

Regression analysis (PC-exercise)

Sebastian Jentschke





Agenda

- Regression analysis
 - **Checking requirements**
 - Correlation to Linear regression (one predictor)
 - Linear regression (multiple predictors)
 - Methods for adding predictors
 - Assessing the quality of your model
 - **Assignment**
- Logistic regression (binary)
 - Introduction
 - **Assignment**





Exam Anxiety.sav

IV: Exam; DV: Revise, Anxiety

**Check for normality, linearity,
and multi-collinearity**

Remove outliers (10 – 15 mins)



From correlation to regression

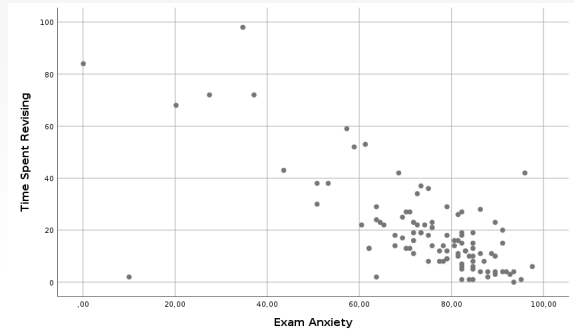
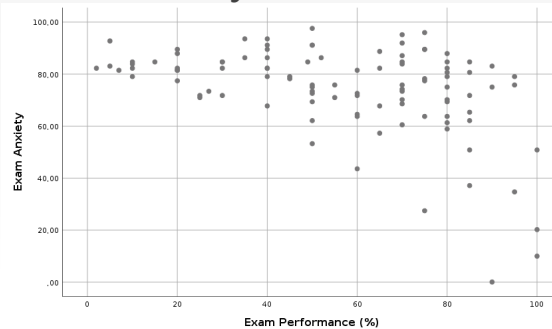
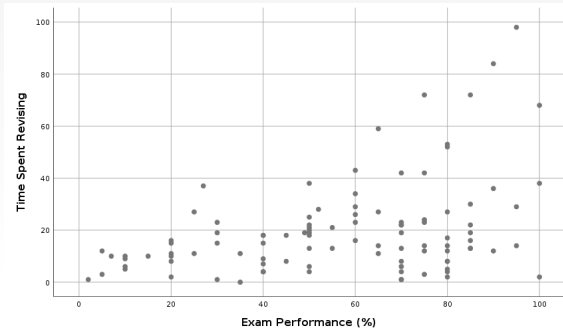
- **check normality:**
 - Analyze → Descriptive statistics → Explore...
 - tick «Normality plots» under «Plots»; assess «Tests of Normality»
- **check multicollinearity:**
 - Analyze → Correlate → Bivariate; enter the three variables «Revise», «Exam», and «Anxiety»
 - all vars. correlate substantially, esp. Revise and Anxiety
 - possibly only include one in the model; check in any case (using hierarchical regression whether adding both improves the prediction)





From correlation to regression

- **check linearity:**
 - Graphs → Legacy dialogs → Scatter / dot..., select «Simple Scatter», click define; in the window that opens, click «Revise» into «Y axis» and «Exam» into «X axis»; copy the syntax either using «Paste», duplicate it twice and create all three possible variable combinations
 - check deviations from linearity in the scatter plots





From correlation to regression

Option 1 – Box-Whisker-plots → manually de-select extreme outliers (stars in the Box-Whisker-plots):

- create a new variable (e.g., selSbj) using
→ Transform → Create variable...
use selSbj as «Target variable» and «1» under
«Numeric expression»
- manually set participants that are outliers (stars) in the
Box-Whisker-plots to «0»
- use → Data → Select cases and choose «selSbj»
under «Select filter variable»





From correlation to regression

Option 1 – Box-Whisker-plots (contd.)

- repeat → Analyze → Descriptive statistics → Explore... and → Analyze → Correlate → Bivariate...
check the results:
 - how do the «Tests of Normality» change?
 - how do the correlations change?
- create another selection variable and further remove the less extreme outliers (circles) and re-check the analyses
- **BEFORE** you deselect participants on the basis of that they might be outliers, ask yourself whether the outliers could be genuine (it might be uncommon but valid to be 2.04 m tall)





From correlation to regression

Option 2 – z-scores (use the table with descriptive stats.):

- create a new variable (e.g., selSbjZ)
- re-arrange the table using the Pivoting trays
- calculate $M \pm 3.3 * SD$ (manually or in Excel)

Option 3 – multivariate outliers (Mahalanobis)

- create a new variable (e.g., selSbjM)
- Analyze → Regression → Linear; «Save»-button, select Distances – Mahalanobis; crit. $\chi^2 = 16.266$





Questions?
Comments?



From correlation to regression

→ Analyze → Correlate → Bivariate...

→ Analyze → Regression → Linear (Revise → Exam)

Correlations

		Time Spent Revising	Exam Performance (%)	Exam Anxiety
Time Spent Revising	Pearson Correlation	1	.322**	-.620**
	Sig. (2-tailed)		.001	.000
	N	97	97	97
Exam Performance (%)	Pearson Correlation	.322**	1	-.317**
	Sig. (2-tailed)	.001		.002
	N	97	97	97
Exam Anxiety	Pearson Correlation	-.620**	-.317**	1
	Sig. (2-tailed)	.000	.002	
	N	97	97	97

** . Correlation is significant at the 0.01 level (2-tailed).

Model Summary^b

Model	R	R Square	Adjusted R Square	Std. Error of the Estimate
1	.322 ^a	.104	.094	11.716

a. Predictors: (Constant), Exam Performance (%)

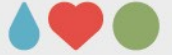
b. Dependent Variable: Time Spent Revising

Coefficients^a

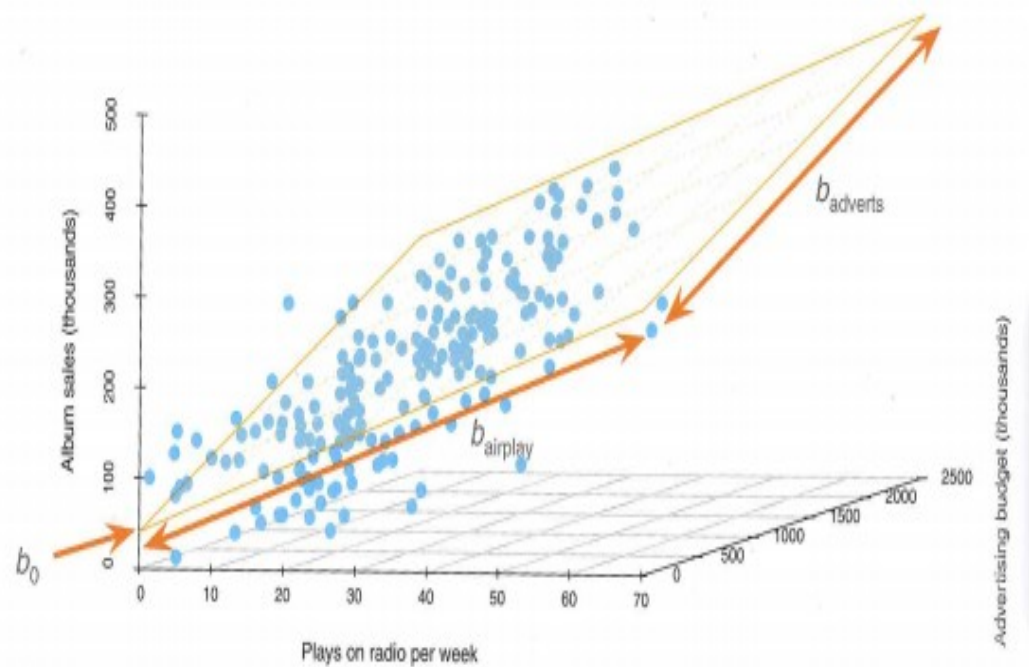
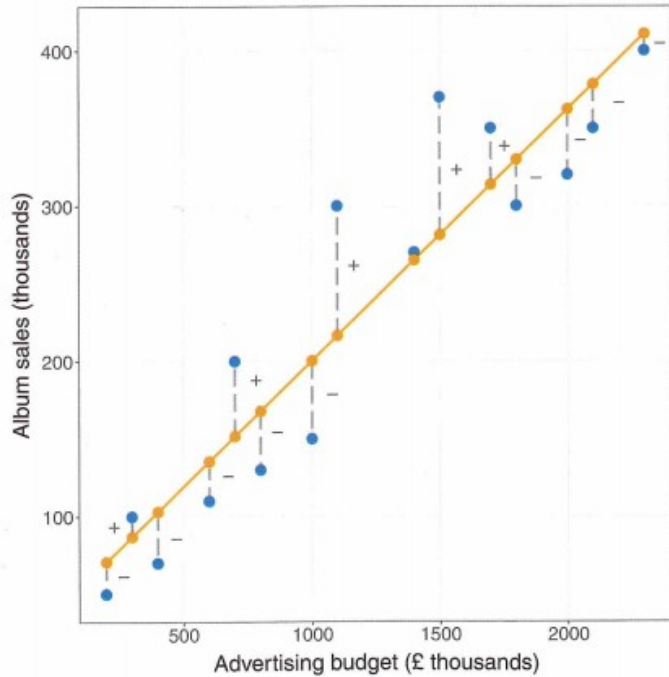
Model		Unstandardized Coefficients		Standardized Coefficients	t	Sig.
		B	Std. Error	Beta		
1	(Constant)	8.414	2.850		2.953	.004
	Exam Performance (%)	.158	.048	.322	3.316	.001

a. Dependent Variable: Time Spent Revising





Regression: From uni- to multivar.





Questions?
Comments?



Regression: From uni- to multivar.

*Exam Anxiety.sav [DataSet3] - IBM SPSS Statistics Data Editor

File Edit View Data Transform Analyze Graphs Utilities Extensions Window Help

87 : selSbjM 0

	selSbjE	selSbjM
83	0	0
84	1	1
85	1	1
86	1	1
87	1	0
88	1	1
89	1	1
90	1	1
91	1	1
92	1	1
93	1	1
94	1	1
95	1	1
96	1	1
97	1	1
98	1	1
99	1	1
100	1	1
101	1	1
102	1	1
103	1	1
104		
105		
106		
107		
108		
109		
110		
111		
112		
113		
114		
115		
116		
117		
118		
119		
120		
121		

Visible: 7 of 7 Variables

Sex	var	var	var	var	var	var	var	var	var	var	var	var	var	var
2														
1														
2														
1														
2														
1														
2														

Regression

- Automatic Linear Modeling...
- Linear...
- Curve Estimation...
- Partial Least Squares...
- Binary Logistic...
- Multinomial Logistic...
- Ordinal...
- Probit...
- Nonlinear...
- Weight Estimation...
- 2-Stage Least Squares...
- Optimal Scaling (CATREG)...

Data View Variable View

Linear...

IBM SPSS Statistics Processor is ready | Unicode:ON | Filter On





Regression: From uni- to multivar.

Linear Regression

Dependent:

Block 1 of 1

Independent(s):

Method: Enter

Selection Variable:

Case Labels:

WLS Weight:

OK Paste Reset Cancel Help

Statistics... Plots... Save... Options... Style... Bootstrap...

Select participants (remove ext...)
Select participants (remove ext...)
Participant Code [Code]
Time Spent Revising [Revise]
Exam Performance (%) [Exam]
Exam Anxiety [Anxiety]
Biological sex of participant [Sex]

Previous Next

Linear Regression

Dependent: Exam Performance (%) [Exam]

Block 1 of 1

Independent(s):
Time Spent Revising [Revise]
Exam Anxiety [Anxiety]

Method: Enter

Selection Variable:

Case Labels:

WLS Weight:

OK Paste Reset Cancel Help

Statistics... Plots... Save... Options... Style... Bootstrap...

Select participants (remove ext...)
Select participants (remove ext...)
Participant Code [Code]
Time Spent Revising [Revise]
Exam Anxiety [Anxiety]
Biological sex of participant [Sex]

Previous Next





Regression: From uni- to multivar.

Model Summary^b

Model	R	R Square	Adjusted R Square	Std. Error of the Estimate
1	.322 ^a	.104	.094	23,928

- a. Predictors: (Constant), Time Spent Revising
 b. Dependent Variable: Exam Performance (%)

ANOVA^a

Model		Sum of Squares	df	Mean Square	F	Sig.
1	Regression	6294,045	1	6294,045	10,993	.001 ^b
	Residual	54393,997	95	572,568		
	Total	60688,041	96			

- a. Dependent Variable: Exam Performance (%)
 b. Predictors: (Constant), Time Spent Revising

Coefficients^a

Model		Unstandardized Coefficients		Standardize	t	Sig.
		B	Std. Error	d Coefficients Beta		
1	(Constant)	43,272	4,157		10,410	.000
	Time Spent Revising	.658	.198	.322	3,316	.001

- a. Dependent Variable: Exam Performance (%)

Model Summary

Model	R	R Square	Adjusted R Square	Std. Error of the Estimate
1	.355 ^a	.126	.107	23,755

- a. Predictors: (Constant), Exam Anxiety, Time Spent Revising

ANOVA^a

Model		Sum of Squares	df	Mean Square	F	Sig.
1	Regression	7641,955	2	3820,977	6,771	.002 ^b
	Residual	53046,086	94	564,320		
	Total	60688,041	96			

- a. Dependent Variable: Exam Performance (%)
 b. Predictors: (Constant), Exam Anxiety, Time Spent Revising

Coefficients^a

Model		Unstandardized Coefficients		Standardize	t	Sig.
		B	Std. Error	d Coefficients Beta		
1	(Constant)	81,471	25,059		3,251	.002
	Time Spent Revising	.417	.251	.204	1,662	.100
	Exam Anxiety	-.440	.284	-.190	-1,545	.126

- a. Dependent Variable: Exam Performance (%)



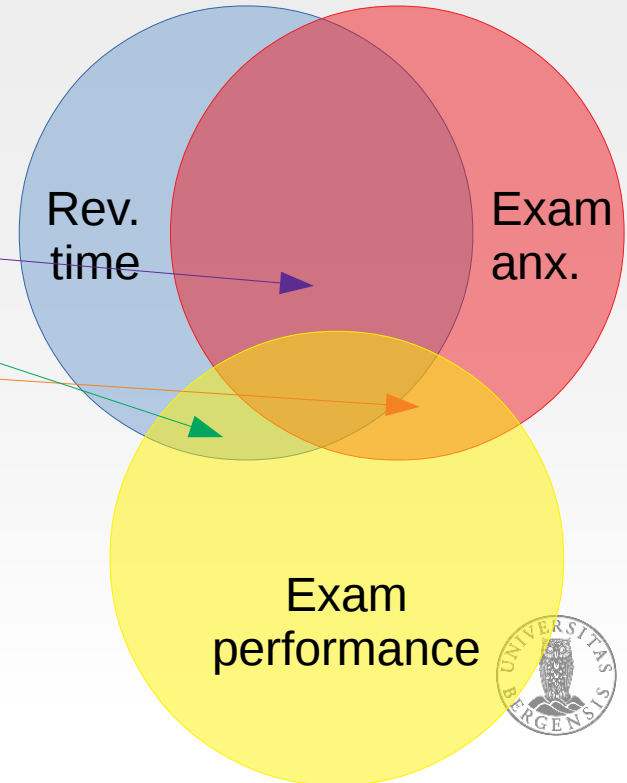


Regression: From uni- to multivar.

Correlations

		Time Spent Revising	Exam Performance (%)	Exam Anxiety
Time Spent Revising	Pearson Correlation	1	.322**	-.620**
	Sig. (2-tailed)		,001	,000
	N	97	97	97
Exam Performance (%)	Pearson Correlation	.322**	1	-.317**
	Sig. (2-tailed)	,001		,002
	N	97	97	97
Exam Anxiety	Pearson Correlation	-.620**	-.317**	1
	Sig. (2-tailed)	,000	,002	
	N	97	97	97

** , Correlation is significant at the 0.01 level (2-tailed).





Regression: Entering predictors

Linear Regression

Dependent: Exam Performance (%) [Exam]

Block 1 of 1

Independent(s): Time Spent Revising [Revise]

Method: Enter

Selection Variable: Rule...

Case Labels:

WLS Weight:

OK Paste Reset Cancel Help

Linear Regression

Dependent: Exam Performance (%) [Exam]

Block 1 of 1 **doesn't count correctly**

Independent(s): Exam Anxiety [Anxiety]

Method: Enter

Selection Variable: Rule...

Case Labels:

WLS Weight:

OK Paste Reset Cancel Help

Linear Regression: Statistics

Regression Coefficient... Model fit

Estimates **R squared change**

Confidence intervals Descriptives

Level(%): 95 Part and partial correlations

Covariance matrix Collinearity diagnostics

Residuals

Durbin-Watson

Casewise diagnostics

Outliers outside: 3 standard deviations

All cases

Continue Cancel Help





Regression: Entering predictors

Model Summary

Model	R	R Square	Adjusted R Square	Std. Error of the Estimate	Change Statistics			Sig. F Change	
					R Square Change	F Change	df1		df2
1	.237 ^a	.056	.045	24,616	.056	5,103	1	86	.026
2	.290 ^b	.084	.063	24,388	.028	2,609	1	85	.110

a. Predictors: (Constant), Time Spent Revising

b. Predictors: (Constant), Time Spent Revising, Exam Anxiety

ANOVA^a

Model		Sum of Squares	df	Mean Square	F	Sig.
1	Regression	3092,020	1	3092,020	5,103	.026 ^b
	Residual	52109,878	86	605,929		
	Total	55201,898	87			
2	Regression	4643,996	2	2321,998	3,904	.024 ^c
	Residual	50557,902	85	594,799		
	Total	55201,898	87			

a. Dependent Variable: Exam Performance (%)

b. Predictors: (Constant), Time Spent Revising

c. Predictors: (Constant), Time Spent Revising, Exam Anxiety

Coefficients^a

Model		Unstandardized Coefficients		Standardized Coefficients Beta	t	Sig.
		B	Std. Error			
1	(Constant)	42,735	5,044		8,472	.000
	Time Spent Revising	.684	.303	.237	2,259	.026
2	(Constant)	92,430	31,168		2,966	.004
	Time Spent Revising	.393	.350	.136	1,122	.265
	Exam Anxiety	-.574	.355	-.196	-1,615	.110

a. Dependent Variable: Exam Performance (%)





Regression: Entering predictors

Methods for entering:

- «Enter»: Enter one new variable
- «Stepwise»: Enter multiple new variables (one step at a time an in order of explained variance) according to F-probability
- «Remove»: Remove one var.
- «Backward»: Remove one var. according to F-probability
- «Forward»: enter one variable according to F-probability

The image shows two overlapping SPSS dialog boxes. The top one is the 'Linear Regression' dialog, and the bottom one is the 'Linear Regression: Options' dialog.

Linear Regression Dialog:

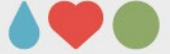
- Dependent:** Exam Performance (%) [Exam]
- Independent(s):** Exam Anxiety [Anxiety]
- Method:** Enter (highlighted in a red box)
- Selection Variable:** (empty)
- Case Labels:** (empty)
- WLS Weight:** (empty)
- Buttons:** Statistics..., Plots..., Save..., Options... (highlighted in a red box), Style..., Bootstrap...

Linear Regression: Options Dialog:

- Stepping Method Criteria:**
 - Use probability of F
 - Entry: .05
 - Removal: .10
 - Use F value
 - Entry: 3,84
 - Removal: 2,71
- Include constant in equation
- Missing Values:**
 - Exclude cases listwise
 - Exclude cases pairwise
 - Replace with mean
- Buttons:** Continue, Cancel, Help

An arrow points from the 'Options...' button in the main dialog to the 'Linear Regression: Options' dialog.





Regression: Entering predictors

Model Summary

Model	R	R Square	Adjusted R Square	Std. Error of the Estimate	R Square Change	Change Statistics			Sig. F Change
						F Change	df1	df2	
1	.266 ^a	.071	.060	24,425	.071	6,530	1	86	.012

a. Predictors: (Constant), Exam Anxiety

ANOVA^a

Model		Sum of Squares	df	Mean Square	F	Sig.
1	Regression	3895,715	1	3895,715	6,530	.012 ^b
	Residual	51306,183	86	596,584		
	Total	55201,898	87			

a. Dependent Variable: Exam Performance (%)

b. Predictors: (Constant), Exam Anxiety

Coefficients^a

Model		Unstandardized Coefficients		Standardized Coefficients	t	Sig.
		B	Std. Error			
1	(Constant)	114,314	24,343		4,696	.000
	Exam Anxiety	-,779	,305	-,266	-2,555	.012

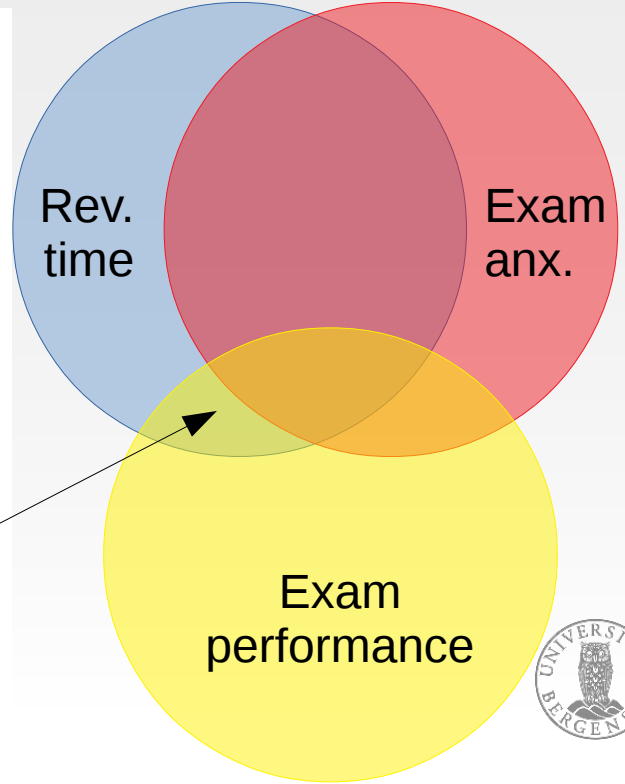
a. Dependent Variable: Exam Performance (%)

Excluded Variables^a

Model		Beta In	t	Sig.	Partial Correlation	Collinearity Statistics
1	Time Spent Revising	.136 ^b	1,122	.265	.121	.734

a. Dependent Variable: Exam Performance (%)

b. Predictors in the Model: (Constant), Exam Anxiety





Questions?
Comments?

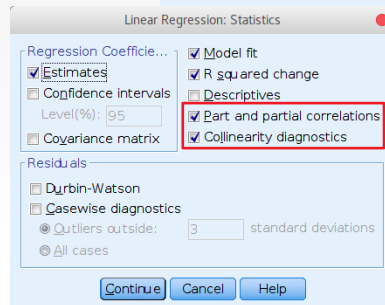
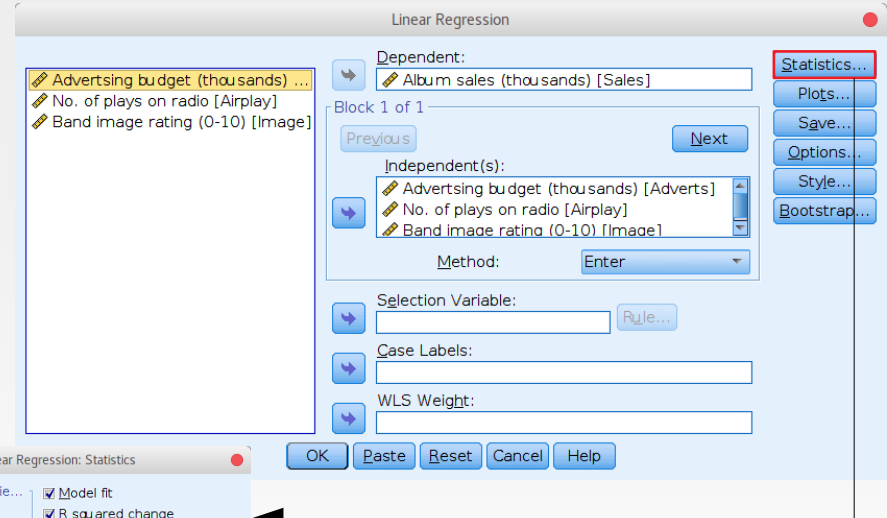


**Use Album Sales.sav
Predict Sales from
Adverts, Airplay, and
Image
(~ 5 – 10 mins)**



Regression: Diagnostics

- use *Album Sales.sav*
- → Analyze
- → Regression
- → Linear...
- click «Statistics» and tick those





Regression: Diagnostics

- collinearity describes a linear association between explanatory variables (i.e. the degree to which one explanatory variable can be predicted by a combination of one or more other explanatory variables)
- tolerance: $1 - R_j^2$ (R_j^2 – what degree of variance of variable j is explained by the other predictor variables)
- variance inflation factor (VIF): $1 / \text{tolerance}$
- (a) $VIF < 5$ and tolerance > 0.2 ; (b) the average of the VIF of all variables should be close to 1





Regression: Diagnostics

- all tolerances are > 0.2 , all VIF < 5
- the average VIF is $(1.015 + 1.043 + 1.038) / 3 = 1.032$ which is close to 1

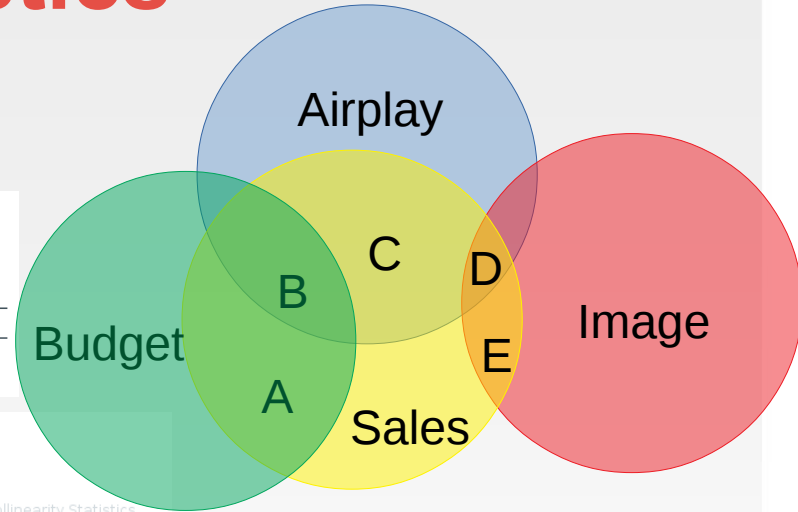
		Coefficients ^a									
		Unstandardized Coefficients		Standardized Coefficients			Correlations			Collinearity Statistics	
Model		B	Std. Error	Beta	t	Sig.	Zero-order	Partial	Part	Tolerance	VIF
1	(Constant)	-26,613	17,350		-1,534	,127					
	Advertising budget (thousands)	,085	,007	,511	12,261	,000	,578	,659	,507	,986	1,015
	No. of plays on radio	3,367	,278	,512	12,123	,000	,599	,655	,501	,959	1,043
	Band image rating (0-10)	11,086	2,438	,192	4,548	,000	,326	,309	,188	,963	1,038

a. Dependent Variable: Album sales (thousands)



Regression: Diagnostics

- $R^2 = A + B + C + D + E$



Model Summary^b

Model	R	R Square	Adjusted R Square	Std. Error of the Estimate	R Square Change	Change Statistics			Sig. F Change
						F Change	df1	df2	
1	,815 ^a	,665	,660	47,087	,665	129,498	3	196	,000

a. Predictors: (Constant), Band image rating (0-10), Advertising budget (thousands), No. of plays on radio
 b. Dependent Variable: Album sales (thousands)

Coefficients^a

Model		Unstandardized Coefficients		Standardized Coefficients Beta	t	Sig.	Correlations			Collinearity Statistics		
		B	Std. Error				Zero-order	Partial	Part	Tolerance	VIF	
1	(Constant)	-26,613	17,350		-1,534	,127						
	Advertising budget (thousands)	,085	,007	,511	12,261	,000	,578	,659	,507	,986	1,015	
	No. of plays on radio	3,367	,278	,512	12,123	,000	,599	,655	,501	,659	1,043	
	Band image rating (0-10)	11,086	2,438	,192	4,548	,000	,326	,309	,188	,663	1,038	

a. Dependent Variable: Album sales (thousands)





Regression: Diagnostics

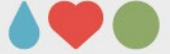
- within the Variance Proportions, for each dimension should only one variable have high loadings

Collinearity Diagnostics^a

Model	Dimension	Eigenvalue	Condition Index	(Constant)	Variance Proportions		
					Advertising budget (thousands)	No. of plays on radio	Band image rating (0-10)
1	1	3,562	1,000	,00	,02	,01	,00
	2	,308	3,401	,01	,96	,05	,01
	3	,109	5,704	,05	,02	,93	,07
	4	,020	13,219	,94	,00	,00	,92

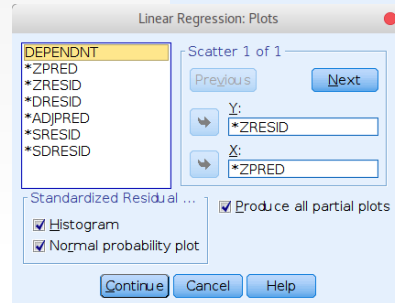
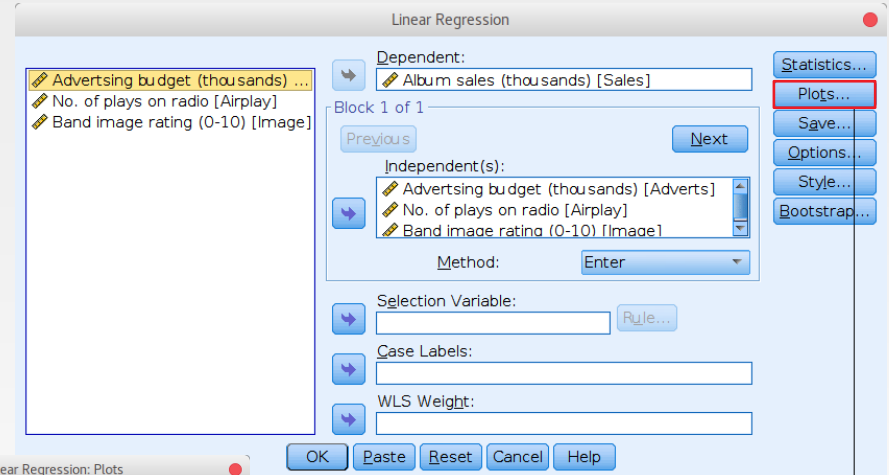
a. Dependent Variable: Album sales (thousands)





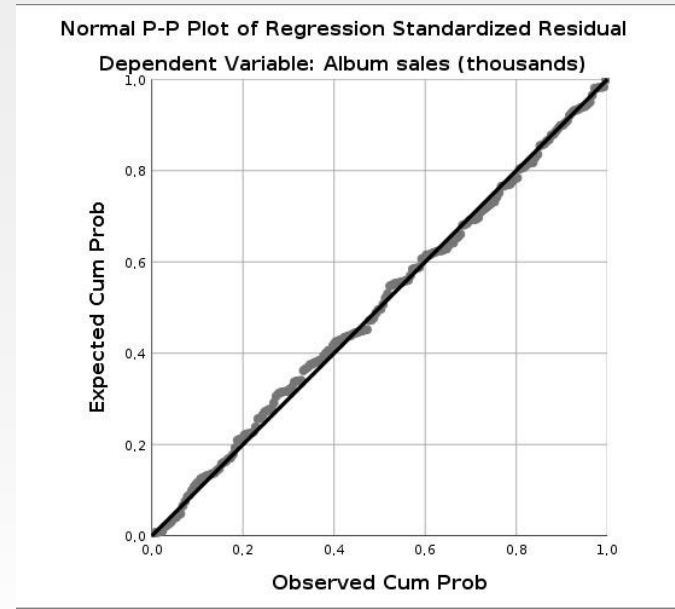
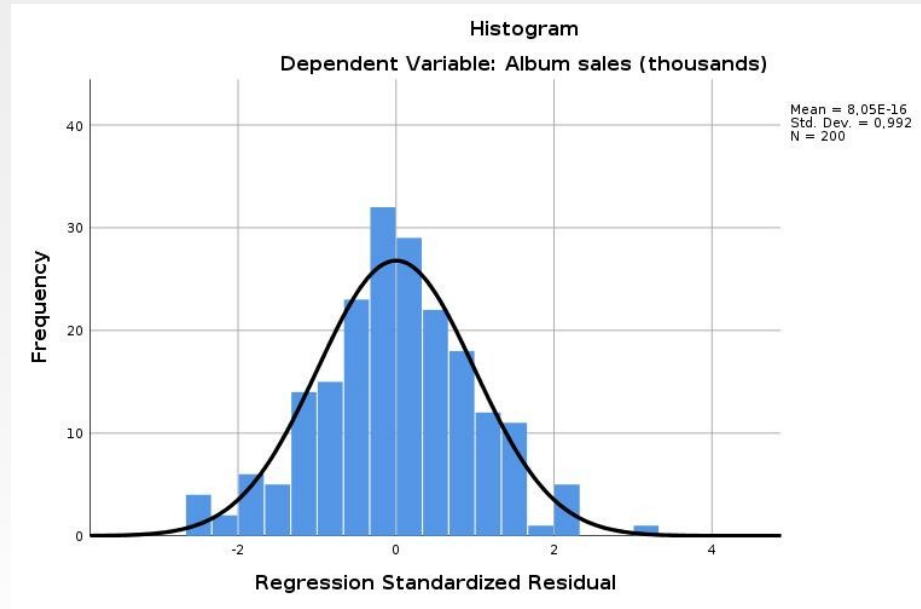
Regression: Diagnostics

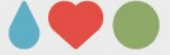
- click «Plots»
 - tick all options
 - click ZPRED to X
click ZRESID to Y





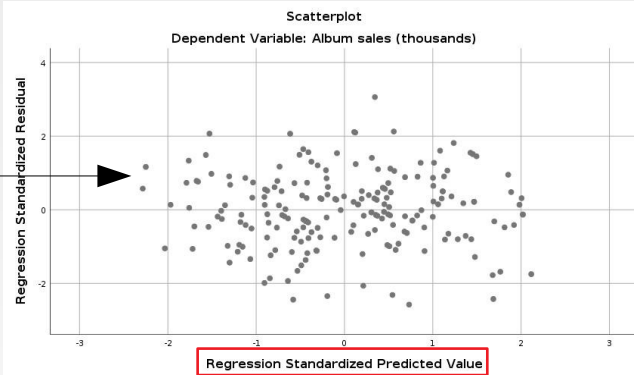
Regression: Diagnostics



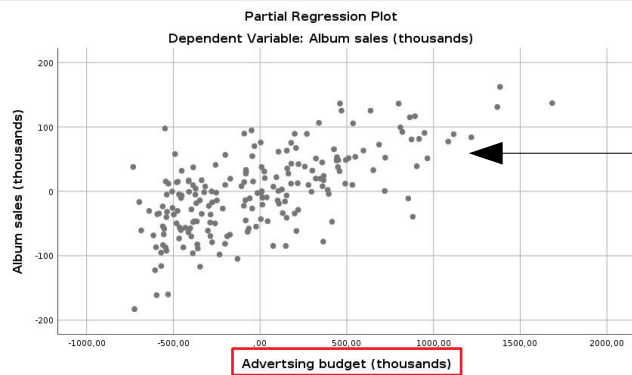


Regression: Diagnostics

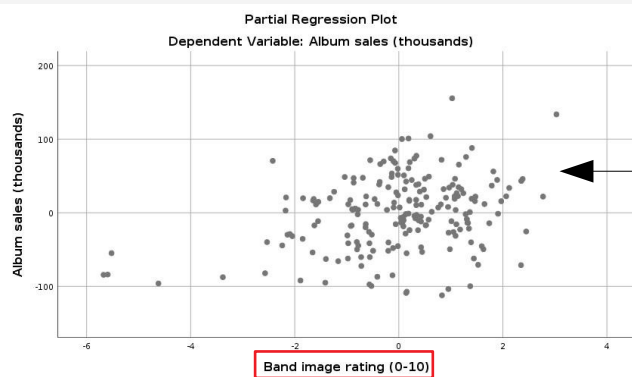
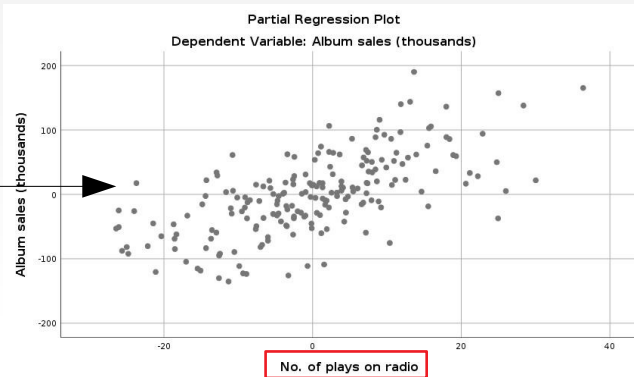
should have no shape ($r \approx 0$)



should show a linear trend ($r \neq 0$)



should show a linear trend ($r \neq 0$)





Questions?
Comments?



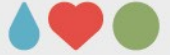
Logistic Regression

- predicting categorical outcomes from categorical and continuous predictors (binary / multinomial)
- log-transform the result of the GLM → probability of the category (e.g., successful treatm.) to occur

$$P(\gamma) = \frac{1}{1 + e^{-(b_0 + b_1 X_{1i} + b_2 X_{2i} + \dots + b_n X_{ni})}}$$

→ maximum-likelihood-estimation





Logistic Regression

assessing the quality of the model:

$$\text{log-likelihood} = \sum_{i=1}^N \left[Y_i \ln(P(Y_i)) + (1 - Y_i) \ln(1 - P(Y_i)) \right]$$

deviance = $-2 \times \text{log-likelihood}$ ($-2LL$); χ^2 -distributed.

$$\chi^2 = (-2LL_{\text{baseline}}) - (-2LL_{\text{model}}) = 2LL_{\text{model}} - 2LL_{\text{baseline}}$$

$$df = k_{\text{model}} - k_{\text{baseline}} \quad (\text{predictors} + 1 \text{ [intercept]})$$



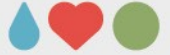


Logistic Regression

assessing the quality of the model:

- $R = \sqrt{(\chi^2 - 2df) / -2LL_{\text{baseline}}}$
- $R^2_{\text{HL}} = (2LL_{\text{model}} - 2LL_{\text{baseline}}) / -2LL_{\text{baseline}}$ (Hosmer & Lemeshow)
- $R^2_{\text{CS}} = 1 - \exp((2LL_{\text{baseline}} - 2LL_{\text{model}}) / n)$ (Cox & Snell)
- Wald (assessing signif. of predictors): $z = b / SE_b$
- odds-ratio = $P(\text{event}) / P(\text{non-event})$
 = odds after unit chg. in pred. / orig. odds
 (if > 1 : if the pred. increases the prob. of outcome incr.
 if < 1 : if the pred. increases the prob. of outcome decr.)





Logistic Regression

assumptions:

- linear relationship between any (continuous) predictor and the logit of the outcome
(can be tested by testing the significance of the interaction of a predictor with its log-transform.)
- independence of errors
- threats to convergence: incomplete information
(not all possible combinations of variables available)

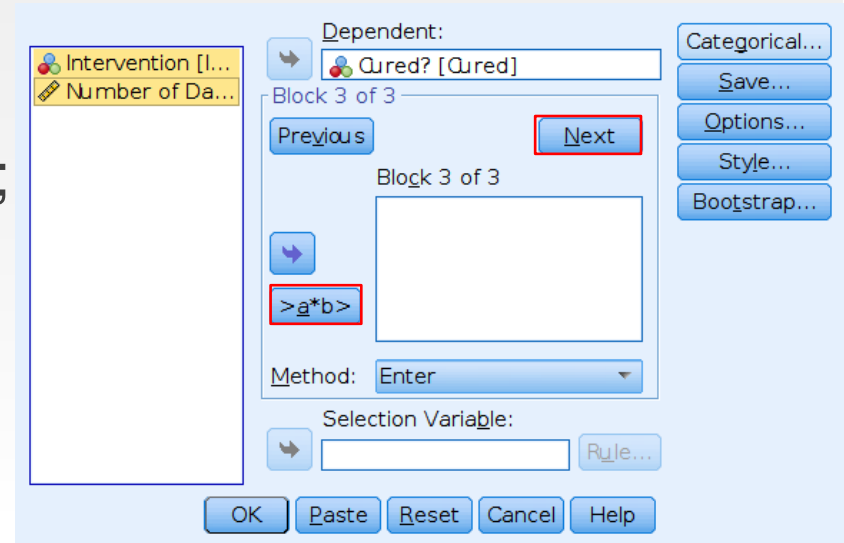
complete separation (outcome can be perfectly predicted by one or a combination of variables)





Logistic Regression

- using Eel.sav:
Analyze → Regression → Binary Logistic
DV: Cured;
IV: Intervention [Block 1;
press Next] Duration [Block 2];
Intervention × Duration
[Block 3; select both by holding
[Ctrl] + Click, press button >a*b<]
use «Enter» as method





Logistic Regression

Predicted Values

Probabilities

Group membership

Influence

Cook's

Leverage values

DfBeta(s)

Residuals

Unstandardized

Logit

Studentized

Standardized

Deviance

Export model information to XML file

Include the covariance matrix

Statistics and Plots

Classification plots

Correlations of estimates

Hosmer-Lemeshow goodness-of-fit

Iteration history

Casewise listing of residuals

CI for exp(B): 95 %

Outliers outside 2 std. dev.

All cases

Display

At each step At last step

Probability for Stepwise

Entry: 0,05 Removal: 0,10

Classification cutoff: 0,5

Maximum Iterations: 20

Conserve memory for complex analyses or large datasets

Include constant in model

there should be max. 5% residuals > 2; max. %1 > 2.5; > 3 is certainly an outlier
 Cook's distance > 1: case influences the model; look out for DfBeta > 1
 Leverage should be with 2-3 times predictors / N ($2 / 113 = 0.018 \rightarrow$ check > 0.036)





Logistic Regression

Block 2: Method = Enter

Omnibus Tests of Model Coefficients

		Chi-square	df	Sig.
Step 1	Step	,002	1	,964
	Block	,002	1	,964
	Model	9,928	2	,007

Model Summary

Step	-2 Log likelihood	Cox & Snell R Square	Nagelkerke R Square
1	144,156 ^a	,084	,113

a. Estimation terminated at iteration number 3 because parameter estimates changed by less than ,001.

Classification Table^a

	Observed	Predicted		Percentage Correct	
		Not Cured	Cured		
Step 1	Cured?	Not Cured	32	16	66,7
		Cured	24	41	63,1
	Overall Percentage				64,6

a. The cut value is ,500

Variables in the Equation

Step 1 ^a		B	S.E.	Wald	df	Sig.	Exp(B)	95% C.I. for EXP(B)	
								Lower	Upper
	Intervention(1)	1,234	,415	8,854	1	,003	3,433	1,523	7,737
	Number of Days with Problem before Treatment	-,008	,176	,002	1	,964	,992	,703	1,401
	Constant	-,235	1,221	,037	1	,848	,791		

a. Variable(s) entered on step 1: Number of Days with Problem before Treatment.





use Penalty.sav

DV: Scored

IVs: PSWQ, Anxious,

Previous

(10 - 15 min)



It's your turn now!



UNIVERSITY OF BERGEN

