

Bi8352: Metody antropologie II  
 jaro 2019  
 Mgr. Mikoš Jurda, Ph.D.

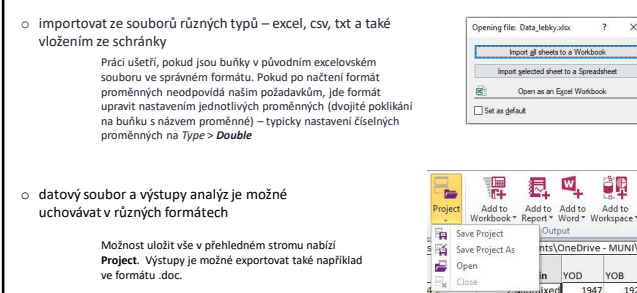
**MUNI  
SCI**

## Základy Statistics

1

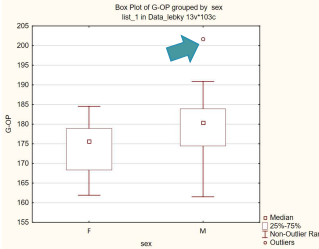
### Základní východiska

- importovat ze souborů různých typů – excel, csv, txt a také vložením ze schránky
  - Práci ušetří, pokud jsou buňky v původním excelovském souboru ve správném formátu. Pokud po načtení formát proměnných neodpovídá našim požadavkům, jde formát upravit nastavením jednotlivých proměnných (dvojitě poklikání na buňku s názvem proměnné) – typicky nastavení číselných proměnných na *Type > Double*
- datový soubor a výstupy analýz je možné uchovávat v různých formátech
  - Možnost uložit vše v přehledném stromu nabízí **Project**. Výstupy je možné exportovat také například ve formátu .doc.



2

### Popisná statistika – vizuální hodnocení – krabicový graf

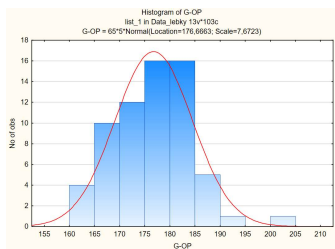


Pokud nezádáte grupovací proměnnou, zobrazí se graf pro celý soubor, pokud ano, pak odděleně pro definované

Krabicový graf pro dvě skupiny – m a f – dobrý pro vyhledávání extrémních případů (například chyb v datech) – při podržení myši nad odlehlou hodnotou se zobrazí její ID

3

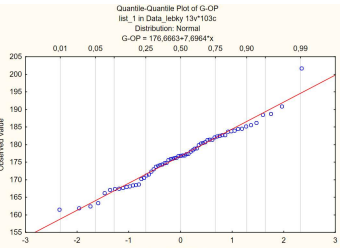
### Popisná statistika – vizuální hodnocení – histogram



Umožňuje posoudit rozložení hodnot a srovnat je s předpokládaným rozložením (linie). Nastavuje se jako *Fit type* dialogové okně.

4

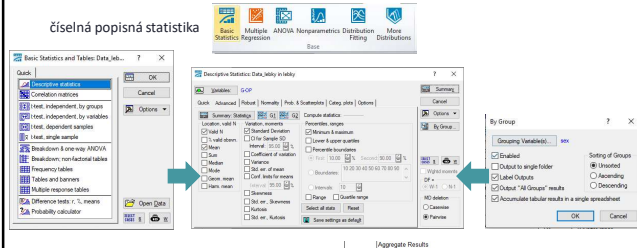
### Popisná statistika – vizuální hodnocení – QQ graf



Alternativní způsob porovnání pozorovaných hodnot s normálním rozložením (pozorovaný kvantil vs. teoretický kvantil).

5

### Popisná statistika – číselná popisná statistika



Variable	Valid N	Mean	Median	Maximum	Variance	Std Dev.
G-OP	82	176.6663	176.8000	181.5500	201.6700	58.86475

pro všechna data zároveň

by group – pro skupiny zvlášť

6

### Popisná statistika – číselná popisná statistika

souhrnné výsledky – histogram, krabicové grafy, zvolené parametry a P-P plot

v základní, přednastavené podobě

7

### Popisná statistika – normalita dat

**Grafické posouzení**  
Srovnání s normálním rozložením – viz předchozí grafy

**Testování**  
Statistické testy  
Statistics > Basic statistics > Descriptive statistics > Normality

Category	Count	Cumulative Count	Percent of Valid	Cumulative % of Valid
155.0000<<=160.0000	4	4	6.91380	6.91380
160.0000<<=165.0000	10	14	15.38462	21.53846
165.0000<<=170.0000	12	26	18.46154	40.00000
170.0000<<=175.0000	16	42	24.61538	64.61538
175.0000<<=180.0000	16	58	24.61538	89.23077
180.0000<<=185.0000	6	64	7.69231	96.92308
185.0000<<=190.0000	1	65	1.53846	98.46154
190.0000<<=195.0000	0	65	0.00000	98.46154
195.0000<<=200.0000	1	66	1.53846	100.00000
200.0000<<=205.0000	0	66	0.00000	100.00000
Missing	3	3	3.46154	3.46154
Total	69	69	100.00000	100.00000

8

### T-test

**Nepárový dvouvýběrový t-test**

**Předpoklady:**  
normální rozložení v rámci porovnávaných skupin  
- již představenými postupy

shoda rozptylů těchto skupin  
- testování je přímo součástí výsledků jako F-statistika

**pokud data nesplňují**

neparametrické alternativy  
v tomto případě  
*Mann-whitney U-test*

V případě různých rozptylů  
možno použít t-test se  
samostatnými odhady  
rozptylů

9

### T-test

Ověření předpokladů testu přímo v dialogovém okně

Shoda rozptylů – graficky  
Advanced > Box & Whiskers plot

Normalita rozložení v rámci skupin  
Advanced > Categorized normal plots

10

### T-test – výstupy

Samotné výstupy testu – lze provést hromadně pro všechny zároven

Variable	Mean F	Mean M	t-value	df	p	Valid N F	Valid N M	Std. Dev. F	Std. Dev. M	F-ratio	p
G-OP	174.22624	179.1825	-2.72022	63	0.008157	33	32	6.417068	8.133095	1.606344	0.187893
EU-EU	135.42591	141.4306	-3.37316	63	0.001275	33	32	5.635122	8.488520	2.258439	0.024809
BA-B	120.94911	132.4063	-2.51013	63	0.014650	33	32	4.981380	6.365717	1.646258	0.162511
ZYG-ZYG	120.5567	128.8028	-6.17618	63	0.000000	33	32	5.473546	5.284873	1.072676	0.846689
D-D	20.2794	21.4384	-1.89441	63	0.062759	33	32	2.369016	2.562275	1.169810	0.661057
RH-RS	33.09533	37.5522	-5.41930	63	0.000001	33	32	3.163414	3.467877	1.201199	0.600593
ZM-ZM	88.1052	92.2650	-2.89237	63	0.005243	33	32	6.449334	5.035607	1.640310	0.171598

skupinové průměry      samotná statistika      směrodatné odchylky      shoda rozptylů

11

### Diskriminační analýza

**Jaké použít proměnné**  
význam mají pouze ty, které mají nějakou souvislost s kategoriální proměnnou

redundantní proměnné snižují stabilitu modelu a mohou vést k nesmyslným výsledkům

**Hodnocení vztahu nezávislých proměnných a kategoriální proměnné**

- o t-test a ANOVA
- o korelační analýza a XY grafy
- o hlavní komponenty a faktorová analýza
- o diskriminační analýza
- o „expertní znalost proměnných“

12

### Diskriminační analýza

**Vztah ke kategoriální proměnné**

Samostatný t-test pro jednotlivé proměnné – pro dvě skupiny!!  
(Basic statistics > t-test, independent, by groups)

ANOVA (Basic statistics > Breakdowns & One-way ANOVA; Analysis of variance) – pro dvě a více skupin (pro dvě skupiny jsou výsledky obdobné jako t-test)

Variable	SS	df	MS	SSE	df	MS	F	p
KaDP	202.028	1	202.028	338.295	63	5.3699	7.45396	0.008153
EU-EU	585.362	1	585.362	3279.338	63	51.41806	11.37815	0.001275
BeH	294.416	1	294.416	2941.956	63	46.6993	6.30074	0.014665
ZYG-ZYG	1104.721	1	1104.721	1024.637	63	16.26460	38.14525	0.000000
D-D	21.825	1	21.825	303.114	63	4.81168	3.58911	0.062709
RH-RH	322.996	1	322.996	692.869	63	10.99793	29.36881	0.000051
ZM-ZM	281.129	1	281.129	2417.083	63	38.30449	8.30582	0.005243

Obě analýzy mohou napovědět, ale diskriminace může být dána i kombinací proměnných

13

### Diskriminační analýza

(Statistics > Mult/Exploratory > Discriminant)

grupovací proměnná – stav, který chceme určovat

nezávislá proměnná – výběr hodnot pro analýzu

14

### Diskriminační analýza – interpretace výsledků

Číselný výstup analýzy

**Celková Wilks Lambda**  
celková kvalita modelu s použitím všech proměnných (0 = nejlepší diskriminace)

Wilks lambda celého modelu při vyřazení dané proměnné

Variable	Wilks' Lambda	Partial Lambda	F	Remove p-value	Toler	1-Toler
NH-E5	0.5112327	1.000000	0.000000	1.000000	0.748612	0.251388
GOP	0.5115327	0.995423	0.028195	0.856678	0.589176	0.410824
EU-EU	0.5114241	0.995640	0.028321	0.856598	0.590895	0.391605
BeH	0.5647865	0.905185	5.970338	0.017668	0.493531	0.506469
ZYG-ZYG	0.5112422	0.999989	0.000612	0.999357	0.825532	0.174468
D-D	0.6090261	0.800088	9.989070	0.002523	0.892669	0.107332
RH-RH	0.5274881	0.969192	1.811848	0.183616	0.852432	0.147568

Unikátní příspěvek dané proměnné k diskriminaci

Variabilita proměnné nevysvětlená ostatními proměnnými

Variabilita proměnné vysvětlená kombinací ostatních proměnných v modelu

15

### Diskriminační analýza – interpretace výsledků

**Klasifikační funkce**

rozebrané funkce pro jednu a pro druhou kategorii

případ je přiřazen do té skupiny, pro kterou je výsledek vyšší

16

### Diskriminační analýza – hodnocení klasifikačního kritéria

Hodnocení úspěšnosti klasifikačního kritéria

**Klasifikační tabulka** – procentuální vyjádření úspěšnosti zařazení objektů do skupin

**Resubstituce** – klasifikační rovnici testujeme na stejném souboru, na kterém byla vytvořena

Group	Percent Correct	F	M
0	84.88493	28	5
1	81.25000	6	26
Total	83.07692	34	31

daleko lépe **křížové ověření (leave-one-out-cross-validation)** aplikace na nezávislý vzorek, případně rozdělení původního vzorku

17

### Diskriminační analýza – podle čeho se dál orientovat?

Co může dál napovědět?

**Mahalanobisova vzdálenost** – popisuje vzdálenost centroidů skupin (bere v úvahu korelaci mezi parametry a je nezávislá na jejich rozsahu)

**Posterior probability** – pravděpodobnost zařazení objektu do skupiny (p toho, že objekt patří do té které skupiny) – vychází z Makalanobisových vzdáleností ke skupinám a *a priori* pravděpodobnosti

18

### Diskriminační analýza – dopředná a zpětná eliminace proměnných

**„Step-wise“ analýza – výběr proměnných samotnou analýzou**

- proměnné jsou přidávány/ubírány podle jejich významu v modelu
- zpravidla je vybrán pouze zlomek původních proměnných

V tomto případě vybrány pouze tři proměnné

Forward stepwise – dopředná  
Backward stepwise – zpětná

19

### Kontingenční tabulky

**Test dobré shody**  
Testuje shodu reálné distribuce hodnot do n skupin s teoretickou distribucí

pozorované

sex	origin admixed	origin European	Row Totals
F	8	25	33
M	0	32	32
Totals	8	57	65

V případě platnosti nulové hypotézy je poměr mezi buňkami jednoho řádku v různých sloupcích nezávislý na výběru tohoto řádku

vs.

očekávané

sex	origin admixed	origin European	Row Totals
F	2.061538	28.93846	33.00000
M	0.938462	28.06154	32.00000
Totals	8.000000	57.00000	65.00000

Statistics > Basic statistics > Tables and banners > Options > Expected frequencies  
Advanced > Detailed Two-way Tables

20

### Korelační analýza

hodnocení vztahu mezi dvěma spojitými veličinami

vizuální posouzení

21

### Korelační analýza

Číselné vyjádření – korelační koeficienty  
Předpokladem použití parametrického testu je normalita rozložení

Vyhovuje? → Pearsonův korelační koeficient  
Basic statistics > Correlation matrices

Nevyhovuje? → Spearmanův korelační koeficient  
Non-parametrics > Correlations – pořadová korelace

22

### Regresní analýza

Vysvětluje, jak vysvětlovaná proměnná závisí na jiných proměnných (prediktorech). Model musí odpovídat typu vztahu – pokud je přímý, můžeme použít lineární model

$$Y = b_0 + b_1X + E$$

Dependent – závislá (vysvětlovaná) proměnná  
Multiple R – koeficient vícerozměrné korelace  
R2 – koeficient determinace – podíl modelem vysvětlované variability  
Adjusted R2 – podobný, ale bere v úvahu počet regresorů  
F, df a p – F test vztahů mezi závislou proměnnou a množinou nezávislých proměnných  
F = regresní průměr čtverců/reziduální průměr čtverců  
Standard error of estimate – směrodatná chyba odhadu – rozptýlení pozorovaných hodnot kolem přímky  
Intercept (Absolutní člen) – hodnota B0  
Std. Error – směrodatná chyba absolutního členu (následují testy Ho – Intercept je roven nule)  
b\* – standardizované koeficienty – umožňují porovnávat vliv jednotlivých proměnných

23

### Regresní analýza

Další výsledky  
Summary: regression results  
První tabulka – statistiky z předchozího souhrnného okna

Druhá tabulka – podrobnější výsledky regrese, včetně nestandardizovaného koeficientu (b) (ten standardizovaný ukazuje relativní příspěvek jednotlivých proměnných)

	B	Std. Err.	t	Sig.	Lower Bound	Upper Bound
(Constant)	-163.097	24.20522	-6.73617	0.00005	-211.524	-114.670
Výška	0.628228	0.064377	9.75612	0.00000	0.50000	0.75645

Pro každý koeficient jsou vypočítány hodnoty t-statistiky a p testující, zda je daný parametrický významně odlišný od 0 (jestli má proměnná v modelu své opodstatnění – součastí verifikace modelu).

24

### Regresní analýza

**Ověření předpokladů**

- 1) Správně specifikovaný model
- 2) Střední hodnota chybové složky je 0
- 3) Chybová složka má konstantní rozptyl
- 4) Jednotlivé složky chybového vektoru jsou nekorelované
- 5) Residuální složka má normální rozdělení

Perform residual analysis > Scatterplots  
> Predicted vs. residuals

Perform residual analysis > Residuals > Histogram of residuals

25

### Regresní analýza – správná podoba výsledků

**Ověření předpokladů**

- 1) Správně specifikovaný model
- 2) Střední hodnota chybové složky je 0
- 3) Chybová složka má konstantní rozptyl
- 4) Jednotlivé složky chybového vektoru jsou nekorelované
- 5) Residuální složka má normální rozdělení

26

### Regresní analýza – správná podoba výsledků

**Ověření předpokladů**

- 1) Správně specifikovaný model
- 2) Střední hodnota chybové složky je 0
- 3) Chybová složka má konstantní rozptyl
- 4) **Jednotlivé složky chybového vektoru jsou nekorelované**
- 5) **Residuální složka má normální rozdělení**

Perform residual analysis > Basics > Normal plot of residuals (Kvantilový graf)

V případě normality musí body ležet na proložené přímce.  
Pokud neleží (dá se dále ověřit testem rezidui) – odhady parametrů modelu a regr. rovnice jsou v pořádku, ale **ne významnost regr. parametrů a konfidenční intervaly**

27

### Regresní analýza – správná podoba výsledků

**Predikce**  
Predict dependent variable

**Compute confidence limits**  
Interval spolehlivosti pro průměrnou hodnotu odezvy  
Udává rozmezí, kde se s 95% spolehlivostí nachází true best fit populace

**Compute prediction limits (interval předpovědi)**  
Interval spolehlivosti pro individuální hodnotu odezvy

pokud použijete stejnou rovnici na další jedince dané výšky, bude se 95% z nich nacházet v daném rozmezí

28