial regression, the response Y is modelled depending on the predictor X_1 with a polynomial function

$$y = \beta_0 + \beta_1 X_1 + \beta_2 X_1^2 + \beta_3 X_1^3 + \cdots + \beta_d X_1^d + \varepsilon, \quad (17)$$

where $d \in \mathbb{N}$ is the degree of the polynomial.

- The degree d describes the flexibility of the model.
- (What does a polynomial of order $d = 2$ look like? Order $d = 3$? $d = 4$?)

EXAMPLE 1: A NON-LINEAR FUNCTION

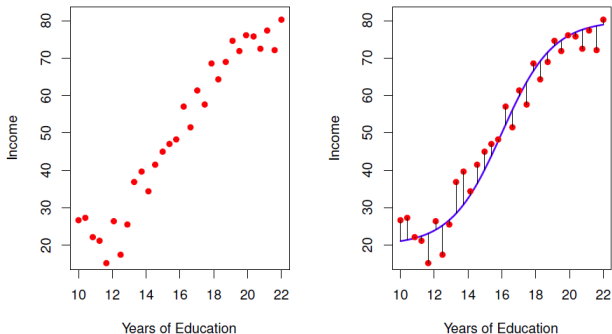


Figure 3: Image by James et al. (2021), based on the Income data set in R. The red dots are the observed values of income in tens of thousand dollars and years of education for 30 individuals.

EXAMPLE 2: DEGREE-4 POLYNOMIAL

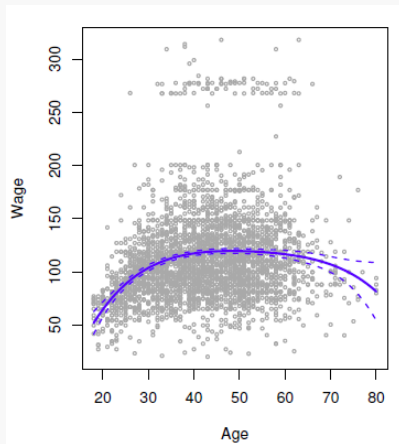


Figure 4: Image by James et al. (2021). The solid blue curve is a degree-4 polynomial of wage (in thousands of dollars) as a function of age, fit by least squares.

EXAMPLE 3: POLYNOMIAL REGRESSION WITH TWO PREDICTORS

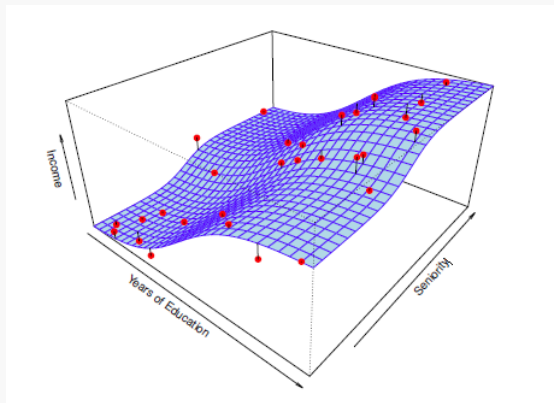


Figure 5: Image by James et al. (2021), based on the Income data set in R. The income is displayed as a function of years of education and seniority, where linearity does not seem appropriate. It might be reasonable to do polynomial regression with two predictors.