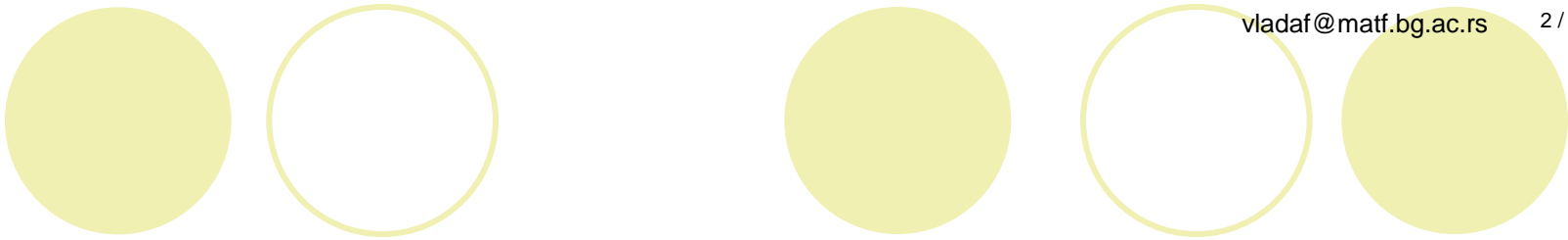


# Primjena računara u biologiji



**Vladimir Filipović**

[vladaf@matf.bg.ac.rs](mailto:vladaf@matf.bg.ac.rs)



# Simple Linear Regression, Multiple Linear Regression and Logistic Regression

# Elementary Statistics with R

- ▶ Qualitative Data

- ▶ Quantitative Data

---

- ▶ Numerical Measures

- ▶ Probability Distributions

- ▶ Interval Estimation

- ▶ Hypothesis Testing

- ▶ Type II Error

---

- ▶ Inference About Two Populations

- ▶ Goodness of Fit

- ▶ Analysis of Variance

- ▶ Non-parametric Methods

---

- ▶ Simple Linear Regression

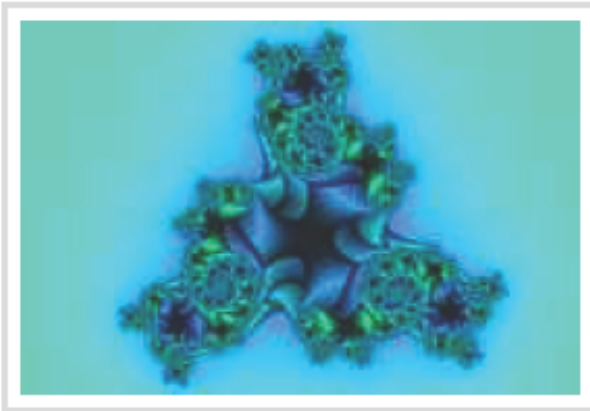
- ▶ Multiple Linear Regression

- ▶ Logistic Regression



# Simple Linear Regression

# Simple Linear Regression



A **simple linear regression model** that describes the relationship between two variables  $x$  and  $y$  can be expressed by the following equation. The numbers  $\alpha$  and  $\beta$  are called **parameters**, and  $\epsilon$  is the **error term**.

$$y = \alpha + \beta x + \epsilon$$

For example, in the data set **faithful**, it contains sample data of two random variables named **waiting** and **eruptions**. The **waiting** variable denotes the waiting time until the next eruptions, and **eruptions** denotes the duration. Its linear regression model can be expressed as:

$$Eruptions = \alpha + \beta * Waiting + \epsilon$$

# Simple Linear Regression

- Estimated Simple Regression Equation
- Coefficient of Determination
- Significance Test for Linear Regression
- Confidence Interval for Linear Regression
- Prediction Interval for Linear Regression
- Residual Plot
- Standardized Residual
- Normal Probability Plot of Residuals

# Estimated Simple Regression Equation

If we choose the parameters  $a$  and  $\beta$  in the **simple linear regression model** so as to minimize the sum of squares of the error term  $\epsilon$ , we will have the so called **estimated simple regression equation**. It allows us to compute **fitted values** of  $y$  based on values of  $x$ .

$$\hat{y} = a + bx$$

## Problem

Apply the simple linear regression model for the data set **faithful**, and estimate the next eruption duration if the waiting time since the last eruption has been 80 minutes.

# Estimated Simple Regression Equation (2)

## Solution

We apply the `lm` function to a formula that describes the variable `eruptions` by the variable `waiting`, and save the linear regression model in a new variable `eruption.lm`.

```
> eruption.lm = lm(eruptions ~ waiting, data=faithful)
```

Then we extract the parameters of the estimated regression equation with the `coefficients` function.

```
> coeffs = coefficients(eruption.lm); coeffs  
(Intercept)      waiting  
  -1.874016      0.075628
```



# Estimated Simple Regression Equation (3)

We now fit the eruption duration using the estimated regression equation.

```
> waiting = 80          # the waiting time
> duration = coeffs[1] + coeffs[2]*waiting
> duration
(Intercept)
  4.1762
```

## Answer

Based on the simple linear regression model, if the waiting time since the last eruption has been 80 minutes, we expect the next one to last 4.1762 minutes.

# Estimated Simple Regression Equation (4)

## Alternative Solution

We wrap the waiting parameter value inside a new **data frame** named newdata.

```
> newdata = data.frame(waiting=80) # wrap the parameter
```

Then we apply the predict function to eruption.lm along with newdata.

```
> predict(eruption.lm, newdata)    # apply predict
      1
4.1762
```

# Coefficient of Determination

The **coefficient of determination** of a **linear regression model** is the quotient of the **variances** of the **fitted values** and observed values of the dependent variable. If we denote  $y_i$  as the observed values of the dependent variable,  $\bar{y}$  as its **mean**, and  $\hat{y}_i$  as the fitted value, then the coefficient of determination is:

$$r^2 = \frac{\sum(\hat{y}_i - \bar{y})^2}{\sum(y_i - \bar{y})^2}$$

## Problem

Find the coefficient of determination for the simple linear regression model of the data set **faithful**.

# Coefficient of Determination

The **coefficient of determination** of a **linear regression model** is the quotient of the **variances** of the **fitted values** and observed values of the dependent variable. If we denote  $y_i$  as the observed values of the dependent variable,  $\bar{y}$  as its **mean**, and  $\hat{y}_i$  as the fitted value, then the coefficient of determination is:

$$r^2 = \frac{\sum(\hat{y}_i - \bar{y})^2}{\sum(y_i - \bar{y})^2}$$

## Problem

Find the coefficient of determination for the simple linear regression model of the data set **faithful**.

# Coefficient of Determination (2)

## Solution

We apply the `lm` function to a formula that describes the variable eruptions by the variable waiting, and save the linear regression model in a new variable `eruption.lm`.

```
> eruption.lm = lm(eruptions ~ waiting, data=faithful)
```

Then we extract the coefficient of determination from the `r.squared` attribute of its summary.

```
> summary(eruption.lm)$r.squared  
[1] 0.81146
```

## Answer

The coefficient of determination of the simple linear regression model for the data set `faithful` is 0.81146.

# Coefficient of Determination (3)

## Note

Further detail of the `r.squared` attribute can be found in the R documentation.

```
> help(summary.lm)
```

# Significance Test for Linear Regression

Assume that the error term  $\epsilon$  in the **linear regression model** is independent of  $x$ , and is **normally distributed**, with zero **mean** and constant **variance**. We can decide whether there is any **significant relationship** between  $x$  and  $y$  by testing the null hypothesis that  $\beta = 0$ .

## Problem

Decide whether there is a significant relationship between the variables in the linear regression model of the data set **faithful** at .05 significance level.

## Solution

We apply the `lm` function to a formula that describes the variable eruptions by the variable waiting, and save the linear regression model in a new variable `eruption.lm`.

```
> eruption.lm = lm(eruptions ~ waiting, data=faithful)
```

# Significance Test for Linear Regression (2)

Then we print out the F-statistics of the significance test with the summary function.

```
> summary(eruption.lm)

Call:
lm(formula = eruptions ~ waiting, data = faithful)

Residuals:
    Min       1Q   Median       3Q      Max
-1.2992 -0.3769  0.0351  0.3491  1.1933

Coefficients:
              Estimate Std. Error t value Pr(>|t|)
(Intercept) -1.87402     0.16014   -11.7   <2e-16 ***
waiting      0.07563     0.00222    34.1   <2e-16 ***
---

```



# Significance Test for Linear Regression (3)

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

```
Residual standard error: 0.497 on 270 degrees of freedom
```

```
Multiple R-squared: 0.811,      Adjusted R-squared: 0.811
```

```
F-statistic: 1.16e+03 on 1 and 270 DF,  p-value: <2e-16
```

## Answer

As the p-value is much less than 0.05, we reject the null hypothesis that  $\beta = 0$ . Hence there is a significant relationship between the variables in the linear regression model of the data set faithful.

## Note

Further detail of the summary function for linear regression model can be found in the R documentation.

```
> help(summary.lm)
```

# Confidence Interval for Linear Regression

Assume that the error term  $\epsilon$  in the **linear regression model** is independent of  $x$ , and is **normally distributed**, with zero **mean** and constant **variance**. For a given value of  $x$ , the interval estimate for the mean of the dependent variable,  $\bar{y}$ , is called the **confidence interval**.

## Problem

In the data set **faithful**, develop a 95% confidence interval of the mean eruption duration for the waiting time of 80 minutes.

## Solution

We apply the `lm` function to a formula that describes the variable eruptions by the variable waiting, and save the linear regression model in a new variable `eruption.lm`.

```
> attach(faithful)      # attach the data frame  
> eruption.lm = lm(eruptions ~ waiting)
```

Then we create a new **data frame** that set the waiting time value.

```
> newdata = data.frame(waiting=80)
```

# Confidence Interval for Linear Regression (2)

We now apply the `predict` function and set the predictor variable in the `newdata` argument. We also set the interval type as "confidence", and use the default 0.95 confidence level.

```
> predict(eruption.lm, newdata, interval="confidence")
      fit      lwr      upr
1 4.1762 4.1048 4.2476
> detach(faithful)      # clean up
```

## Answer

The 95% confidence interval of the mean eruption duration for the waiting time of 80 minutes is between 4.1048 and 4.2476 minutes.

## Note

Further detail of the `predict` function for linear regression model can be found in the R documentation.

```
> help(predict.lm)
```

# Prediction Interval for Linear Regression

Assume that the error term  $\epsilon$  in the **simple linear regression model** is independent of  $x$ , and is **normally distributed**, with zero **mean** and constant **variance**. For a given value of  $x$ , the interval estimate of the dependent variable  $y$  is called the **prediction interval**.

## Problem

In the data set **faithful**, develop a 95% prediction interval of the eruption duration for the waiting time of 80 minutes.

## Solution

We apply the `lm` function to a formula that describes the variable eruptions by the variable waiting, and save the linear regression model in a new variable `eruption.lm`.

```
> attach(faithful)      # attach the data frame  
> eruption.lm = lm(eruptions ~ waiting)
```

# Prediction Interval for Linear Regression (2)

Then we create a new **data frame** that set the waiting time value.

```
> newdata = data.frame(waiting=80)
```

We now apply the predict function and set the predictor variable in the newdata argument. We also set the interval type as "predict", and use the default 0.95 confidence level.

```
> predict(eruption.lm, newdata, interval="predict")
      fit      lwr      upr
1 4.1762 3.1961 5.1564
> detach(faithful)      # clean up
```

## Answer

The 95% prediction interval of the eruption duration for the waiting time of 80 minutes is between 3.1961 and 5.1564 minutes.

# Prediction Interval for Linear Regression (3)

## Note

Further detail of the predict function for linear regression model can be found in the R documentation.

```
> help(predict.lm)
```

# Residual Plot

The **residual** data of the **simple linear regression model** is the difference between the observed data of the dependent variable  $y$  and the **fitted values**  $\hat{y}$ .

$$\text{Residual} = y - \hat{y}$$

## Problem

Plot the residual of the simple linear regression model of the data set **faithful** against the independent variable waiting.

## Solution

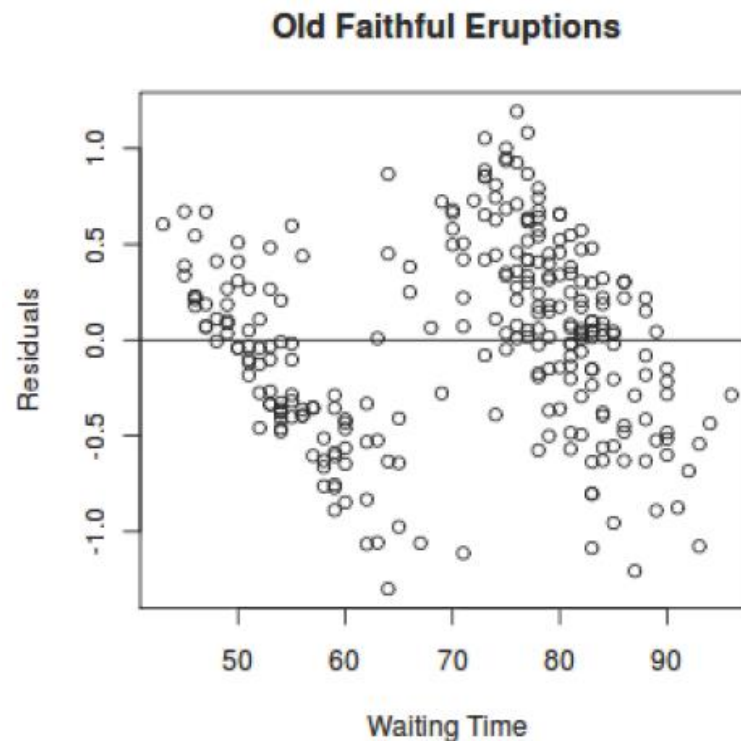
We apply the `lm` function to a formula that describes the variable eruptions by the variable waiting, and save the linear regression model in a new variable `eruption.lm`. Then we compute the residual with the `resid` function.

```
> eruption.lm = lm(eruptions ~ waiting, data=faithful)
> eruption.res = resid(eruption.lm)
```

# Residual Plot (2)

We now plot the residual against the observed values of the variable waiting.

```
> plot(faithful$waiting, eruption.res,  
+      ylab="Residuals", xlab="Waiting Time",  
+      main="Old Faithful Eruptions")  
> abline(0, 0)           # the horizon
```





# Residual Plot (3)

## Note

Further detail of the resid function can be found in the R documentation.

```
> help(resid)
```

# Standardized Residual

The **standardized residual** is the **residual** divided by its **standard deviation**.

$$\text{Standardized Residual } i = \frac{\text{Residual } i}{\text{Standard Deviation of Residual } i}$$

## Problem

Plot the standardized residual of the simple linear regression model of the data set **faithful** against the independent variable waiting.

## Solution

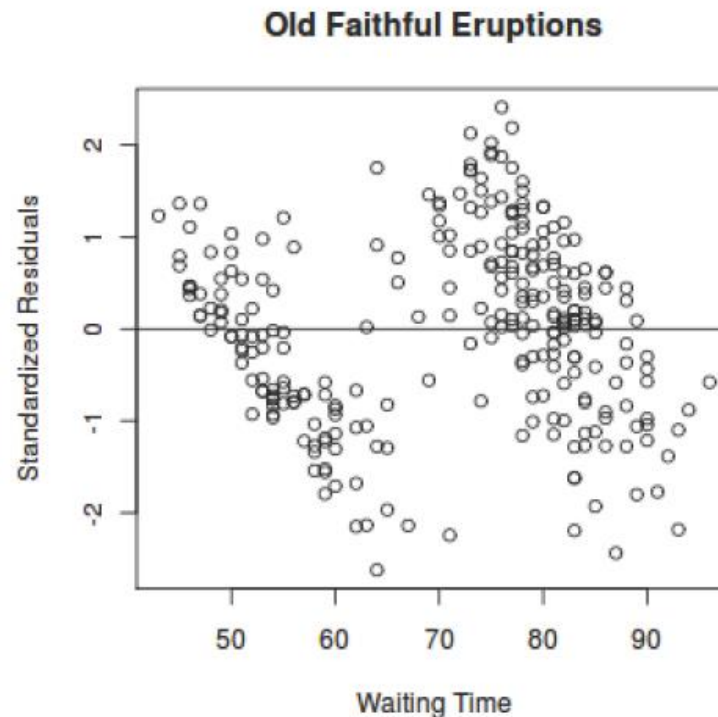
We apply the `lm` function to a formula that describes the variable eruptions by the variable waiting, and save the linear regression model in a new variable `eruption.lm`. Then we compute the standardized residual with the `rstandard` function.

```
> eruption.lm = lm(eruptions ~ waiting, data=faithful)
> eruption.stdres = rstandard(eruption.lm)
```

# Standardized Residual (2)

We now plot the standardized residual against the observed values of the variable waiting.

```
> plot(faithful$waiting, eruption.stdres,  
+      ylab="Standardized Residuals",  
+      xlab="Waiting Time",  
+      main="Old Faithful Eruptions")  
> abline(0, 0) # the horizon
```



# Standardized Residual (3)

## Note

Further detail of the `rstandard` function can be found in the R documentation.

```
> help(rstandard)
```

# Normal Probability Plot of Residuals

The **normal probability plot** is a graphical tool for comparing a data set with the **normal distribution**. We can use it with the **standardized residual** of the **linear regression model** and see if the error term  $\epsilon$  is actually normally distributed.

## Problem

Create the normal probability plot for the standardized residual of the data set **faithful**.

## Solution

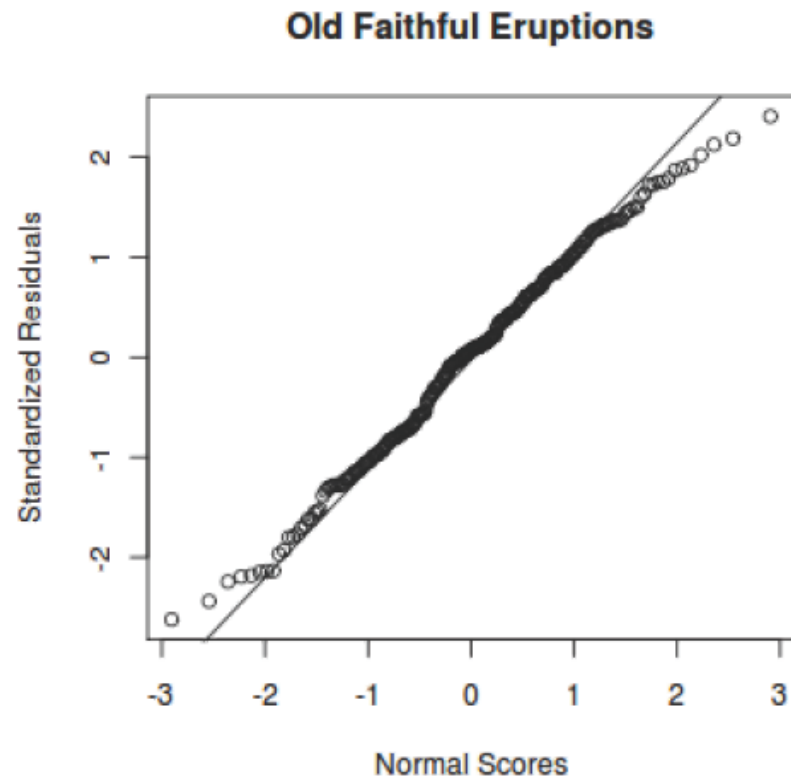
We apply the `lm` function to a formula that describes the variable eruptions by the variable waiting, and save the linear regression model in a new variable `eruption.lm`. Then we compute the standardized residual with the `rstandard` function.

```
> eruption.lm = lm(eruptions ~ waiting, data=faithful)
> eruption.stdres = rstandard(eruption.lm)
```

# Normal Probability Plot of Residuals (2)

We now create the normal probability plot with the `qqnorm` function, and add the `qqline` for further comparison.

```
> qqnorm(eruption.stdres,  
+        ylab="Standardized Residuals",  
+        xlab="Normal Scores",  
+        main="Old Faithful Eruptions")  
> qqline(eruption.stdres)
```



# Normal Probability Plot of Residuals (3)

## Note

Further detail of the `qqnorm` and `qqline` functions can be found in the R documentation.

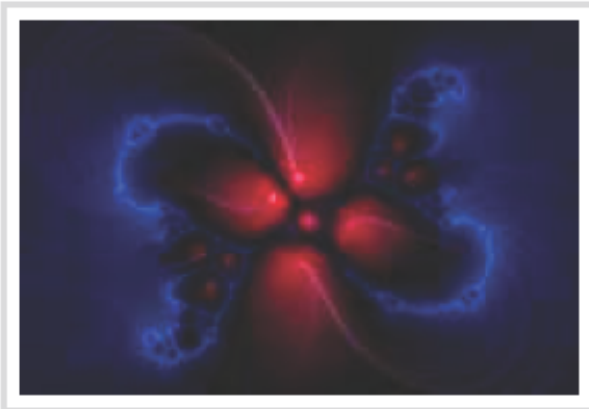
```
> help(qqnorm)
```



# Multiple Linear Regression



# Multiple Linear Regression



A **multiple linear regression** (MLR) model that describes a dependent variable  $y$  by independent variables  $x_1, x_2, \dots, x_p$  ( $p > 1$ ) is expressed by the equation as follows, where the numbers  $\alpha$  and  $\beta_k$  ( $k = 1, 2, \dots, p$ ) are the **parameters**, and  $\epsilon$  is the **error term**.

$$y = \alpha + \sum_k \beta_k x_k + \epsilon$$

For example, in the built-in data set `stackloss` from observations of a chemical plant operation, if we assign `stackloss` as the dependent variable, and assign `Air.Flow` (cooling air flow), `Water.Temp` (inlet water temperature) and `Acid.Conc.` (acid concentration) as independent variables, the multiple linear regression model is:

$$\text{Stack.Loss} = \alpha + \beta_1 * \text{Air.Flow} + \beta_2 * \text{Water.Temp} + \beta_3 * \text{Acid.Conc.} + \epsilon$$

Further detail of the `stackloss` data set can be found in the R documentation.

```
> help(stackloss)
```

# Multiple Linear Regression (2)

- Estimated Multiple Regression Equation
- Multiple Coefficient of Determination
- Adjusted Coefficient of Determination
- Significance Test for MLR
- Confidence Interval for MLR
- Prediction Interval for MLR

# Estimated Multiple Regression Equation

If we choose the parameters  $a$  and  $\beta_k$  ( $k = 1, 2, \dots, p$ ) in the **multiple linear regression model** so as to minimize the sum of squares of the error term  $\epsilon$ , we will have the so called **estimated multiple regression equation**. It allows us to compute **fitted values** of  $y$  based on a set of values of  $x_k$  ( $k = 1, 2, \dots, p$ ).

$$\hat{y} = a + \sum_k b_k x_k$$

## Problem

Apply the multiple linear regression model for the data set **stackloss**, and predict the stack loss if the air flow is 72, water temperature is 20 and acid concentration is 85.

# Estimated Multiple Regression Equation (2)

## Solution

We apply the `lm` function to a formula that describes the variable `stack.loss` by the variables `Air.Flow`, `Water.Temp` and `Acid.Conc.`. And we save the linear regression model in a new variable `stackloss.lm`.

```
> stackloss.lm = lm(stack.loss ~  
+   Air.Flow + Water.Temp + Acid.Conc.,  
+   data=stackloss)
```

We also wrap the parameters inside a new **data frame** named `newdata`.

```
> newdata = data.frame(Air.Flow=72, # wrap the parameters  
+   Water.Temp=20,  
+   Acid.Conc.=85)
```

Lastly, we apply the `predict` function to `stackloss.lm` and `newdata`.

```
> predict(stackloss.lm, newdata)  
1  
24.582
```

# Estimated Multiple Regression Equation (3)

## Answer

Based on the multiple linear regression model and the given parameters, the predicted stack loss is 24.582.

# Multiple Coefficient of Determination

The **coefficient of determination** of a **multiple linear regression model** is the quotient of the **variances** of the **fitted values** and observed values of the dependent variable. If we denote  $y_i$  as the observed values of the dependent variable,  $\bar{y}$  as its **mean**, and  $\hat{y}_i$  as the fitted value, then the coefficient of determination is:

$$R^2 = \frac{\sum(\hat{y}_i - \bar{y})^2}{\sum(y_i - \bar{y})^2}$$

## Problem

Find the coefficient of determination for the multiple linear regression model of the data set **stackloss**.

# Multiple Coefficient of Determination (2)

## Solution

We apply the `lm` function to a formula that describes the variable `stack.loss` by the variables `Air.Flow`, `Water.Temp` and `Acid.Conc.`. And we save the linear regression model in a new variable `stackloss.lm`.

```
> stackloss.lm = lm(stack.loss ~  
+   Air.Flow + Water.Temp + Acid.Conc.,  
+   data=stackloss)
```

Then we extract the coefficient of determination from the `r.squared` attribute of its summary.

```
> summary(stackloss.lm)$r.squared  
[1] 0.91358
```

# Multiple Coefficient of Determination (3)

## Answer

The coefficient of determination of the multiple linear regression model for the data set `stackloss` is 0.91358.

## Note

Further detail of the `r.squared` attribute can be found in the R documentation.

```
> help(summary.lm)
```



# Adjusted Coefficient of Determination

The **adjusted coefficient of determination** of a **multiple linear regression model** is defined in terms of the **coefficient of determination** as follows, where  $n$  is the number of observations in the data set, and  $p$  is the number of independent variables.

$$R_{adj}^2 = 1 - (1 - R^2) \frac{n - 1}{n - p - 1}$$

## Problem

Find the adjusted coefficient of determination for the multiple linear regression model of the data set **stackloss**.

# Adjusted Coefficient of Determination (2)

## Solution

We apply the `lm` function to a formula that describes the variable `stack.loss` by the variables `Air.Flow`, `Water.Temp` and `Acid.Conc.`. And we save the linear regression model in a new variable `stackloss.lm`.

```
> stackloss.lm = lm(stack.loss ~  
+   Air.Flow + Water.Temp + Acid.Conc.,  
+   data=stackloss)
```

Then we extract the coefficient of determination from the `adj.r.squared` attribute of its summary.

```
> summary(stackloss.lm)$adj.r.squared  
[1] 0.89833
```

# Adjusted Coefficient of Determination (3)

## Answer

The adjusted coefficient of determination of the multiple linear regression model for the data set `stackloss` is 0.89833.

## Note

Further detail of the `adj.r.squared` attribute can be found in the R documentation.

```
> help(summary.lm)
```

# Significance Test for MLR

Assume that the error term  $\epsilon$  in the **multiple linear regression (MLR) model** is independent of  $x_k$  ( $k = 1, 2, \dots, p$ ), and is **normally distributed**, with zero **mean** and constant **variance**.

We can decide whether there is any **significant relationship** between the dependent variable  $y$  and any of the independent variables  $x_k$  ( $k = 1, 2, \dots, p$ ).

## Problem

Decide which of the independent variables in the multiple linear regression model of the data set **stackloss** are statistically significant at .05 significance level.

## Solution

We apply the `lm` function to a formula that describes the variable `stack.loss` by the variables `Air.Flow`, `Water.Temp` and `Acid.Conc.`. And we save the linear regression model in a new variable `stackloss.lm`.

```
> stackloss.lm = lm(stack.loss ~  
+   Air.Flow + Water.Temp + Acid.Conc.,  
+   data=stackloss)
```

# Significance Test for MLR (2)

The t values of the independent variables can be found with the summary function.

```
> summary(stackloss.lm)
```

Call:

```
lm(formula = stack.loss ~ Air.Flow + Water.Temp + Acid.Conc.,  
    data = stackloss)
```

Residuals:

Min	1Q	Median	3Q	Max
-7.238	-1.712	-0.455	2.361	5.698

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )	
(Intercept)	-39.920	11.896	-3.36	0.0038	**
Air.Flow	0.716	0.135	5.31	5.8e-05	***
Water.Temp	1.295	0.368	3.52	0.0026	**
Acid.Conc.	-0.152	0.156	-0.97	0.3440	

# Significance Test for MLR (3)

```
---
```

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

```
Residual standard error: 3.24 on 17 degrees of freedom
```

```
Multiple R-squared: 0.914,    Adjusted R-squared: 0.898
```

```
F-statistic: 59.9 on 3 and 17 DF,  p-value: 3.02e-09
```

## Answer

As the p-values of Air.Flow and Water.Temp are less than 0.05, they are both statistically significant in the multiple linear regression model of stackloss.

## Note

Further detail of the summary function for linear regression model can be found in the R documentation.

```
> help(summary.lm)
```

# Confidence Interval for MLR

Assume that the error term  $\epsilon$  in the **multiple linear regression (MLR) model** is independent of  $x_k$  ( $k = 1, 2, \dots, p$ ), and is **normally distributed**, with zero **mean** and constant **variance**.

For a given set of values of  $x_k$  ( $k = 1, 2, \dots, p$ ), the interval estimate for the mean of the dependent variable,  $\bar{y}$ , is called the **confidence interval**.

## Problem

In data set stackloss, develop a 95% confidence interval of the stack loss if the air flow is 72, water temperature is 20 and acid concentration is 85.

## Solution

We apply the `lm` function to a formula that describes the variable `stack.loss` by the variables `Air.Flow`, `Water.Temp` and `Acid.Conc`. And we save the linear regression model in a new variable `stackloss.lm`.

```
> attach(stackloss)    # attach the data frame
> stackloss.lm = lm(stack.loss ~
+   Air.Flow + Water.Temp + Acid.Conc.)
```

# Confidence Interval for MLR (2)

Then we wrap the parameters inside a new **data frame** variable newdata.

```
> newdata = data.frame(Air.Flow=72,  
+   Water.Temp=20,  
+   Acid.Conc.=85)
```

We now apply the predict function and set the predictor variable in the newdata argument. We also set the interval type as "confidence", and use the default 0.95 confidence level.

```
> predict(stackloss.lm, newdata, interval="confidence")  
      fit      lwr      upr  
1 24.582 20.218 28.945  
> detach(stackloss)    # clean up
```



# Confidence Interval for MLR (3)

## Answer

The 95% confidence interval of the stack loss with the given parameters is between 20.218 and 28.945.

## Note

Further detail of the `predict` function for linear regression model can be found in the R documentation.

```
> help(predict.lm)
```

# Prediction Interval for MLR

Assume that the error term  $\epsilon$  in the **multiple linear regression (MLR) model** is independent of  $x_k$  ( $k = 1, 2, \dots, p$ ), and is **normally distributed**, with zero **mean** and constant **variance**.

For a given set of values of  $x_k$  ( $k = 1, 2, \dots, p$ ), the interval estimate of the dependent variable  $y$  is called the **prediction interval**.

## Problem

In data set **stackloss**, develop a 95% prediction interval of the stack loss if the air flow is 72, water temperature is 20 and acid concentration is 85.

## Solution

We apply the `lm` function to a formula that describes the variable `stack.loss` by the variables `Air.Flow`, `Water.Temp` and `Acid.Conc`. And we save the linear regression model in a new variable `stackloss.lm`.

```
> attach(stackloss)    # attach the data frame
> stackloss.lm = lm(stack.loss ~
+   Air.Flow + Water.Temp + Acid.Conc.)
```

# Prediction Interval for MLR (2)

Then we wrap the parameters inside a new **data frame** variable newdata.

```
> newdata = data.frame(Air.Flow=72,  
+   Water.Temp=20,  
+   Acid.Conc.=85)
```

We now apply the predict function and set the predictor variable in the newdata argument. We also set the interval type as "predict", and use the default 0.95 confidence level.

```
> predict(stackloss.lm, newdata, interval="predict")  
      fit      lwr      upr  
1 24.582 16.466 32.697  
> detach(stackloss)    # clean up
```

# Prediction Interval for MLR (3)

## Answer

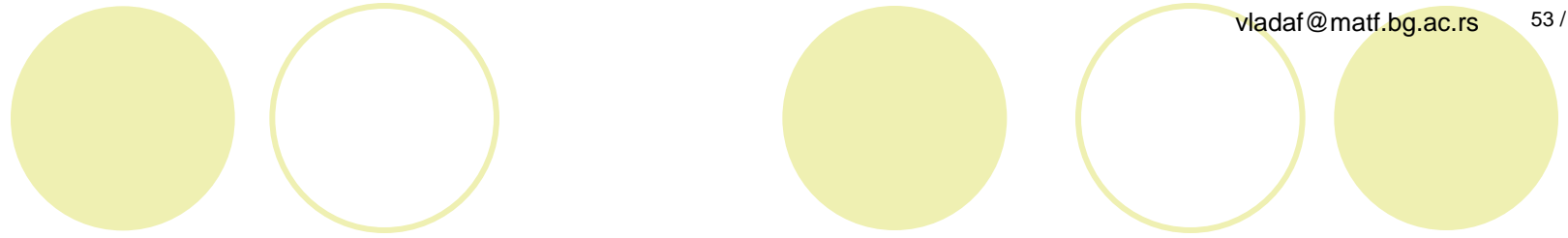
The 95% confidence interval of the stack loss with the given parameters is between 16.466 and 32.697.

## Note

Further detail of the predict function for linear regression model can be found in the R documentation.

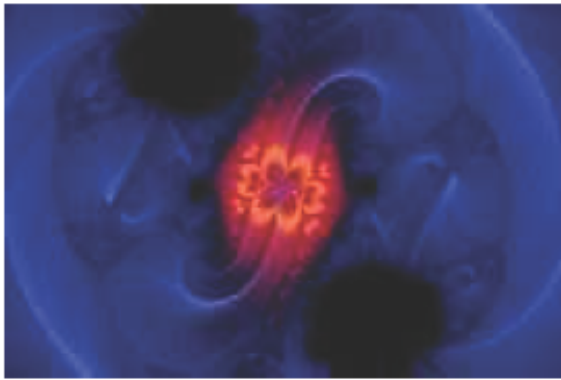
```
> help(predict.lm)
```

---



# Logistic Regression

# Logistic Regression



We use the **logistic regression equation** to predict the probability of a dependent variable taking the dichotomy values 0 or 1. Suppose  $x_1, x_2, \dots, x_p$  are the independent variables,  $\alpha$  and  $\beta_k$  ( $k = 1, 2, \dots, p$ ) are the parameters, and  $E(y)$  is the expected value of the dependent variable  $y$ , then the logistic regression equation is:

$$E(y) = 1 / (1 + e^{-(\alpha + \sum_k \beta_k x_k)})$$

For example, in the built-in data set **mtcars**, the data column **am** represents the transmission type of the automobile model (0 = automatic, 1 = manual). With the logistic regression equation, we can model the probability of a manual transmission in a vehicle based on its engine horsepower and weight data.

$$P(\text{Manual Transmission}) = 1 / (1 + e^{-(\alpha + \beta_1 * \text{Horsepower} + \beta_2 * \text{Weight})})$$

# Logistic Regression (2)

- Estimated Logistic Regression Equation
- Significance Test for Logistic Regression

# Estimated Logistic Regression Equation

Using the generalized linear model, an **estimated logistic regression equation** can be formulated as below. The coefficients  $a$  and  $b_k$  ( $k = 1, 2, \dots, p$ ) are determined according to a maximum likelihood approach, and it allows us to estimate the probability of the dependent variable  $y$  taking on the value 1 for given values of  $x_k$  ( $k = 1, 2, \dots, p$ ).

$$\text{Estimate of } P(y = 1 \mid x_1, \dots, x_p) = 1 / (1 + e^{-(a + \sum_k b_k x_k)})$$

## Problem

By use of the **logistic regression equation of vehicle transmission** in the data set **mtcars**, estimate the probability of a vehicle being fitted with a manual transmission if it has a 120hp engine and weights 2800 lbs.



# Estimated Logistic Regression Equation (2)

## Solution

We apply the function `glm` to a formula that describes the transmission type (`am`) by the horsepower (`hp`) and weight (`wt`). This creates a generalized linear model (GLM) in the binomial family.

```
> am.glm = glm(formula=am ~ hp + wt,  
+              data=mtcars,  
+              family=binomial)
```

We then wrap the test parameters inside a **data frame** `newdata`.

```
> newdata = data.frame(hp=120, wt=2.8)
```

Now we apply the function `predict` to the generalized linear model `am.glm` along with `newdata`. We will have to select *response* prediction type in order to obtain the predicted probability.

```
> predict(am.glm, newdata, type="response")  
1  
0.64181
```

# Estimated Logistic Regression Equation (3)

## Answer

For an automobile with 120hp engine and 2800 lbs weight, the probability of it being fitted with a manual transmission is about 64%.

## Note

Further detail of the function `predict` for generalized linear model can be found in the R documentation.

```
> help(predict.glm)
```

# Significance Test for Logistic Regression

We can decide whether there is any significant relationship between the dependent variable  $y$  and the independent variables  $x_k$  ( $k = 1, 2, \dots, p$ ) in the **logistic regression equation**. In particular, if any of the null hypothesis that  $\beta_k = 0$  ( $k = 1, 2, \dots, p$ ) is valid, then  $x_k$  is statistically insignificant in the logistic regression model.

## Problem

At .05 significance level, decide if any of the independent variables in the logistic regression model of vehicle transmission in data set **mtcars** is statistically insignificant.

## Solution

We apply the function `glm` to a formula that describes the transmission type (`am`) by the horsepower (`hp`) and weight (`wt`). This creates a generalized linear model (GLM) in the binomial family.

```
> am.glm = glm(formula=am ~ hp + wt,  
+               data=mtcars,  
+               family=binomial)
```

# Significance Test for Logistic Regression (2)

We then print out the summary of the generalized linear model and check for the p-values of the hp and wt variables.

```
> summary(am.glm)
```

Call:

```
glm(formula = am ~ hp + wt, family = binomial, data = mtcars)
```

Deviance Residuals:

Min	1Q	Median	3Q	Max
-2.2537	-0.1568	-0.0168	0.1543	1.3449

Coefficients:

	Estimate	Std. Error	z value	Pr(> z )	
(Intercept)	18.8663	7.4436	2.53	0.0113	*
hp	0.0363	0.0177	2.04	0.0409	*
wt	-8.0835	3.0687	-2.63	0.0084	**

# Significance Test for Logistic Regression (3)

```
---  
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1  
  
(Dispersion parameter for binomial family taken to be 1)  
  
Null deviance: 43.230  on 31  degrees of freedom  
Residual deviance: 10.059  on 29  degrees of freedom  
AIC: 16.06  
  
Number of Fisher Scoring iterations: 8
```

## Answer

As the p-values of the hp and wt variables are both less than 0.05, neither hp or wt is insignificant in the logistic regression model.

## Note

Further detail of the function summary for the generalized linear model can be found in the R documentation.

```
> help(summary.glm)
```

# Acknowledgments

Material in this presentation is taken from  
<http://www.r-tutor.com/>