

## Session 3: Implementing linear regression in R and interpreting the output


Slide 1:	Introduction .....	2
Slide 2:	Section 1: Case Study 1: Using Multiple Linear Regression and Model Building .....	2
Slide 3:	Case Study 1.: Student Grades in Statistics .....	3
Slide 4:	The Dataset .....	4
Slide 5:	Identifying Relationships Between Variables and Predictors .....	4
Slide 6:	Creating the Scatterplot .....	5
Slide 7:	Obtaining and Interpreting the Regression Equation .....	6
Slide 8:	Excluding Predictors .....	7
Slide 9:	Excluding a Second Predictor.....	8
Slide 10:	Model Building .....	8
Tab 1:	Adjusted $R^2$ .....	9
Tab 2:	Akaike's Information Criterion .....	10
Slide 11:	Residuals .....	13
Slide 12:	Predictions.....	13
Slide 13:	Section 2: Case Study 2: Using Polynomial Regression .....	14
Slide 14:	Case Study 2: Salary From Years of Experience .....	15
Slide 15:	Scatterplot.....	16
Slide 16:	Fit SLR Model .....	16
Slide 17:	Polynomial Regression Models .....	17
Slide 18:	Implementing a Quadratic Polynomial Regression Model in R.....	18
Slide 19:	Implementing a Cubic Polynomial Regression Model in R.....	19
Slide 20:	Conclusions Derived from the Quadratic Model .....	19
Slide 21:	Summary.....	20



Slide 1: **Introduction**


Trinity College Dublin  
Coláiste na Tríonóide, Baile Átha Cliath  
The University of Dublin

## Implementing Linear Regression in R and Interpreting the Output

 Presenter: Athanasios Georgiadis  
Duration: 22:00  
School: School of Computer Science and Statistics


Hello and welcome to this presentation on Implementing Linear Regression in R and Interpreting the Output. My name is Athanasios Georgiadis and I'm the instructor for this presentation where we will look at two case studies and implement linear regression in R. In the first case study, we will implement multiple linear regression and model building. In the second, we'll implement polynomial regression.

Slide 2: **Section 1: Case Study 1: Using Multiple Linear Regression and Model Building**

 2 of 21

## Using Multiple Linear Regression and Model Building

## Slide 3: Case Study 1.: Student Grades in Statistics



**Case Study 1: Student Grades in Statistics**


3 of 21

### Case Study

- To model student grades in Statistics, based on their performance in other modules
- Dataset:
  - Grades of  $n=10$  students for:
    - Statistics
    - Mathematics
    - Computer Science
    - Foreign Languages

### Challenges in Evaluating the Model

- Are all the modules used as predictors, affecting the grade in Statistics?
- How do we finalise the model to be used for predictions?



In our first case study, we would like to model Students' grades in Statistics, based on their performance in other modules. We have a dataset with the grades of  $n=10$  Students in Statistics and some other modules; namely Mathematics, Computer Science and Foreign Languages.

We will use R to fit a multiple linear regression model. To this end, we would like to evaluate the obtained model.

Are all the modules used as predictors effecting the grade in Statistics?

How do we finalise the model to be used for predictions?

All these are just a few challenges that a model developer or user has to face when evaluating or creating a model..

In this presentation, we will address these challenges using a dataset and R language.


## Slide 4: The Dataset

Grades for $n = 10$ students			
Statistics (S)	Mathematics (M)	Computer Science (C)	Foreign Languages (L)
70	60	81	70
72	61	82	60
74	65	80	50
75	65	85	80
77	66	85	90
80	68	85	50
82	70	86	60
85	70	87	90
88	75	87	70
90	75	90	80

Take time to view the information on this slide.

We collect the grades of  $n=10$  Students in the modules of Statistics (S), Mathematics (M), Computer Science (C) and Foreign Languages (L). The dataset is presented in this table.


## Slide 5: Identifying Relationships Between Variables and Predictors



### Identifying Relationships Between Variables and Predictors

5 of 21

- **Variable Y:** (Statistics grade) should be explained by the predictors:
  - $X_1$ : Mathematics grade
  - $X_2$ : Computer Science grade
  - $X_3$ : Foreign Languages grade
- Potential relationships between the response variable and the predictors may be instinctively identified.
  - These relationships must be justified statistically.
- The data analyst may not be able to instinctively identify potential relationships in cases where their knowledge is limited.



Let's collect some initial thoughts.

The variable  $Y$ := the grade in Statistics, should be hopefully explained by the predictors  $X_1$ := grade in Mathematics,  $X_2$ := grade in Computers and  $X_3$ := grade in Languages.


In the specific case study, we may have some experience or instinct about potential relationships between the response variable and the predictors. For example, we may expect the dependence of  $Y$  on  $X_1$ : the grade in Mathematics. On the other hand, the connection between  $Y$  and  $X_3$  (the grade in Foreign languages), may be poor or non-existent.

These relationships must be justified statistically.

There may also be occasions when a data analyst has to deal with cases where his or her understanding of the data is limited and therefore, they do not instinctively recognise potential relationships.

Over the following slides, we are going to create a scatterplot and determine multiple linear regressions for different groups of variables to help determine which is the best to use.

### Slide 6: Creating the Scatterplot



## Creating the Scatterplot

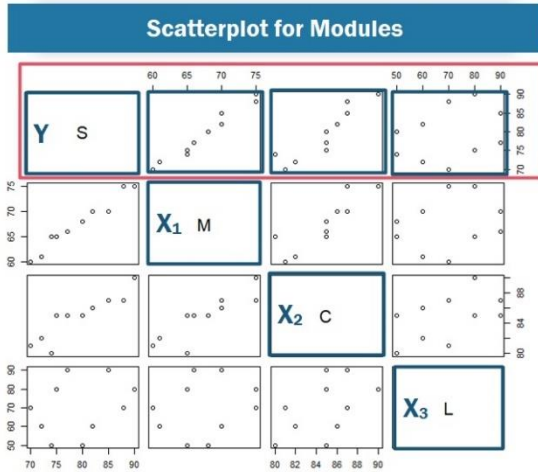
6 of 21

- To create the scatterplot in R:

- Enter the dataset in R
- Save the .txt file to your PC
- Upload it to R
- Use the "plot" command:

```
1 mods = read.table(file="C :\\Users\\nasge\\Desktop\\
\\External\\ 1 driveold\\R+Latex\\ lectdata\\
modules.txt",header=TRUE)
2 plot ( mods [, c("S", "M", "C", "L")])
```

### Scatterplot for Modules



We start with the usual first step. We enter the dataset in R and view the scatterplot matrix. We save the .txt file in our PC, upload it in R and use the "plot" command shown and obtain the scatterplot.

Let us read the scatter-matrix.


The scatter-matrix contains the response variable  $Y$  and the three predictors  $X_1, X_2, X_3$ .

Let's focus on the first row. We can observe some linear dependence between the grade in Statistics  $Y$  and the grade in Mathematics  $X_1$ . Something similar holds for the grade  $X_2$  in Computers too.

We can see a very sparse scatterplot between  $(X_3, Y)$ . This supports our instinct, that Students' skills in Languages, may not affect their performance in Statistics.



## Slide 7: Obtaining and Interpreting the Regression Equation



### Obtaining and Interpreting the Regression Equation

7 of 21

- Apply the `lm()` command in R:

```
1 modelF <- lm(S ~ . , data = mods )
2 summary (modelF)
```
- Include all the predictors and extract:  
Coefficients:

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	-34.675760	18.282690	-1.897	0.10666
M	1.042325	0.187688	5.554	0.00144 **
C	0.511827	0.360640	1.419	0.20564
L	0.003084	0.038032	0.081	0.93801

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

Residual standard error: 1.243 on 6 degrees of freedom

Multiple R-squared: 0.978 Adjusted R-squared: 0.9671

F-statistic: 89.11 on 3 and 6 DF, p-value: 2.295e-05

### Conclusion

!  $H_0: \beta_i = 0, i \neq 1$ , cannot be rejected.

- The significance of  $\beta_3$  is very low.
- There is a lack of satisfaction due to the low significance of the regression coefficients.
- The obtained model may be reconsidered.

We apply the "`lm()`" command in R as follows:


We start by including all the predictors that we have at our disposal. We then extract the following summary table.

Let's interpret the output carefully.

- The fit presents a very high  $R^2$  around 97.8%.
- The coefficient  $\beta_1$  of the  $X_1$ , is very significant. The p-value is around 1%, which is very low.
- There doesn't seem to be anything significant in the rest of the coefficients.
- Rigorously this means that the hypothesis  $H_0: \beta_i = 0$ , cannot be rejected, when "i" is not 1.
- The significance of  $\beta_3$ , which corresponds to the grade in Language, is really low.

Overall we are not completely satisfied, because of the low significance of the regression coefficients. We may reconsider the obtained model.

## Slide 8: Excluding Predictors



### Excluding Predictors

8 of 21

- Non-significant predictors can be excluded from the study.
- Re-fit, excluding  $X_3$  (Language grade):

```
1 modelF <- lm(S ~ M+C , data = mods )
2 summary (modelNoL)
```
- Output:  
Coefficients:

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	-35.5657	13.5447	-2.626	0.034120 *
M	1.0343	0.1480	6.987	0.000214 ***
C	0.5312	0.2500	2.124	0.071249 .

—


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

Residual standard error: 1.151 on 7 degrees of freedom

Multiple R-squared: 0.978, Adjusted R-squared: 0.9717

F-statistic: 155.8 on 2 and 7 DF, p-value: 1.573e-06

### Conclusion



The grade in Statistics highly depends on the grade in Mathematics.

- Further analysis is required to determine whether or not the grade for Computers should be excluded.

A predictor which ends up being non-significant could be excluded from the study. In our example, we re-fit excluding  $X_3$  (the grade in Languages) and we obtain the following summary table in R.

Interpreting the output, we have the following:

- The fit still presents a very high  $R^2 \approx 97.8\%$ .
- All the coefficients appear to be significant.
- The coefficient  $\beta_1$  of the  $X_1$ , is the most significant. This confirms statistically our initial speculations that the grade in Statistics, highly depends on the grade in Mathematics.
- The p-value for the test  $H_0: \beta_2 = 0$ , which corresponds to the grade in Computers, is around 7%, which is a low one, but still leaves the door open for further analysis. Is it a good idea to exclude the predictor  $X_2$ ?

## Slide 9: Excluding a Second Predictor

**Excluding A Second Predictor**

9 of 21

• Excluding the predictor  $X_2$ , results in a simple linear regression between  $Y$  and  $X_1$ .

• Re-fit, excluding  $X_2$  (Computer grade):

```
1 modelF <- lm(S ~ M, data = mods)
2 summary(modelF)
```

• Output:

Coefficients:

	Estimate	Std. Error	t value	Pr(>  t )
(Intercept)	-8.86038	6.05161	-1.464	0.181
M	1.30608	0.08942	14.606	4.74e-07 ***

---

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

Residual standard error: 1.381 on 8 degrees of freedom

Multiple R-squared: 0.9639, Adjusted R-squared: 0.9593

F-statistic: 213.3 on 1 and 8 DF, p-value: 4.735e-07

? What final model will be used for predictions?

We re-fit excluding  $X_2$ . Therefore, we have a simple linear regression between  $Y$  and  $X_1$ .

Interpreting the output of the SLR model obtained, we observe that:

- The fit still presents a very high  $R^2 \approx 96.4\%$ .
- The coefficient  $\beta_1$  of the  $X_1$ , is very significant. The intercept presents a lower significance. More carefully, there is a large standard error for it.

The last two obtained models look satisfactory. But what is the final model that we will use for predictions?

## Slide 10: Model Building

**Model Building**

10 of 21

Adjusted  $R^2$

Akaike's Information Criterion

**Introduction**

? Which predictor variable will be used for the final model?

- Let  $m \in \mathbb{N}$  and  $X_1, X_2, \dots, X_m$ , the potential predictors.
- Model building (variable selection) aims to find the predictors that will build the best possible model.

**Overfitting**  
(Too many predictors)

**Balanced**

**Underfitting**  
(Insufficient predictors)

Click each tab to learn more. Then, click Next to continue.

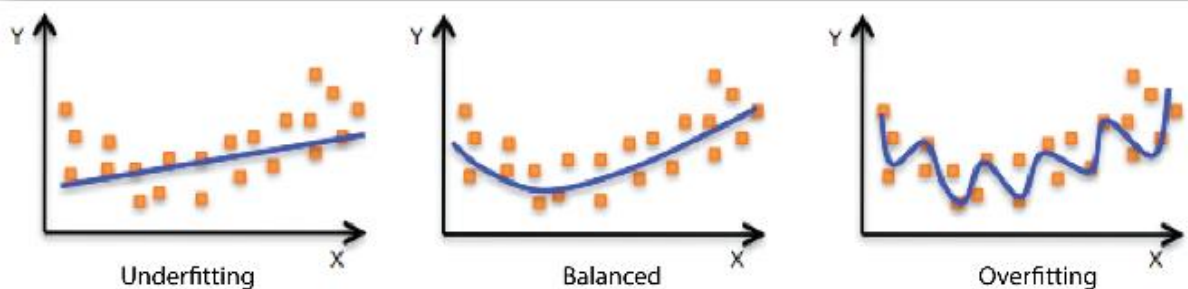


Nowadays, datasets can be easily found. In most cases, there are many candidates for predictor variables available. The natural question is which of them should be used for our final model?

Formally, let  $m \in \mathbb{N}$  and  $X_1, X_2, \dots, X_m$ , the potential predictors.

*Model building or Variable Selection*, aims to find the predictors that will build the best possible model.

Including too many predictors in the model, is referred to as over-fitting while the opposite case is called under-fitting.



There are several measures that help us arrive at the “best” model. We will look at two measures that are most commonly used: Adjusted  $R^2$  and Akaike’s Information Criterion.

Click the tabs to learn about each of these measures. When you are ready, click “Next” to continue.

#### Tab 1: Adjusted $R^2$

### Adjusted $R^2$

#### Determining the Best Variable

- $R^2$  increases when predictors are added to a model, even if those predictors are irrelevant.
  - Adjusted  $R^2$  is used because of this.
- Let  $p \in \mathbb{N}$ ,  $p \leq m$  is the number of predictors in a model.

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

It is common practice to use the model with the highest Adjusted  $R^2$ .

#### R Values for Adjusted $R^2$ in Case Study 1

<b>Model 1:</b> $Y; X_1, X_2, X_3$
• Multiple R-squared: 0.978, Adjusted R-squared: 0.9671
<b>Model 2:</b> $Y; X_1, X_2$
• Multiple R-squared: 0.978, Adjusted R-squared: 0.9717
<b>Model 3:</b> $Y; X_1$
• Multiple R-squared: 0.9639, Adjusted R-squared: 0.9593

We start with the so-called Adjusted  $R^2$ .

When we add predictors to a model,  $R^2$  increases - even if the predictors are completely irrelevant.

For this reason, we use the so-called Adjusted  $R^2$ . Let  $p \in \mathbb{N}$ ,  $p \leq m$ , be the number of the predictors included in a model. R provides the value of Adjusted  $R^2$  in the summary table.

The precise value of  $R_{adj}^2$  is as follows:

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

It is common practice to use the model with the highest Adjusted  $R^2$ .

If we go back to our case study on the student grades, the value for the Adjusted  $R^2$  was featured in the output beside the  $R^2$  value for all three of the models we ran. Those values are displayed on your screen now. Between the three models we examined, the one with the highest Adjusted  $R^2$  is the model including predictors  $X_1, X_2$ . Although, we can observe small differences in Adjusted  $R^2$  between the three models.

**Tab 2: Akaike's Information Criterion**

### Akaike's Information Criterion (1/3)

#### Introduction

- Akaike's information criterion (AIC) is the most common measure for variable selection.
- The smaller the AIC value, the better the model.

#### Formula

$$AIC = n \log \left( \frac{RSS}{n} \right) + 2p + n \log (2\pi e)$$

- As  $n \log (2\pi e)$  is independent of the model, R computes a simplified AIC:

$$AIC_R = n \log \left( \frac{RSS}{n} \right) + 2p$$

#### Implementing AIC in R

- R finds the model with the lowest AIC.
- The full model, using all available predictors, is:
$$Y = \beta_0 + \beta_1 X_1 + \dots + \beta_m X_m + e$$
- The null model, without predictors is:
$$Y = \beta_0 + e$$

#### stepAIC command can be run in two ways:

##### Forward Method

- Start with the null model and add predictors

##### Backward Method

- Start with the full model and exclude predictors

The most common measure for variable selection is Akaike's Information Criterion (AIC).

Its definition passes via the likelihood and it is based on balancing how well the model fits with a penalty for model complexity. *The smaller the value of AIC, the better the model.*

Precisely, for a model with  $p$  predictors, in a sample of  $n$ , the AIC formula is

$$AIC = n \log \left( \frac{RSS}{n} \right) + 2p + n \log (2\pi e)$$

The last term,  $n \log (2\pi e)$ , is independent of the model under consideration, therefore R computes the simplified version of AIC, by excluding this term as below.

Let us now implement AIC in R.

R finds the model with the lowest possible AIC and this is the one we finally work with.

Let's look at some of the terminology.

We refer to the model:

$Y = \beta_0 + \beta_1 X_1 + \dots + \beta_m X_m + e$ , as the full model (using all the available predictors) and to the model  $Y = \beta_0 + e$  as the null model (using no predictors).

There are two methods of finding AIC in R:

- The first is known as the forward method. stepAIC command in R starts with the null model and adds predictors until it finds the model with the lowest AIC.
- The second is the backward method, which is performed in the opposite way; starting from the full model and excluding predictors.

**Tab 2.1: Performing Model Selection**

### Akaike's Information Criterion (2/3)

#### Performing Model Selection in R

- Run the "stepAIC" command in R:

```
1 null = lm(S ~ 1, data = mods)
2 stepAIC(null, scope = list(lower = null, upper = modelF),
  data = mods, direction = 'forward')
```
- Output:

Start: AIC=39.43

	Df	Sum of Sq	RSS	AIC
S ~ 1				
+ M	1	406.84	15.26	8.224
+ C	1	348.14	73.96	24.010
<none>			422.10	39.427
+ L	1	36.45	385.65	

Step: AIC=8.22

S ~ M

	Df	Sum of Sq	RSS	AIC
+ C	1	5.9806	9.2756	5.2480
+ L	1	2.8804	12.3758	8.1316
<none>			15.2562	8.2240

Step: AIC=5.25

S ~ M + C

	Df	Sum of Sq	RSS	AIC
<none>			9.2756	5.248
+ L	1	0.010153	9.2654	7.237

Call:  
lm(formula = S ~ M + C, data = mods)

Coefficients:  
(Intercept)            M            C  
-35.5657            1.0343            0.5312

Preferred Model

Let us perform the Model selection in our case study.

We run the command stepAIC in R, deriving a large table that I have split in two to allow for better commenting on it.

R starts with the null model because we chose to work with a forward method. This model presents AIC = 39.43. Next R gives a new AIC for every single predictor, if we add new predictors to the model. You can see these highlighted on the screen.

R then adds the predictor which leads to the new model with the smallest AIC and continues the algorithm in the same way.

R continues this way until the addition of more predictors cannot offer a lower AIC.

R concludes to the preferred model and presents its coefficients at the end of the table.

Tab 2.2: Conclusion

### Akaike's Information Criterion (3/3)

#### Conclusion

- The final model based on the lowest AIC is:

$$\hat{S} = -35.567 + 1.0343 * M + 0.5312 * C$$

Confidence Intervals of Regression Coefficients		
> confint(modelNoL)		
	2.5 %	97.5 %
(Intercept)	-67.59379926	-3.537527
M	0.68428439	1.384406
C	-0.06004872	1.122487

#### Final Comments

- Predictor  $X_3$  has been justifiably excluded.
- Predictors  $X_1$  and  $X_2$  are significant.
- The grade in Mathematics appears to be more effective than the grade in Computers, due to the larger regression coefficient.
- The confidence intervals for the corresponding regression coefficients overlap.

The *final model* based on the lowest-AIC criterion is:

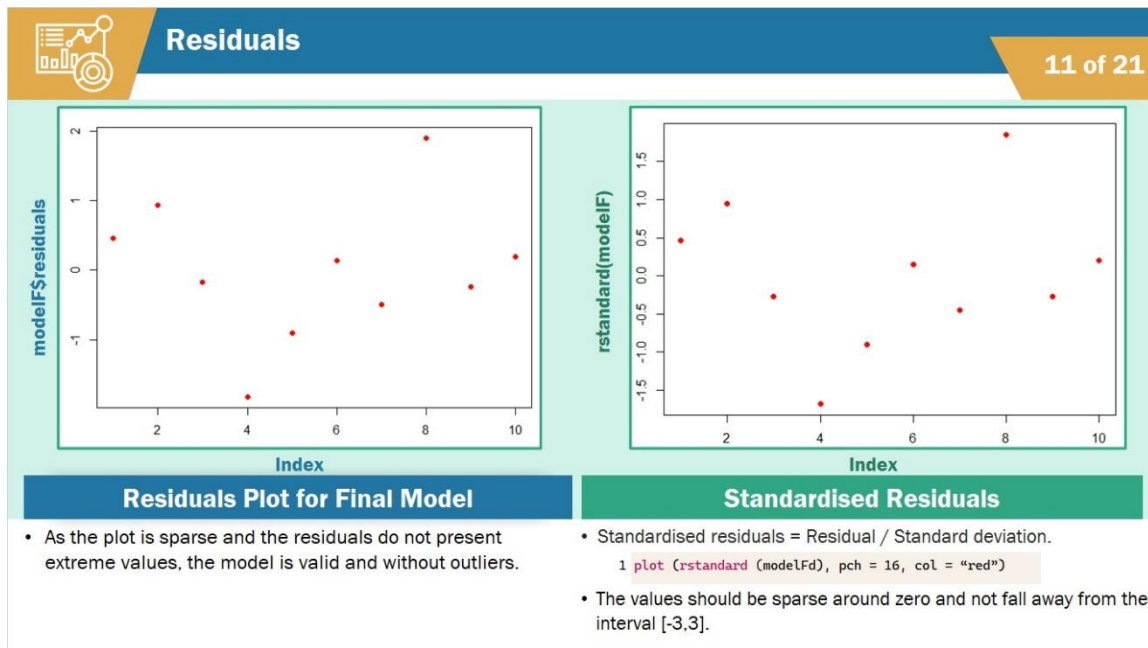
$\hat{S} = -35.567 + 1.0343 * M + 0.5312 * C$ , where S, M and C are the grades in Statistics, Mathematics and Computers respectively.

Having the final model in hand, we can find for example, the confidence intervals of the regression coefficients presented here.

Let us further comment on the final regression model:

- The predictor  $X_3$  has been justified statistically to be excluded.
- The predictors  $X_1$  and  $X_2$  are significant.
- The grade in Mathematics appears to be more effective than the grade in Computers because of the larger regression coefficient.
- The confidence intervals for the corresponding regression coefficients, overlap.

## Slide 11: Residuals



Let us extract the Residuals' plot for the final model, to validate it. As we can see, the plot is sparse, as it should be and the residuals do not present extreme values. This means that the model is valid and without outliers.

When we divide the residual by its standard deviation, we obtain the standardised residuals. Their values should be sparse around zero and not fall away from the interval [-3,3].

Here we can see that the standardised residuals are behaving properly, which validates the final model.

## Slide 12: Predictions

**Predictions** 12 of 21

- The model can be used for predictions.
- Example R Code:**  

```
1 a <- data.frame(M=70, C=75)
2 result <- predict(ModelNoL,a)
3 print(result)
```

**Math (X.) = 70%** **Computers (X.) = 75%** **Predicted Grade in Statistics (Y) = 76.67992%**

- Predictions only make sense for predictor values close to those of the dataset.
- Example R Code:**  

```
1 a <- data.frame(M=20, C=25)
2 result <- predict(ModelNoL,b)
3 print(result)
```

**Math (M) = 20%** **Computers (C) = 25%** **Predicted Grade in Statistics (Y) = -1.598287%**

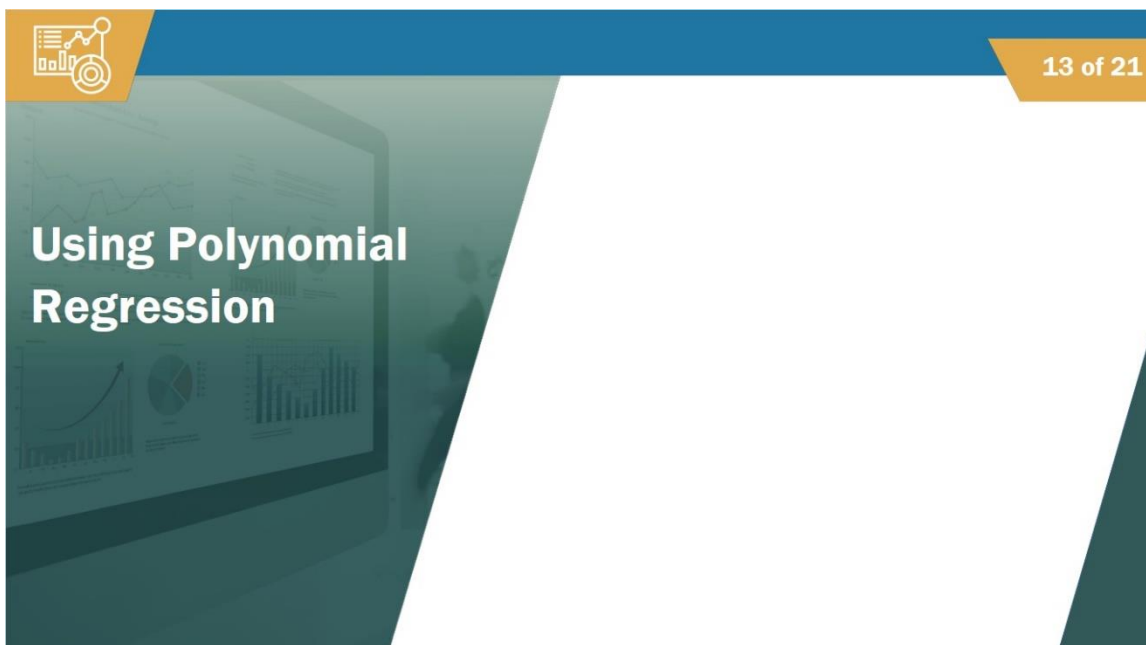


In practice, we can use the obtained model for predictions. For example, having a student with grades  $X_1=70$  (in Math) and  $X_2=75$  (in Computers), the model predicts a grade in Statistics as: around 76.7.

As a remark here, we have to mention that the predictions make sense for values of the predictors, “close” to those of the dataset. For example, if we want to predict the grade in Statistics, of a Student with grades  $M=20$  and  $C=25$  (far from the observed data values), we end up with an estimated grade in Statistics of around -1.6.


This negative grade in Statistics, makes no sense of course. The inconsistency is because we used predictors' values, that are far from the original dataset. This has to be avoided.

## Slide 13: Section 2: Case Study 2: Using Polynomial Regression



The slide features a blue header with a white icon of a bar chart and a circular arrow on the left, and the text "13 of 21" on the right. The main content area has a dark teal background on the left with a faint image of a computer monitor displaying various charts and graphs. The title "Using Polynomial Regression" is written in white, bold, sans-serif font over this background. The right side of the slide is a plain white area.


## Slide 14: Case Study 2: Salary From Years of Experience



### Case Study 2: Salary From Years of Experience

14 of 21

- **Case Study:** To model the salary of employees from their years of experience
- **Variable Y:** = Salary, should be explained by **Predictor X:** = Years of experience.
- There was initial salary growth, followed by a decline as the years of experience increased.
  - It is not a linear relationship so, what can we do?



Years	Salary
5	45
10	59
15	65
20	70
22	72
25	70
28	68
30	67
32	64
35	60

In our second case study, we would like to model the Salary of employees from their years of experience. We collected a dataset of  $n=10$  couples in the form (Years,Salary). The dataset is the following:


Let us collect some initial thoughts.

The variable  $Y$ := the Salary, should be explained by the predictor  $X$ := the years of experience. We have only one predictor variable at our disposal.

From the dataset, we may observe an initial growth of the Salary and later, a decline when more years of experience are added. This is, of course, not a case with a linear relationship.

So what can we do?

## Slide 15: Scatterplot



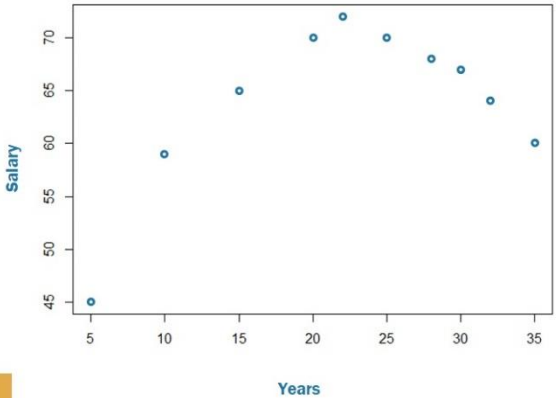
### Scatterplot

15 of 21

- To **create** the **scatterplot** in **R**:

```
1 mods = read.table(file="C:\\Users\\nasge\\Desktop\\External\\1 driveold\\R+Latex\\lectdata\\modules.txt", header=TRUE)  
2 plot(mods)
```

#### Years, Salary Scatterplot




! The scatterplot does not look like a line, so fit an SLR model.

We act as always: we enter the dataset in R and view the scatterplot. We save the .txt file in our PC, upload it in R and use the "plot" command.

The scatterplot does not look like a line.

What if we fit a SLR model?

## Slide 16: Fit SLR Model



### Fit SLR Model

16 of 21

- Fit the SLR model in R:

```
1 lm=lm(Salary ~ Years, data=mods)  
2 summary(lm)
```

- Output:

Residuals:					
	Min	1Q	Median	3Q	Max
	-11.4715	-3.3206	0.9006	4.6187	8.0875

Coefficients:				
	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	54.2830	5.7803	9.391	1.35e-05 ***
Years	0.4377	0.2402	1.822	0.106

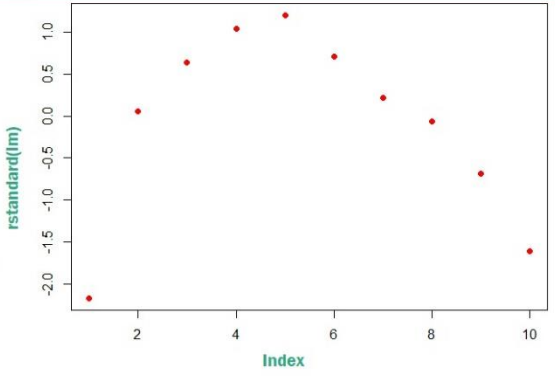
Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 7.058 on 8 degrees of freedom

Multiple R-squared: 0.2934, Adjusted R-squared: 0.205

F-statistic: 3.321 on 1 and 8 DF, p-value: 0.1059

#### Standardised Residuals Scatterplot



! The SLR is rejected.

We fit the SLR model in R, deriving this summary table. We ask R to plot the standardised residuals too, obtaining this graph.


Let us comment on the SLR:

- The fit presents a disappointing  $R^2 \approx 29\%$ .
- The predictor Years, the only one we have at our disposal, appears to be non-significant.
- The residuals present a very large range and the standard error is large.
- The standardised residuals are not sparse. They present a pattern. Not an appropriate residuals' curve.

For all these reasons, The SLR is rejected.

What can we do?




## Slide 17: Polynomial Regression Models



### Polynomial Regression Models

17 of 21

- When the scatterplot presents a non-linear pattern, try to fit a **polynomial regression model**.
  - You may attempt to involve  $X^k$  as well as  $X$ .


Types of Models	
Form	Name
$Y = \beta_0 + \beta_1 X + \beta_2 X^2 + e$	 Quadratic Model
$Y = \beta_0 + \beta_1 X + \beta_2 X^2 + \beta_3 X^3 + e$	 Cubic Model
$Y = \beta_0 + \beta_1 X + \dots + \beta_p X^p + e$	 Polynomial Model of degree $p \in \mathbb{N}$

When the scatterplot presents a non-linear pattern, we may try to fit a polynomial regression model. Except with the predictor  $X$ , we may attempt to involve its powers  $X^k$  too.

A model in the form of  $Y = \beta_0 + \beta_1 X + \beta_2 X^2 + e$ , is referred to as a quadratic model.

A model of the form  $Y = \beta_0 + \beta_1 X + \beta_2 X^2 + \beta_3 X^3 + e$ , is referred to as a cubic model and more generally, a model  $Y = \beta_0 + \beta_1 X + \dots + \beta_p X^p + e$ , is referred to as a polynomial model of degree (or order)  $p$ .

## Slide 18: Implementing a Quadratic Polynomial Regression Model in R



Implementing a Quadratic Polynomial Regression Model in R

18 of 21

- Obtain a **quadratic regression model** for the dataset:

```
1 lm2=lm(Salary ~ Years+I(Years^2) , data=exper)
2 summary (lm2)
```
- Output:

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	29.626275	1.355336	21.86	1.06e-07 ***
Years	3.590507	0.148033	24.25	5.15e-08 ***
I(Years^2)	-0.078273	0.003594	-21.78	1.09e-07 ***

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

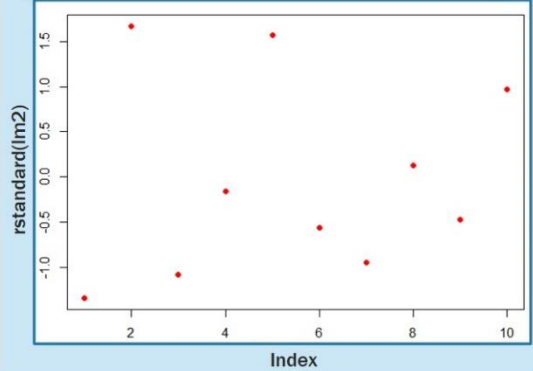
Residual standard error: 0.9099 on 7 degrees of freedom  
Multiple R-squared: 0.9897 Adjusted R-squared: 0.9868

**The quadratic model looks ideal:**

!  $\widehat{Salary} \approx 29.6 + 3.6Years - 0.08Years^2$
- Plot the standardise residuals in R:

```
1 plot(rstandard (lm2), pch = 16, col = "red")
```

Standardise Residuals Scatterplot



We use R as follows to obtain a quadratic regression model for the Salary-Experience dataset. Note that in the `lm()` command, we add a new predictor  $Years^2$ . Then R runs a multiple regression with predictors  $X_1:=Years$  and  $X_2:=Years^2$ . We derive the following summary table.

We then ask R to provide us with the plot of the standardised residuals and get the next graph.

Let's comment on the polynomial fit.


- The summary table presents an ideal  $R^2 \approx 99\%$  and high significance in all the regression coefficients.
- The residuals are sparse around zero - an optimal residuals' plot.

The *quadratic model* looks ideal and has the form:

$$\widehat{Salary} \approx 29.6 + 3.6Years - 0.08Years^2$$



## Slide 19: Implementing a Cubic Polynomial Regression Model in R



Implementing a Cubic Polynomial Regression Model in R

19 of 21


- Fit a **cubic regression model** using:

```
1 lm3=lm(Salary ~ Years+I(Years^2)+I(Years^3) , data=exper)
2 summary (lm3)
```
- Output:

Coefficients:				
	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	27.3894385	2.5183228	10.876	3.58e-05 ***
Years	4.0893725	0.4969974	8.228	0.000174 ***
I(Years^2)	-0.1068284	0.0274101	-3.897	0.008008 **
I(Years^3)	0.0004714	0.0004487	1.051	0.333847

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

Residual standard error: 0.9032 on 6 degrees of freedom  
Multiple R-squared: 0.9913. Adjusted R-squared: 0.987


 The cubic model is not suitable.

Let us attempt to fit a cubic polynomial regression model with this command. We obtain the following table.

We can see that the  $R^2$  is very large, as before, but the coefficient  $\beta_3$  is not significant. It may well be zero. We do not use the cubic model.

We will stay with the quadratic one obtained before. We can now use it for predictions.

## Slide 20: Conclusions Derived from the Quadratic Model



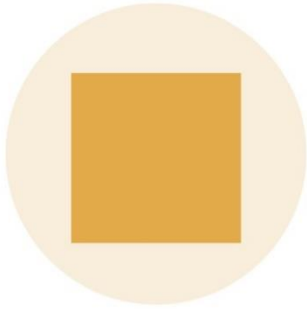
Conclusions Derived From the Quadratic Model

20 of 21

- The **quadratic model** is:
$$\widehat{Salary} \approx 29.6 + 3.6\text{Years} - 0.08\text{Years}^2$$
- To predict the salary of a person with  $X=7$  years experience, the model estimates:
  - $\widehat{Salary} \approx 51$

**Confidence Intervals of Regression Coefficients**

	2.5 %	97.5 %
(Intercept)	26.42141373	32.83113556
Years	3.24046500	3.94054934
I(Years^2)	-0.08677072	-0.06977463




For a prediction of the salary of a person with  $X=7$  years of experience, the model estimates:

$$\widehat{\text{Salary}} \approx 51$$

The confidence intervals for the regression coefficients are listed here.

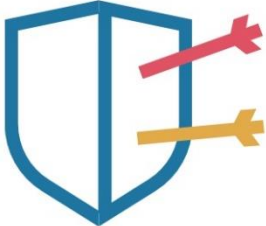
## Slide 21: Summary



Summary

21 of 21

- Having completed this presentation, you should be able to:
  - Use R for multiple linear regression
  - Conclude the final model using AIC
  - Validate the obtained model using residuals
  - Predict values of the response variable
  - Perform polynomial regression in R



Developed by Trinity Online Services CLG with the School of Computer Science and Statistics, Trinity College Dublin, The University of Dublin

Having completed this presentation, you should be able to:

- Use R for multiple linear regression
- Conclude the final model using AIC
- Validate the obtained model using residuals
- Predict values of the response variable
- Perform polynomial regression in R