

Aula 7:

'Que fatores explicam a variação nos salários na organização?

Validação de Modelos de Regressão Linear

Docente: Daniela Craveiro

dcraveiro@iseg.ulisboa.pt



No final desta aula,

@s alun@s deverão:

- Perceber qual a necessidade de fazermos diagnósticos aos pressupostos do nosso modelo de regressão
- Saber quais são os pressupostos do modelo de regressão linear
- Saber como, com a ajuda de gráficos e testes estatísticos, podemos conferir se os pressupostos do modelo estão a ser cumpridos
- Saber implementar os estudo dos pressupostos do modelo no SPSS



	PRESSUPOSTOS	DEFINIÇÃO	FORMA DE VALIDAÇÃO
I	Linearidade	O efeito das variáveis independentes na variável dependente é linear e aditivo.	 Análise gráfica (Matriz de Dispersão, por exemplo)



	PRESSUPOSTOS	DEFINIÇÃO	FC	DRMA DE VALIDAÇÃO
I	Linearidade	O efeito das variáveis independentes na variável dependente é linear e aditivo.	•	Análise gráfica (Matriz de Dispersão, por exemplo)
II	Normalidade da Distribuição dos Erros	Os erros seguem uma distribuição normal.	•	Análise de Resíduos Gráfico de Q-Q



	PRESSUPOSTOS	DEFINIÇÃO	FORMA DE VALIDAÇÃO
I	Linearidade	O efeito das variáveis independentes na variável dependente é linear e aditivo.	 Análise gráfica (Matriz de Dispersão, por exemplo)
II	Normalidade da Distribuição dos Erros	Os erros seguem uma distribuição normal.	Análise de ResíduosGráfico de Q-Q
Ш	Média Condicional Zero dos Erros	O termo de erro aleatório tem valor esperado igual a zero.	• Análise de Resíduos



	PRESSUPOSTOS	DEFINIÇÃO	FORMA DE VALIDAÇÃO
I	Linearidade	O efeito das variáveis independentes na variável dependente é linear e aditivo.	 Análise gráfica (Matriz de Dispersão, por exemplo)
II	Normalidade da Distribuição dos Erros	Os erros seguem uma distribuição normal.	Análise de ResíduosGráfico de Q-Q
Ш	Média Condicional Zero dos Erros	O termo de erro aleatório tem valor esperado igual a zero.	Análise de Resíduos
IV	Homocedasticidade (ou Igual Variância)	A distribuição dos erros apresenta uma variância constante (hipótese da homocedasticidade).	Análise de Resíduos



	•		
	PRESSUPOSTOS	DEFINIÇÃO	FORMA DE VALIDAÇÃO
I	Linearidade	O efeito das variáveis independentes na variável dependente é linear e aditivo.	 Análise gráfica (Matriz de Dispersão, por exemplo)
II	Normalidade da Distribuição dos Erros	Os erros seguem uma distribuição normal.	Análise de ResíduosGráfico de Q-Q
III	Média Condicional Zero dos Erros	O termo de erro aleatório tem valor esperado igual a zero.	Análise de Resíduos
IV	Homocedasticidade (ou Igual Variância)	A distribuição dos erros apresenta uma variância constante (hipótese da homocedasticidade).	Análise de Resíduos
V	Independência dos Erros	Os erros não estão correlacionados, i.e., o valor de um erro não depende de qualquer outro erro.	 Dublin-Watson



	PRESSUPOSTOS	DEFINIÇÃO	FORMA DE VALIDAÇÃO
I	Linearidade	O efeito das variáveis independentes na variável dependente é linear e aditivo.	 Análise gráfica (Matriz de Dispersão, por exemplo)
II	Normalidade da Distribuição dos Erros	Os erros seguem uma distribuição normal.	Análise de ResíduosGráfico de Q-Q
Ш	Média Condicional Zero dos Erros	O termo de erro aleatório tem valor esperado igual a zero.	Análise de Resíduos
IV	Homocedasticidade (ou Igual Variância)	A distribuição dos erros apresenta uma variância constante (hipótese da homocedasticidade).	Análise de Resíduos
V	Independência dos Erros	Os erros não estão correlacionados, i.e., o valor de um erro não depende de qualquer outro erro.	 Dublin-Watson
VI	Ausência de multicolinearidade perfeita	As variáveis independents não estão perfeitamente correlacionadas entre si.	 Diagnósticos de Colinearidade



	PRESSUPOSTOS	DEFINIÇÃO	FORMA DE VALIDAÇÃO
1	Linearidade	O efeito das variáveis independentes na variável dependente é linear e aditivo.	 Análise gráfica (Matriz de Dispersão, por exemplo)
I	Normalidade da Distribuição dos Erros	Os erros seguem uma distribuição normal.	Análise de ResíduosGráfico de Q-Q
II	Média Condicional Zero dos Erros	O termo de erro aleatório tem valor esperado igual a zero.	Análise de Resíduos
V	Homocedasticidade (ou Igual Variância)	A distribuição dos erros apresenta uma variância constante (hipótese da homocedasticidade).	Análise de Resíduos
/	Independência dos Erros	Os erros não estão correlacionados, i.e., o valor de um erro não depende de qualquer outro erro.	• Dublin-Watson
/ I	Ausência de multicolinearidade perfeita	As variáveis independents não estão perfeitamente correlacionadas entre si.	 Diagnósticos de Colinearidade
/11	Ausência de Observações Influentes	Não existem observações que tenham uma influência anormal nos resultados do modelo.	Cook's Distance



E qual é o problema se estes pressupostos não se verificarem?

 Os intervalos de confiança ou os p-values podem estar a ser subestimados (i.e. mais pequenos do que na realidade são) ...

ou seja: estamos a atribuir significância estatística a uma estimativa que na realidade não a terá!



Validação do Modelo de Regressão Linear

Avaliação do Pressuposto I: Linearidade

Ver aula anterior: análise gráfica (Matriz de Dispersão)



Validação do Modelo de Regressão Linear

- 1. Estimar o modelo de regressão com os diagnósticos
- 2. Avaliação do Pressuposto II: Normalidade da Distribuição dos Erros
- 3. Avaliação do Pressuposto III: Média Condicional Zero dos Erros
- 4. Avaliação do Pressuposto IV: Homocedasticidade
- 4. Avaliação do Pressuposto V: Independência dos Erros
- 5. Avaliação do Pressuposto VI: Ausência de Multicolinearidade Perfeita
- 6. Avaliação do Pressuposto VII: Ausência de Observações Influentes



Validação do Modelo de Regressão Linear

Estimar o modelo de regressão com os diagnósticos

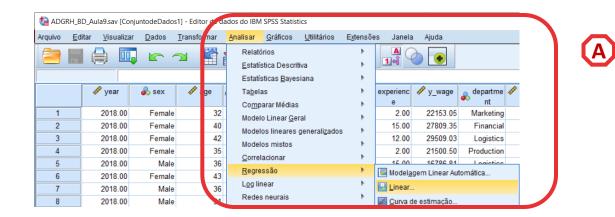
Objetivo: Estimar o modelo com informação adicional para avaliar os pressupostos (VIs: 'sex_female', 'education2' e 'evaluation')

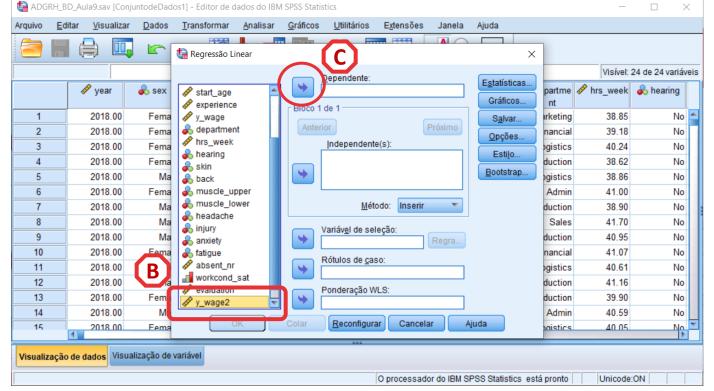
 Selecionar 'Analisar' / 'Regressão' / 'Linear'

- A
- Selecionar a variável 'y_wage2'
- B
- Colocar na caixa 'Dependente'



Exercício: Colocar as variáveis 'sex_female', 'education2' e 'evaluation' na caixa 'Independente(s)'





 Selecionar 'Analisar' / 'Regressão' / 'Linear' A

Selecionar a variável 'y_wage2'



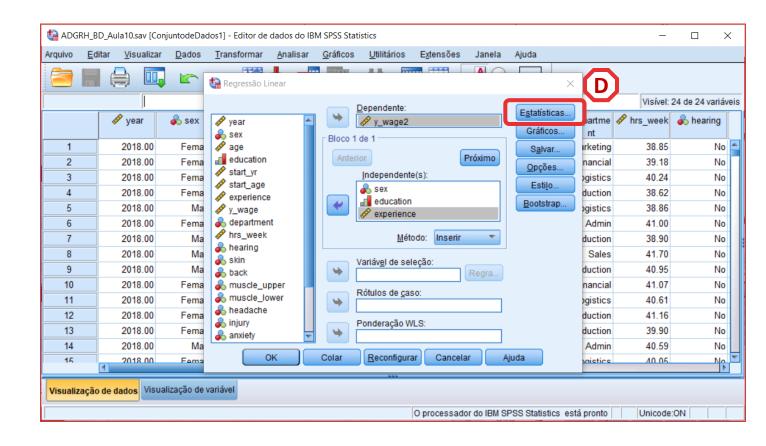
Colocar na caixa 'Dependente'



Exercício: Colocar as variáveis 'sex_female', 'education2' e 'evaluation' na caixa 'Independente(s)'

Selecionar botão 'Estatísticas'





 Selecionar 'Analisar' / 'Regressão' / 'Linear' A

Selecionar a variável 'y_wage2'



Colocar na caixa 'Dependente'



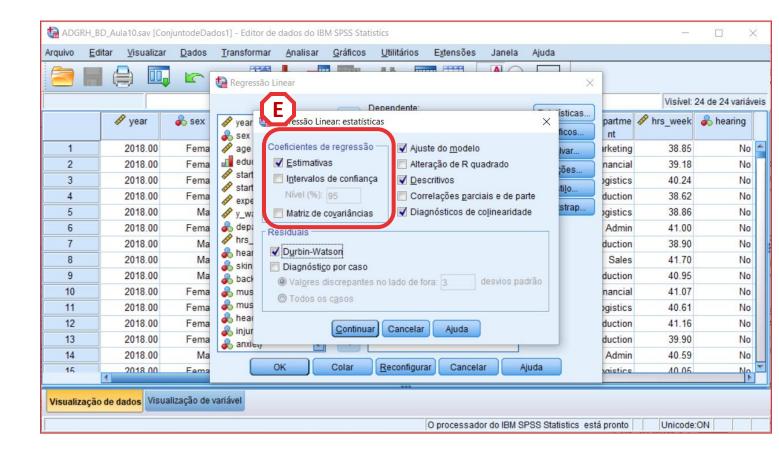
Exercício: Colocar as variáveis 'sex_female', 'education2' e 'evaluation' ' na caixa 'Independente(s)'

Selecionar botão 'Estatísticas'



Selecionar 'Estimativas'





Selecionar 'Ajuste do modelo'

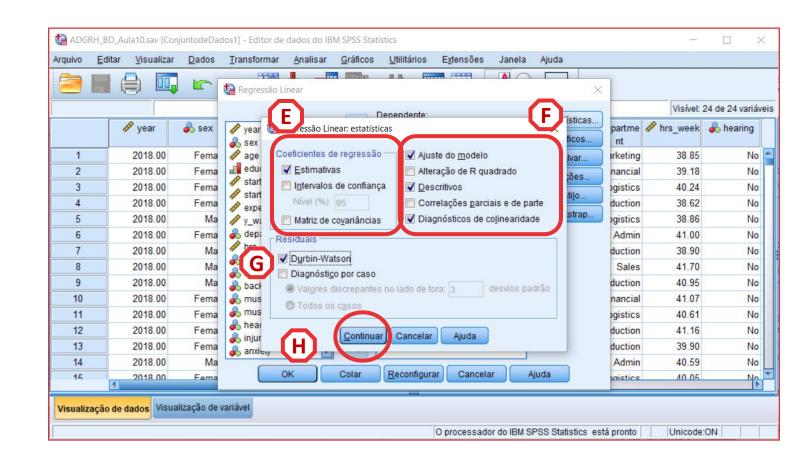


- Selecionar 'Descritivos'
- Selecionar 'Diagósticos de colinearidade'
- Selecionar 'Dublin-Watson'

G

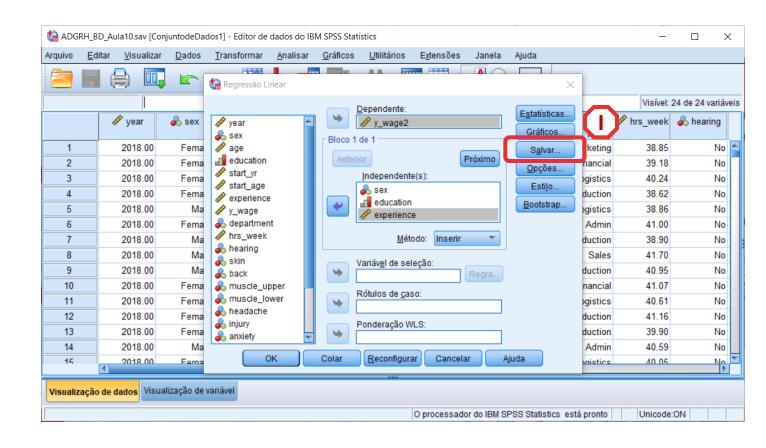
Selecionar 'Continuar'





Selecionar botão 'Salvar'

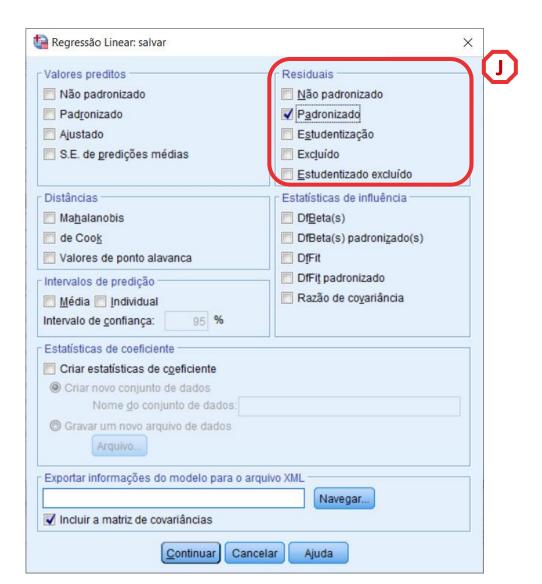




Selecionar botão 'Salvar'

Selecionar 'Padronizado'



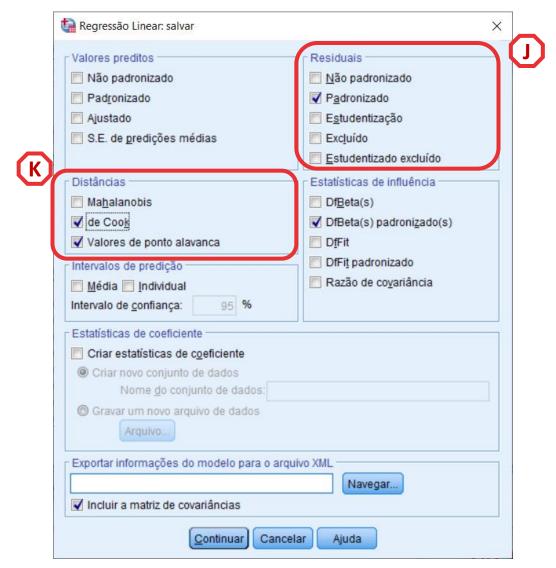


Selecionar botão 'Salvar'

Selecionar 'Padronizado'

(J

 Selecionar 'de Cook' e 'Valores de ponto alavanca'



Selecionar botão 'Salvar'

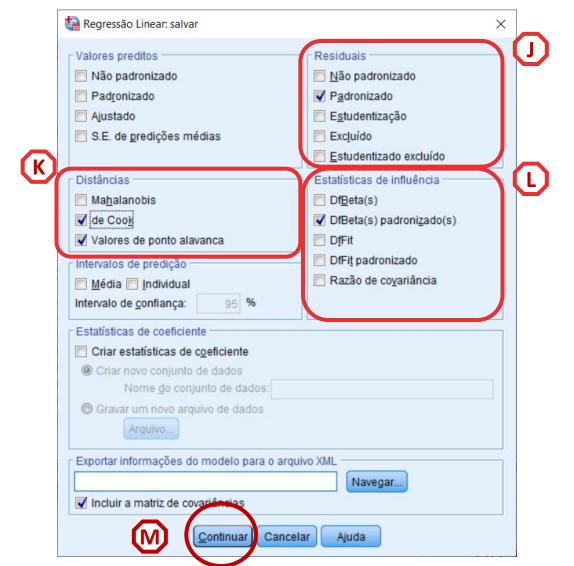
Selecionar 'Padronizado'

- (J)
- Selecionar 'de Cook' e 'Valores de ponto alavanca'
- K

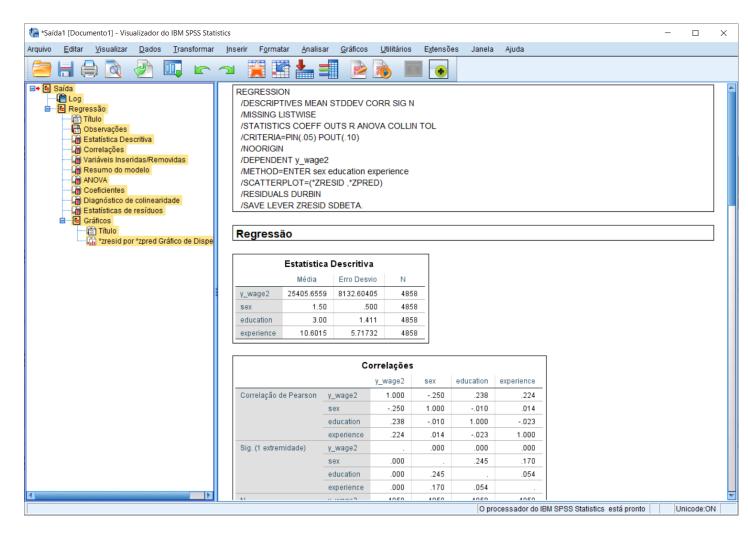
 Selecionar 'DfBeta(s) padronizado(s) (L)

(M)

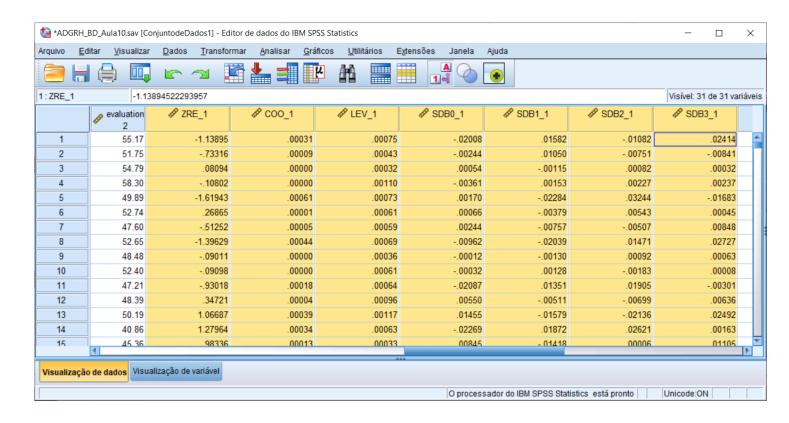
Selecionar 'Continuar'/ 'OK'



 Os resultados são publicados no 'Visualizador de Resultados'

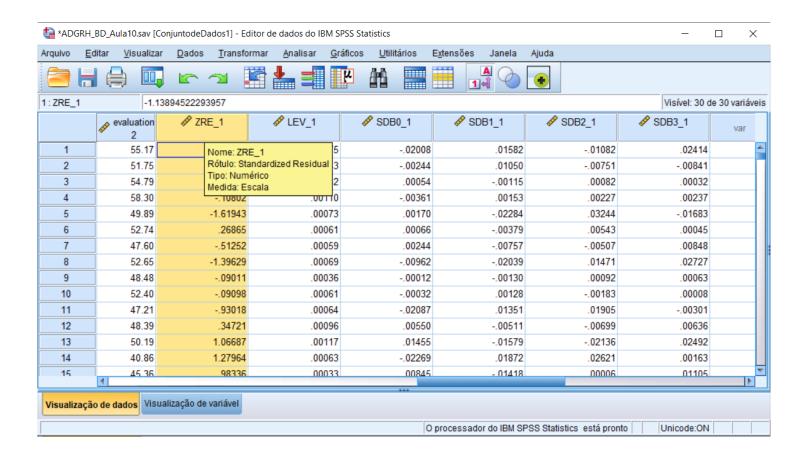


 Quando instruímos o SPSS para produzir os diagnósticos, é criado um conjunto variáveis



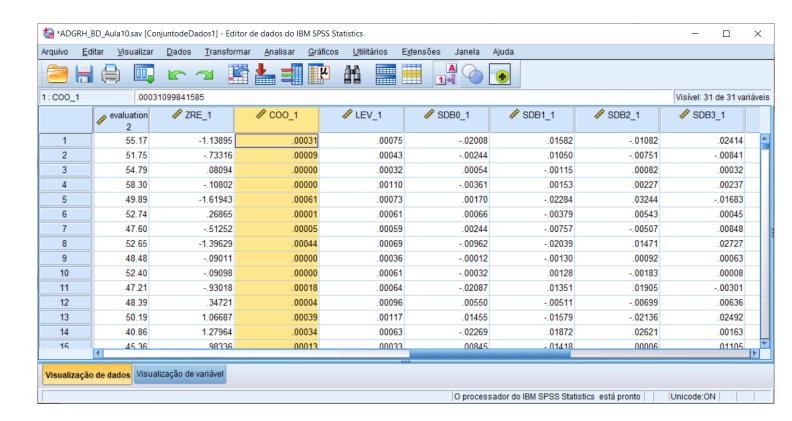


- Quando instruímos o SPSS para produzir os diagnósticos, é criado um conjunto variáveis
 - Uma variável com os 'Resíduos Padronizados' da variável dependente (ZRE_1) para cada observação



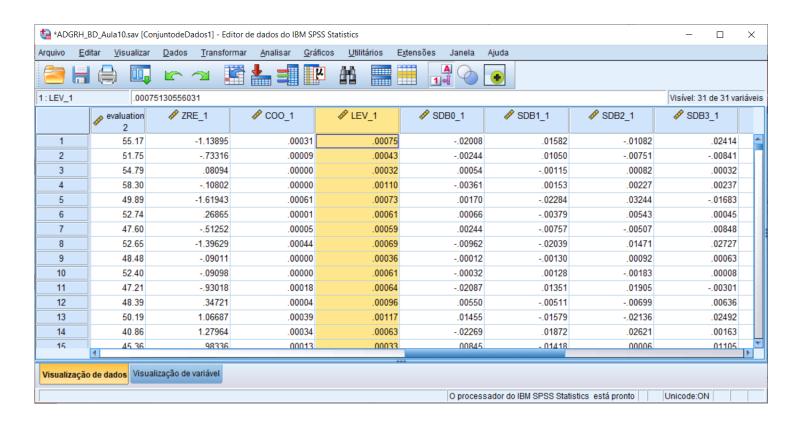


- Quando instruímos o SPSS para produzir os diagnósticos, é criado um conjunto variáveis
 - Uma variável com os 'Resíduos Padronizados' da variável dependente (ZRE_1) para cada observação
 - Uma variável que mede a distancia de Cook associada a cada observação (COO 1)



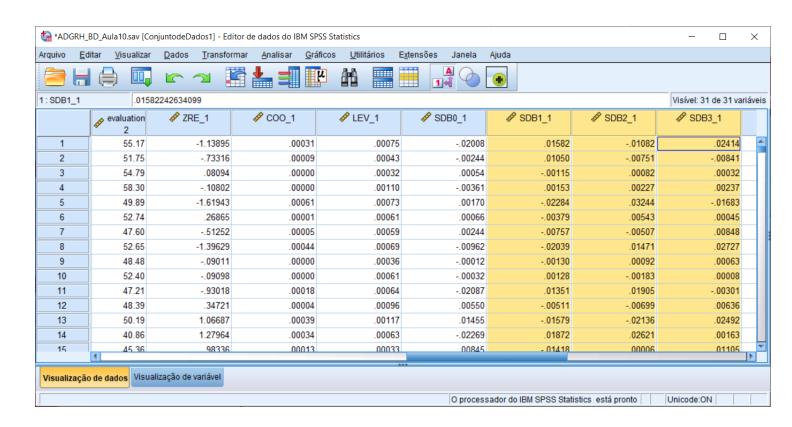


- Quando instruímos o SPSS para produzir os diagnósticos, é criado um conjunto variáveis
 - Uma variável com os 'Resíduos Padronizados' da variável dependente (ZRE_1) para cada observação
 - Uma variável que mede a distancia de Cook associada a cada observação (COO 1)
 - Uma variável que mede influência relativa de cada observação no ajuste do modelo (LEV_1).





- Quando instruímos o SPSS para produzir os diagnósticos, é criado um conjunto variáveis
 - Uma variável com os 'Resíduos Padronizados' da variável dependente (ZRE_1) para cada observação
 - Uma variável que mede a distância de Cook associada a cada observação (COO 1)
 - Uma variável que mede influência relativa de cada observação no ajuste do modelo (LEV_1).
 - Por cada variável independente é criada uma variável com os DFBETA Padronizado, mede a influência de uma dada observação na estimação dos parâmetros.





Validação do Modelo de Regressão Linear

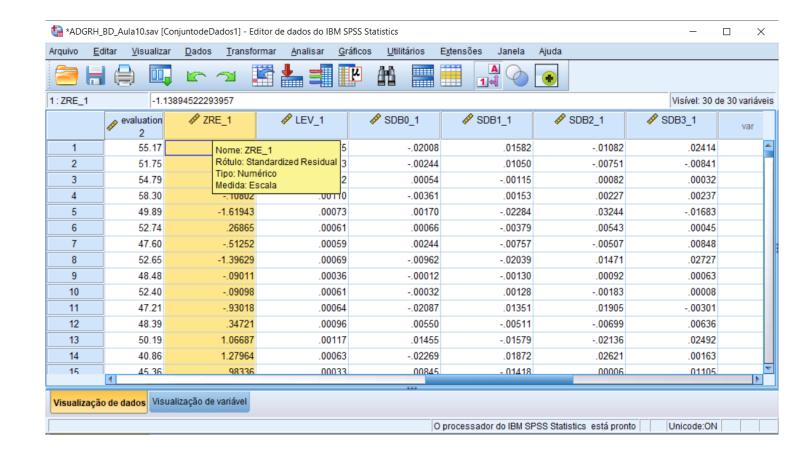
2. Avaliação do Pressuposto II: Normalidade da Distribuição dos Erros



Validação do Modelo de Regressão Linear

Avaliação do Pressuposto II: Normalidade da Distribuição dos Erros

- Para avaliarmos se os erros seguem uma distribuição normal, vamos usar a variável com os 'Resíduos Padronizados' da VD (ZRE_1) que acabamos de criar.
- Vamos então criar usar um gráfico Q-Q para representar a distribuição dos resíduos padronizados



Selecionar 'Analisar' /
 'Estatística Descritiva' /
 'Gráficos Q-Q'



Selecionar a variável 'ZRE_1'

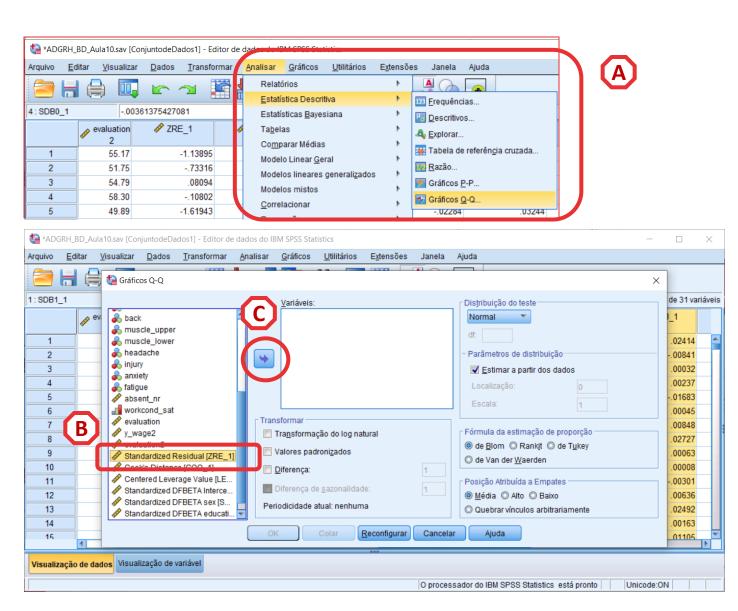
B

Colocar na caixa 'Variáveis'

C

Selecionar 'OK'





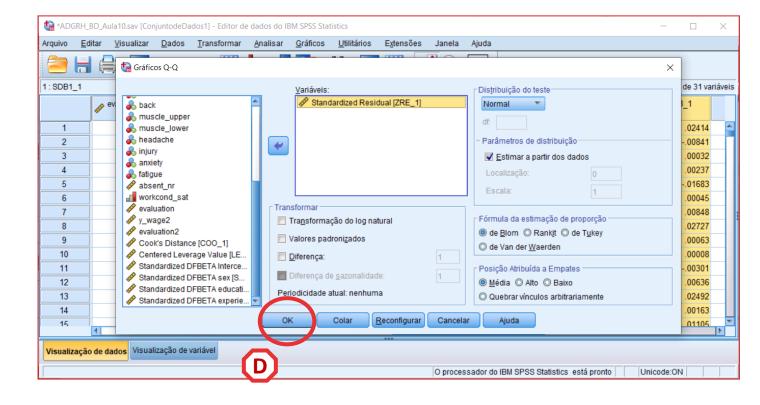
- Selecionar 'Analisar' /
 'Estatística Descritiva' /
 'Gráficos Q-Q'
- Selecionar a variável 'ZRE_1'
- Colocar na caixa 'Variáveis'
- Selecionar 'OK'



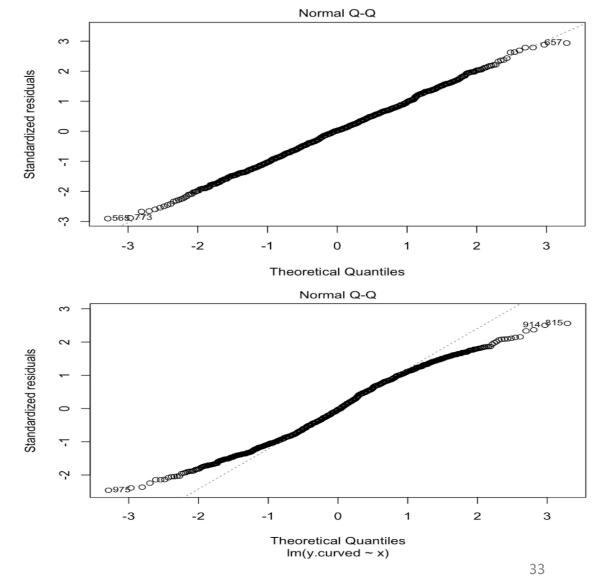




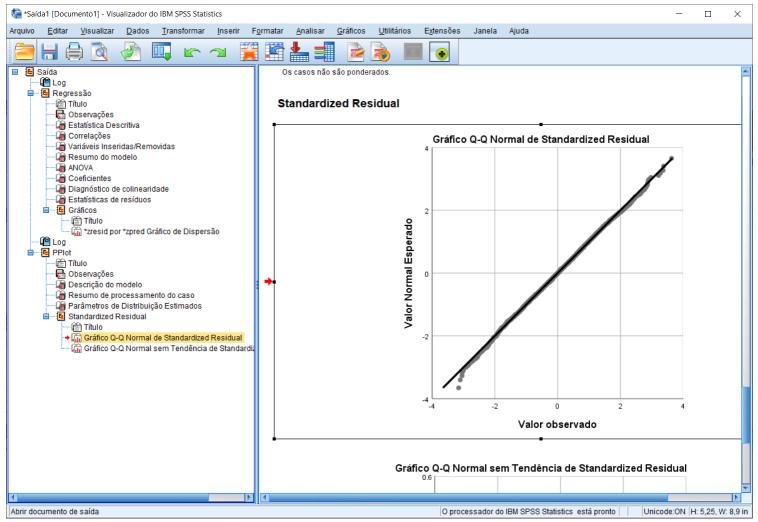




- Linha diagonal reflecte uma distribuição normal
- Os resíduos sobrepõe-se quase totalmente com a linha de diagonal
- Os resíduos parecem estar normalmente distribuídos
- Neste, caso os as caudas da distribuição dos resíduos afasta-se da diagonal, o que sugere que a distribuição dos erros pode não ser normal



- O gráfico é publicado no 'Visualizador de Resultados'
- Neste caso podemos concluir que os erros seguem uma distribuição normal!



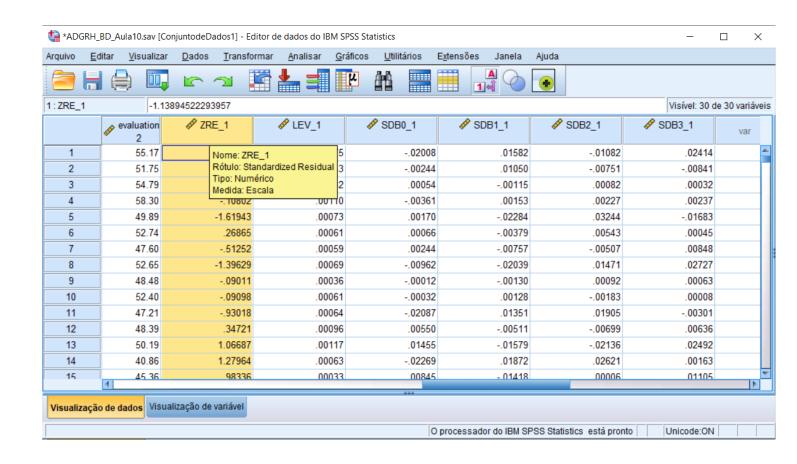


Validação do Modelo de Regressão Linear

Avaliação do Pressuposto III: Média Condicional Zero dos Erros

Média Condicional Zero dos Erros

- Para avaliarmos se o termo de erro aleatório tem valor esperado igual a zero, vamos usar a variável com os 'Resíduos Padronizados' da VD (ZRE_1) que acabamos de criar.
- Mas neste caso, vamos olhar para as estatísticas descritivas desta variável.



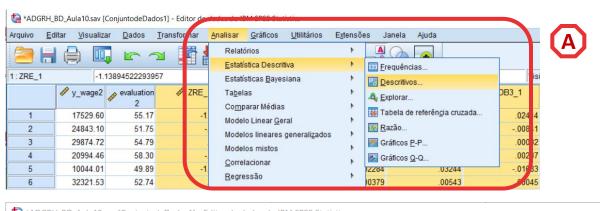
Média Condicional Zero dos Erros

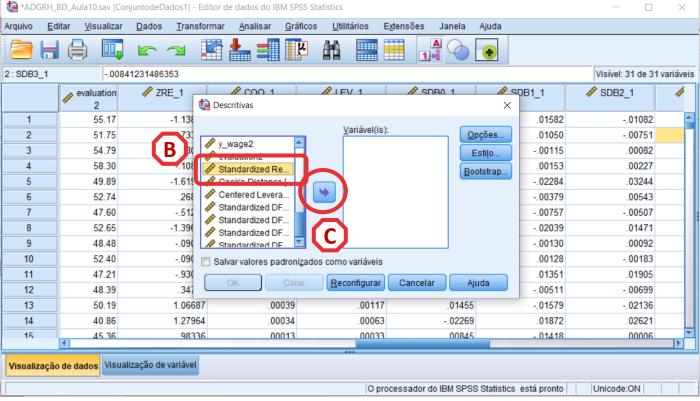
- Selecionar 'Analisar' /
 'Estatística Descritiva' /
 'Descritivos'
 - / ' /

- Selecionar a variável 'ZRE_1'
- B
- Colocar na caixa 'Variável(is)'
- **(C)**

Selecionar 'OK'







Média Condicional Zero dos Erros

Selecionar 'Analisar' /
 'Estatística Descritiva' /
 'Descritivos'



Selecionar a variável 'ZRE_1'

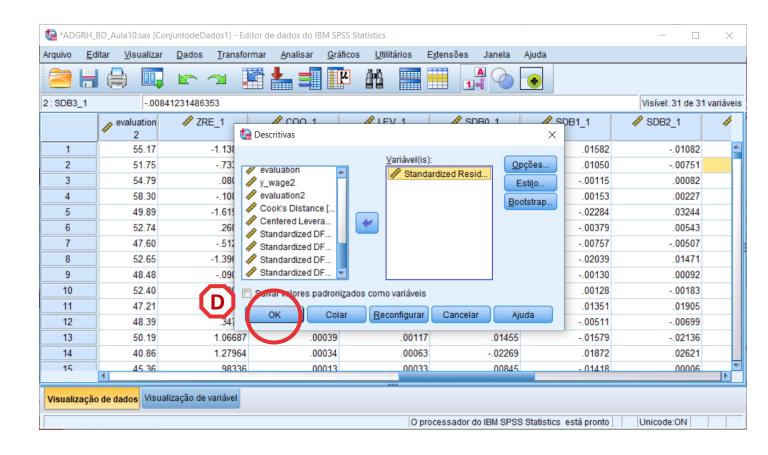
B

Colocar na caixa 'Variável(is)'

C

Selecionar 'OK'

(D)

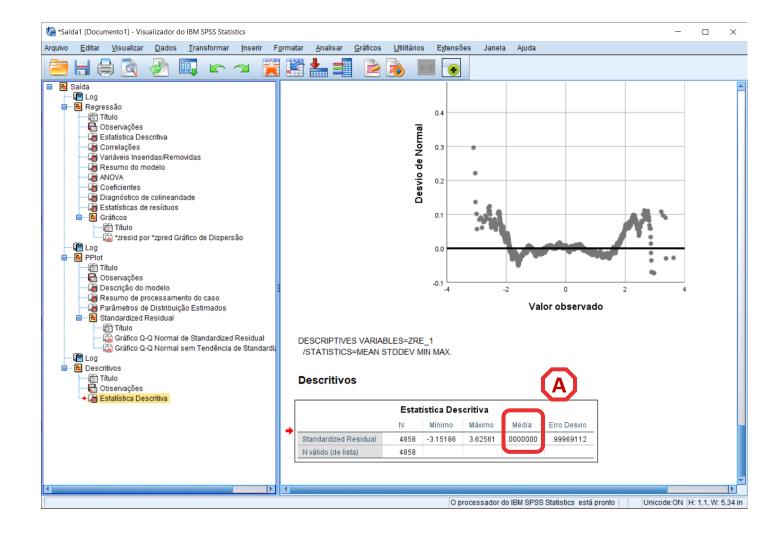


Média Condicional Zero dos Erros

- O gráfico é publicado no 'Visualizador de Resultados'
- Os 'Resíduos Padronizados' da VD (ZRE_1) tem uma média muito próximo de 0,



 Neste caso podemos concluir que se cumpre o pressuposto da Média Condicional Zero dos Erros.





Validação do Modelo de Regressão Linear

Avaliação do Pressuposto IV: Homocedasticidade



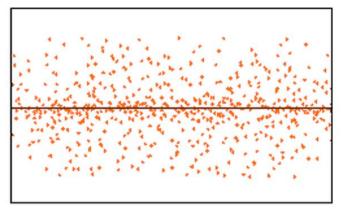
Homocedasticidade (ou Igual Variância)

 A distribuição dos resíduos apresenta uma variância constante ao longo dos valores previstos da variável dependente. Não há indicação de variação nãoconstante.

 Neste, o valor dos resíduos aproxima-se de 0 para os valores mais baixos da predição, mas aumentam à medida que os valores previstos também aumenta.

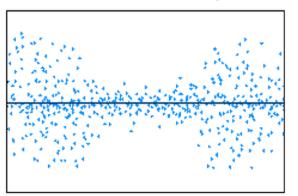
ou seja, a variação não é constante.

Homoscedasticity



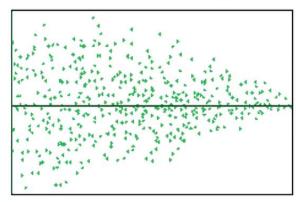
Random Cloud (No Discernible Pattern)

Heteroscedasticity



Bow Tie Shape (Pattern)

Heteroscedasticity

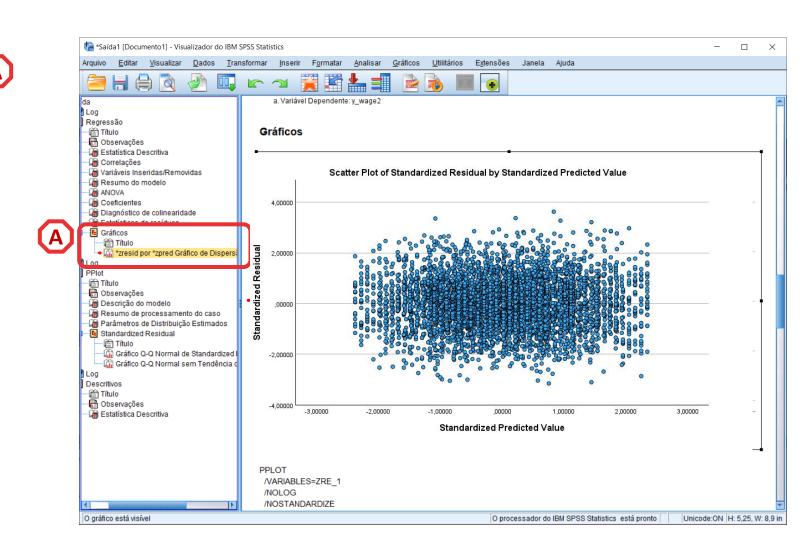


Fan Shape (Pattern)



Homocedasticidade

- Para avaliar se se cumpre este pressuposto, temos de olhar para o Gráfico de Dispersão que compara a distribuição dos 'Residuos Padronizados' com os 'Valores Preditos Padronizados' - que o SPSS produz automaticamente.
- Neste caso, a representação da distribuição parece sugerir que a variação dos resíduos é relativamente constante.
- Ou seja, cumpre-se o pressuposto da Homocedasticidade



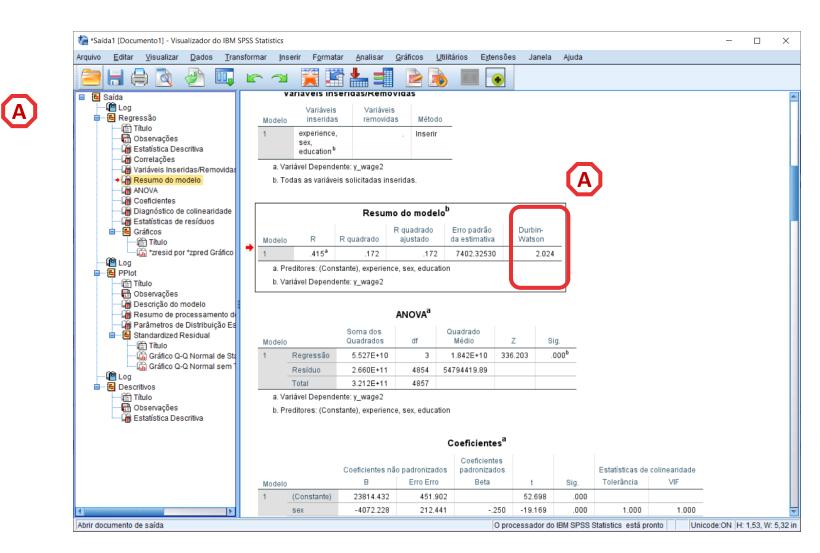


Validação do Modelo de Regressão Linear

Avaliação do Pressuposto V: Independência dos Erros

Independência dos Erros

- Para avaliar se se cumpre este pressuposto, temos de olhar para o resultado do teste Durbin-Watson - que pedimos ao SPSS para produzir.
- Interpretação:
- = 2 -> Erros <u>são</u> independentes
- > 2 / <2 -> Erros <u>não são</u> independentes
- Neste caso aqui apresentado os erros são independentes



Independência dos Erros

 Na nossa base de dados contudo, o pressuposto não é assegurado

teste Durbin-Watson <2

Variáveis Inseridas/Removidasa

Modelo	Variáveis inseridas	Variáveis removidas	Método
1	education2, sex=Female, experience ^b		Inserir

- a. Variável Dependente: y_wage2
- b. Todas as variáveis solicitadas inseridas.

Resumo do modelo^b

	Modelo	R	R quadrado	R quadrado ajustado	Erro padrão da estimativa	Durbin-V	Vatson
•	1	,415 ^a	,172	,172	7402,32530		,528
	a Proditoroe: (Constanto) oducation2 cov-Formalo ovnoriones						

- a. Preditores: (Constante), education2, sex=Female, experience
- b. Variável Dependente: y_wage2

ANOVA^a

Modelo		Soma dos Quadrados	df	Quadrado Médio	F	Sig.
1	Regressão	55266216704	3	18422072235	336,203	<,001 ^b
	Resíduo	2,660E+11	4854	54794419,887		
	Total	3,212E+11	4857			

a. Variável Dependente: y_wage2



Validação do Modelo de Regressão Linear

Avaliação do Pressuposto VI: Ausência de Multicolinearidade Perfeita



Ausência de Multicolienaridade Perfeita

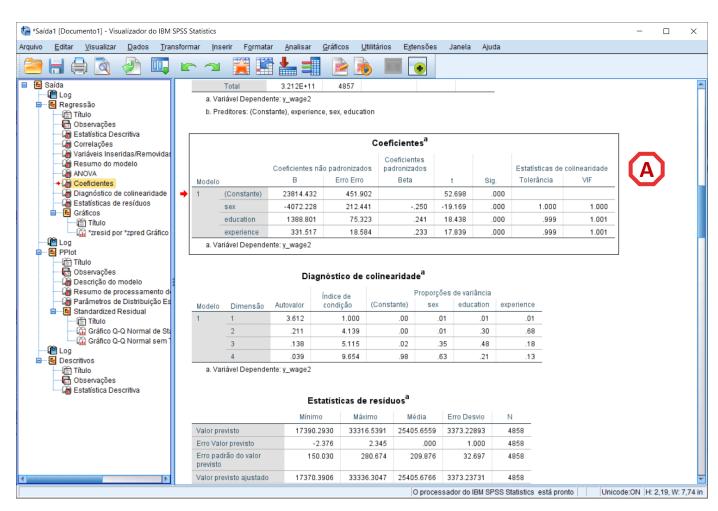
- Quando há fortes relações lineares entre os preditores numa regressão, a precisão dos coeficientes de regressão diminui em comparação com o que teria sido se os preditores não se correlacionassem entre si
- Um valor de VIF > 3 sugere a existência de colinearidade no modelo
- Um valor de VIF > 10 sugere a existência de colinearidade séria

Deve repensar-se as variáveis a incluir no modelo



Ausência de Multicolinearidade

- Para testarmos este pressuposto, temos de olhar para a Tabela de Coeficientes - que o SPSS produz automaticamente.
- Interpretação
- VIF > 3 -> presença de colinearidade
- Neste caso, não se identifica a presença de colinearidade...
- Portanto, cumpre-se o pressuposto da ausência de Multicolinearidade





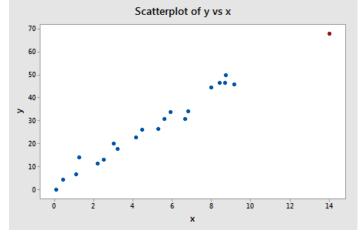
Validação do Modelo de Regressão Linear

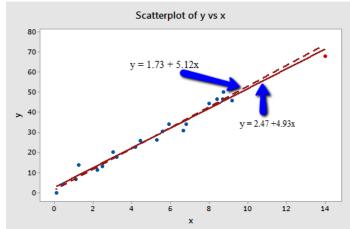
Avaliação do Pressuposto VII: Ausência de Observações Influentes

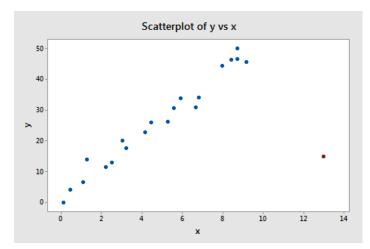


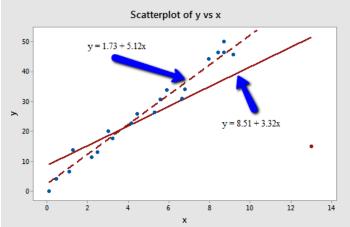
Ausência de Observações Influentes

- A existência de 'Outliers' (valores extremos) não é um problema em si.
- Mas torna-se um problema quando os Outlier têm influência sobre os resultados do modelo
- Nos painéis de baixo, o Outlier é uma 'Observação Influente'

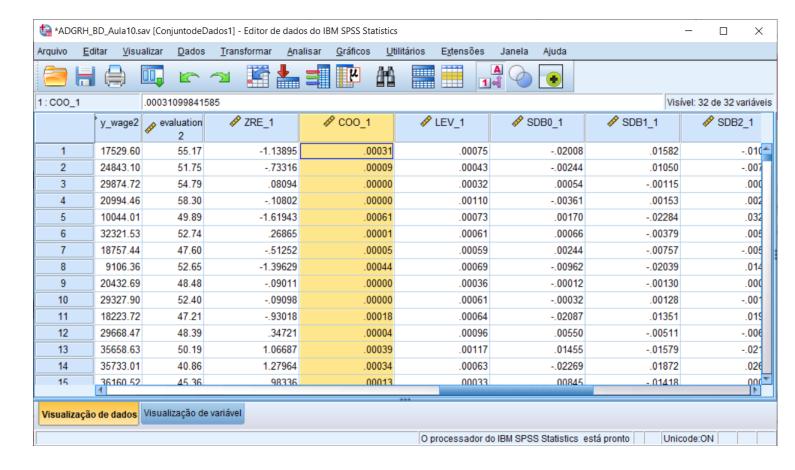








 Para testarmos a presença de observações influentes vamos usar a variável com os 'Distância de Cook' (COO_1) que acabamos de criar.



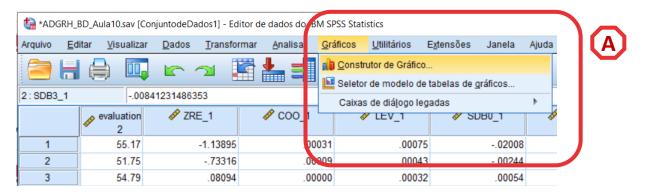
 Selecionar 'Gráficos' / 'Construtor de Gráfico'



Selecionar 'DispersãoPontos'



Selecionar 'Dispersão (Simples)'





 Selecionar 'Gráficos' / 'Construtor de Gráfico' A

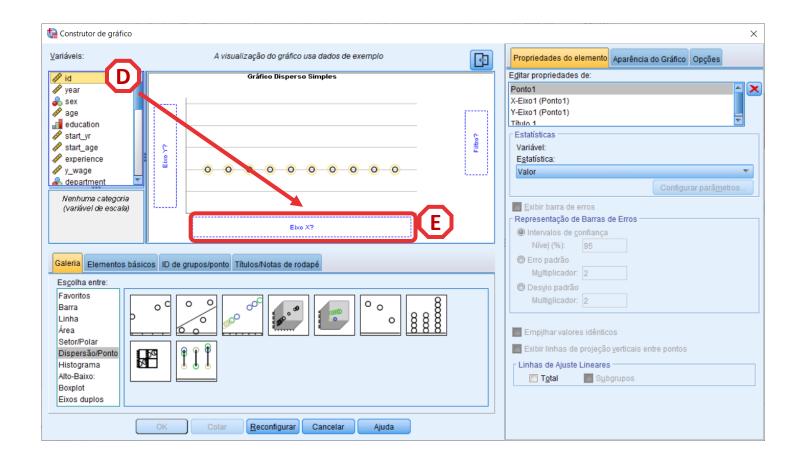
- Selecionar 'DispersãoPontos'
- B
- Selecionar 'Dispersão (Simples)'
- C

Selecionar Variável 'id'

(

Colocar no eixo 'x'

E



Selecionar 'Gráficos' / 'Construtor de Gráfico'

- Selecionar 'DispersãoPontos'
- (B)

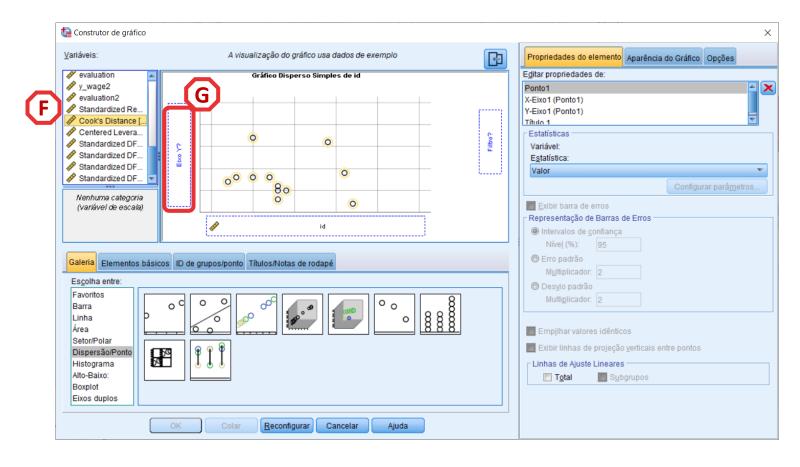
(c)

- Selecionar 'Dispersão (Simples)'
- Selecionar Variável 'id'
- Colocar no eixo 'x'

Selecionar Variável 'Cooks Distance'







- Selecionar 'ID de grupos/ponto'
- H
- Selecionar 'Rótulo da ID do Ponto'

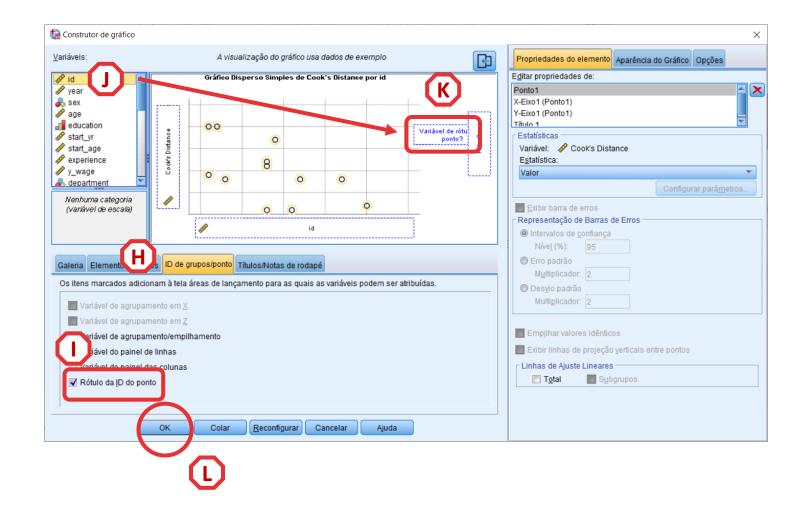
①

Selecionar Variável 'id'

- Colocar na caixa 'Variável do rótulo do ponto'
- K

Selecionar 'OK'

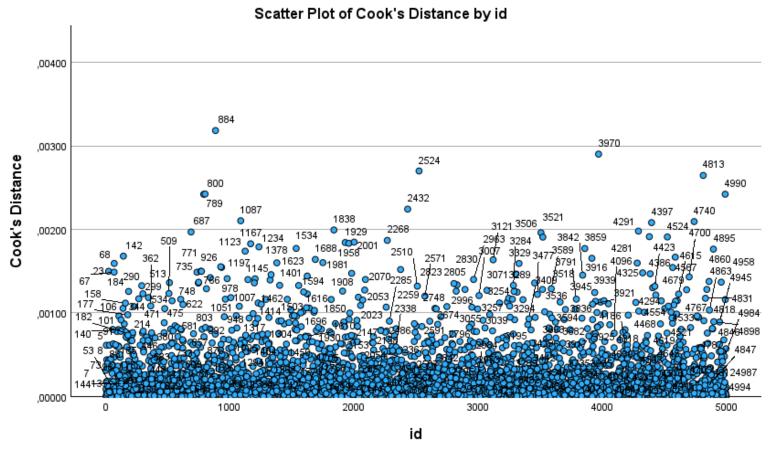




- Reparem que o gráfico permite identificar o ID dos outliers
- Interpretação

CD > 4/n -> Caso Influente

- Neste caso, o valor de corte é
 0.008 (4 / 5000)
- Neste caso não há observações acima do valor de corte.
- Cumpre-se o pressuposto da ausência de observações influentes





Reportar os resultados do estudo dos pressupostos [adaptar]

Foram realizados uma série de análises para averiguar a adequabilidade do modelo de regressão linear para o estudo destas relações, a maior parte dos pressupostos assumidos com a aplicação deste técnica foram validados. Em primeiro lugar, analisou-se graficamente a linearidade das relações entre as variáveis independentes (experiência e desempenho) com a variável dependente, tendo sido possível observar relações tendencialmente lineares, especialmente entre desempenho e rendimento (Figura x). Apurou-se também a ausência de multicolinearidade entre as variáveis independentes e de controlo com recurso às medidas VIF (<3).

Posteriormente, analisou-se a distribuição dos resíduos do modelo, observando-se uma distribuição normal, com um média em torno do valor zero, e com uma variância relativamente constante ao longo dos valores previstos do modelo (Figura x). Apurou-se ainda a existência de observações influentes com a Distância de Cook, admitindo os valores acima de 0,008 (4/N) como indicadores de observações influentes, não tendo sido detetados casos potencialmente problemáticos à estimação do modelo. Os resultados do teste Durbin-Watson, contudo sugerem autocorreção significativa nos resíduos (D-W<2), sugerindo reserva na interpretação dos resultados do modelo.



Exercício em autonomia

Repetir o exercício com o modelo da aula passada, incluído a experience no modelo.