

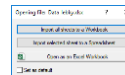
# Statistica

Mgr. Mikołáš Jurda, Ph.D.

## Ovládání programu

### Otevření/import dat

Data je možné importovat ze souborů různých typů – excel, csv, txt a také vložením ze schránky – je možné je také editovat



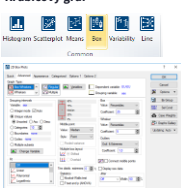
Práci ušetří, pokud jsou buňky v původním excelovském souboru ve správném formátu. Pokud po načtení formát proměnných neodpovídá našim požadavkům, jde formát upravit nastavením jednotlivých proměnných (dvojitě poklikání na buňku s názvem proměnné) – typicky nastavení grupovacích proměnných na *Type > Double*

Datový soubor a výstupy analýz je možné uchovávat v různých formátech. Možnost uložit vše v přehledném stromu nabízí *Workbook*. Výstupy je možné také ukládat do wordu

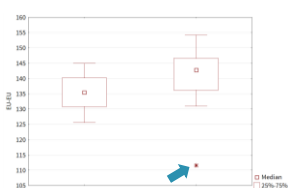


## Popisná statistika

### Vizuální posouzení dat Krabicový graf



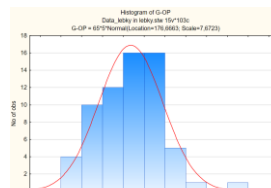
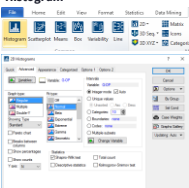
Pokud nezádáte grupovací proměnnou, zobrazí se graf pro celý soubor, pokud ano, pak



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

## Popisná statistika

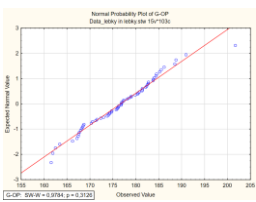
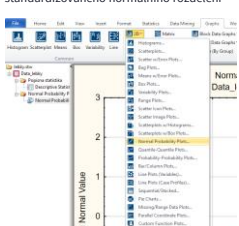
### Vizuální posouzení dat Histogram



Umožňuje posoudit rozložení a srovnat je se předpokládaným rozložením (linie). Nastavuje se jako *Fit type* v předchozím grafu.

## Popisná statistika

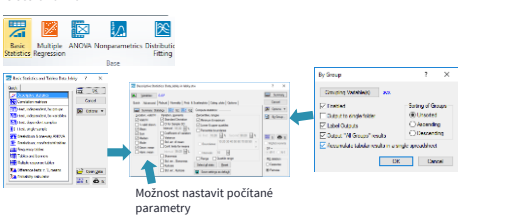
### Vizuální posouzení dat Normální pravděpodobnostní graf



Dovoluje rovněž identifikovat extrémní hodnoty

## Popisná statistika

### Číselná popisná statistika Číselná forma



Možnost nastavit počítané parametry

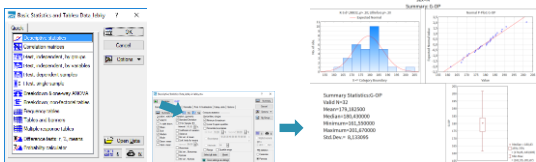
Variable	Descriptive Statistics (Data: volby in lobby stat)	Aggregate Results				
	Valid N	Mean	Median	Minimum	Maximum	Std Dev
G-GOP	33	174.2264	175.6400	161.6100	184.1100	8.471765
G-GOP	32	179.1825	180.4300	161.6500	201.8700	8.133895
G-GOP	0					

Pro všechna data zároveň

By group – pro skupiny zvlášť

## Popisná statistika

Souhrnné výsledky – histogram, krabicový graf, zvolené parametry a P-P plot



V základní, přednastavené podobě

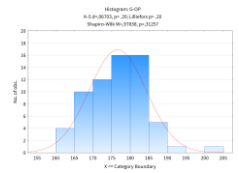
## Normalita dat

Grafické posouzení  
srovnání s normálním rozdělením (viz předchozí grafy)

Testování  
**Shapiro-Wilksův test**  
Statistics > Basic statistics > Descriptive statistics > Normality

Průběhy testů: Chi-Square, Kolmogorov-Smirnov, Lilliefors

Test	Signifikantní	Pravděpodobnost	Pravděpodobnost	Pravděpodobnost
Chi-Square	0,000	0,000	0,000	0,000
Kolmogorov-Smirnov	0,000	0,000	0,000	0,000
Lilliefors	0,000	0,000	0,000	0,000



Záhlaví Frequency table

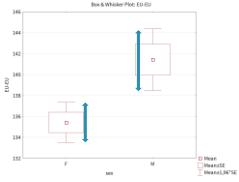
Záhlaví Histograms

## T-test

Nepárový dvouvýběrový t-test

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

Shoda rozptylů  
- testování je přímo součástí výsledků



Pokud data nesplňují

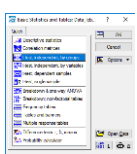


**Neparametrické alternativy**  
např. Mann-Whitney U-test

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

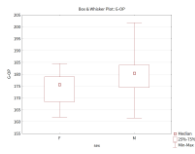
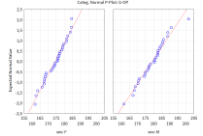
## T-test

Ověření předpokladů testu přímo v jeho dialogovém okně



Normalita rozložení ve skupinách  
Advanced > Categorized normal plots

Shoda rozptylů - graficky  
Advanced > Box & Whiskers plot



## T-test

Samotné výstupy testu – lze provést hromadně pro všechny zároven

Tasks: Grouping: sex (Data:\_hody\_in\_hodiny\_akt)

Variable	Mean	Mean	t-value	df	p	Valid N	Valid N	Std. Dev.	Std. Dev.	F-ratio	p
GDP	112,226	179,1625	-2,7262	63	0,00815	32	32	4,47769	6,13296	5,90134	0,19193
CS-EU	110,429	141,4396	-3,7715	63	0,00187	32	32	4,63532	6,46620	2,26439	0,14489
Skp-B	120,947	132,4963	-2,4193	63	0,01866	32	32	4,96192	6,36577	1,64638	0,16521
ZYU-ZYU	120,569	128,8208	-5,1763	63	0,00000	32	32	4,72546	5,26875	1,07291	0,30689
SD	22,2734	21,0364	1,8844	63	0,06279	32	32	2,36916	2,56275	1,16983	0,28497
SH-MS	33,0903	37,5622	-4,4193	63	0,00001	32	32	3,76134	3,46727	1,20199	0,26983
SH-M	38,9952	52,2660	-2,8207	63	0,00563	32	32	4,44934	6,03667	1,64019	0,17198
N-B	107,4394	119,7098	-2,3398	63	0,02326	32	32	3,34986	6,16487	1,48541	0,22678
N-L	109,1227	111,6422	-1,6447	63	0,11169	32	32	5,89169	6,12667	1,08169	0,30219



## Diskriminační analýza

Závislost jedné kvalitativní proměnné na několika kvantitativních proměnných

**Pro**  
• určení proměnných, které diskriminují mezi dvěma nebo více skupinami  
• ke klasifikaci objektů do různých skupin

**Předpoklady**  
• mnohozměrné normální rozdělení (především citlivé na odlehle hodnoty)  
• shoda skupinových kovariančních matic  
• proměnné nejsou redundantní

## Diskriminační analýza

### Jaké použít proměnné?

Význam mají **pouze ty proměnné**, které mají souvislost s kategoriální proměnnou

Redundantní proměnné snižují stabilitu modelu a mohou vést k nesmyslným výsledkům

### Hodnocení vztahu nezávislých proměnných a kategoriální proměnné

Korelační analýza a XY grafy

Hlavní komponenty a faktorová analýza

Diskriminační analýza

„Expertní“ znalost proměnných – pokud jsou redundantní, můžeme vyřadit ty proměnné zatížené chybami nebo vysokým počtem chybějících hodnot

## Diskriminační analýza

### Hledání proměnných

- Vztah ke kategoriální proměnné
- samostatný t-test pro jednotlivé proměnné – pro dvě skupiny vždy (*Basic statistics > t-test, independent, by groups*)



Variable	Mean	Std. Dev.	N	Significance (2-tailed)	df	F	p
EU-EU	15.428	1.975	12	.000	11	11.844	0.00257
BA-B	15.428	1.975	12	.000	11	11.844	0.00257
ZYG-ZYG	15.428	1.975	12	.000	11	11.844	0.00257
EU-EU	15.428	1.975	12	.000	11	11.844	0.00257
BA-B	15.428	1.975	12	.000	11	11.844	0.00257
ZYG-ZYG	15.428	1.975	12	.000	11	11.844	0.00257
EU-EU	15.428	1.975	12	.000	11	11.844	0.00257
BA-B	15.428	1.975	12	.000	11	11.844	0.00257
ZYG-ZYG	15.428	1.975	12	.000	11	11.844	0.00257
EU-EU	15.428	1.975	12	.000	11	11.844	0.00257
BA-B	15.428	1.975	12	.000	11	11.844	0.00257
ZYG-ZYG	15.428	1.975	12	.000	11	11.844	0.00257

- ANOVA – pro dvě a více skupin (*Basic statistics > Breakdown & One-way ANOVA; Analysis of Variance*) – výsledky jsou ekvivalentní

Variable	Effect	Sum of Squares	df	Mean Square	F	Significance
EU-EU	Corrected Total	11.844	11	1.0767		
	Between Groups	11.844	1	11.844	11.844	.00257
	Within Groups	0.000	10	.000		
BA-B	Corrected Total	11.844	11	1.0767		
	Between Groups	11.844	1	11.844	11.844	.00257
	Within Groups	0.000	10	.000		
ZYG-ZYG	Corrected Total	11.844	11	1.0767		
	Between Groups	11.844	1	11.844	11.844	.00257
	Within Groups	0.000	10	.000		

Může napovědět, ale diskriminace může být dána pouze kombinací proměnných

## Diskriminační analýza

### Hledání proměnných

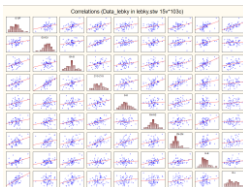
Korelační analýza – vztah proměnných mezi sebou

(*Statistics > Basic statistics and tables > Correlation matrices; Summary: Correlations*)

Variable	Mean	Std. Dev.	EU-EU	BA-B	ZYG-ZYG	EU-EU	BA-B	ZYG-ZYG	EU-EU	BA-B	ZYG-ZYG	N
EU-EU	15.428	1.975	1.000			1.000			1.000			12
BA-B	15.428	1.975	.143	1.000		.143	1.000		.143			12
ZYG-ZYG	15.428	1.975	.143	.143	1.000	.143	.143	1.000	.143	.143		12
EU-EU	15.428	1.975	1.000			1.000			1.000			12
BA-B	15.428	1.975	.143	1.000		.143	1.000		.143	1.000		12
ZYG-ZYG	15.428	1.975	.143	.143	1.000	.143	.143	1.000	.143	.143	1.000	12

### XY grafy

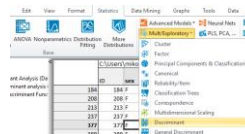
(*Statistics > Basic statistics and tables > Correlation matrices; Scatterplot matrix for selected correlations*)



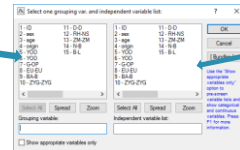
## Diskriminační analýza

### Diskriminační analýza

(*Statistics > Mult/Exploratory > Discriminant*)



### Variables



Grupovací proměnná – kategorie, do kterých bychom případně chtěli klasifikovat

Nezávislé proměnné

## Diskriminační analýza

### Číselný výstup analýzy

**Celkové Wilks lambda** – popisuje celkovou kvalitu modelu všech proměnných (0 = nejlepší diskriminace)

Wilks lambda celého modelu při vyloučení dané proměnné

Discriminant Function Analysis Summary (Data_Lobby in lobby.sav)	Wilk's Lambda	Partial Wilk's Lambda	F	Sig.	Lower Bound	Upper Bound
EU-EU	0.948177	0.999122	0.64746	0.82361	0.326232	0.463748
BA-B	0.484136	0.999760	0.93987	0.97345	0.364732	0.326820
ZYG-ZYG	0.527997	0.946436	3.06136	0.08115	0.429520	0.574480
EU-EU	0.494447	0.996920	0.46206	0.92180	0.719564	0.250544
BA-B	0.169044	0.834231	15.53348	0.01192	0.776250	0.223791
ZYG-ZYG	0.524703	0.941656	3.39194	0.07546	0.466161	0.324840
EU-EU	0.492712	0.996620	0.18132	0.67041	0.627469	0.472891
BA-B	0.562629	0.982807	0.94021	0.33654	0.496252	0.607448

Unikátní příspěvek proměnné k diskriminaci

Variabilita proměnné vysvětlená ostatními proměnnými

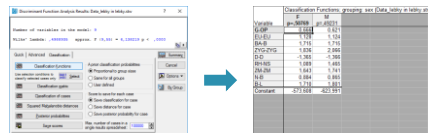
Var. proměnné vysvětlená kombinací ostatních proměnných v modelu

## Diskriminační analýza

### Číselný výstup analýzy

### Klasifikační funkce

Sada rovnic – objekt je zařazen do té skupiny, jejíž klasifikační funkce nabývá nejvyšší hodnoty



$$F = 0,666 \cdot G-OP + 1,128 \cdot EU-EU + 1,715 \cdot BA-B + \dots + (-573,608)$$

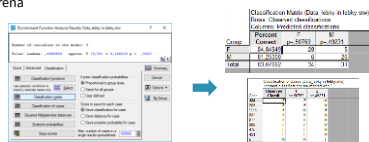
$$M = 0,621 \cdot G-OP + 1,124 \cdot EU-EU + 1,715 \cdot BA-B + \dots + (-623,991)$$

Jindy jako jedna rovnice, jejíž výsledek se porovnává se dělicím bodem

## Diskriminační analýza

**Hodnocení úspěšnosti klasifikačního kritéria**  
**Klasifikační tabulka** – procentuální vyjádření úspěšnosti zařazení objektů do skupin

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



**Křížové ověření (leave-one-out cross-validation)** – vybereme n-1 objektů, z nich vytvoříme kritérium a to pak aplikujeme na vypuštěný případ. Postup opakujeme se všemi dalšími případy

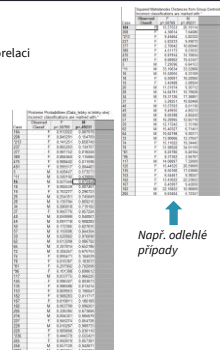
Aplikace na nezávislý vzorek – kritérium vytvoříme například pouze na části případů a ověříme na tom zbytku

## Diskriminační analýza

**Co může dále napovědět?**

Mahanobisova vzdálenost - popisuje vzdálenost centroidů skupin (bere v úvahu korelaci mezi parametry a je nezávislá na rozsahu parametrů)

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

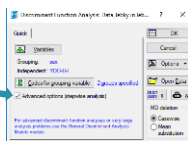


↑  
Např. odlehle případy

## Diskriminační analýza

„Step-wise“ analýza (dopředná a zpětná eliminace) – výběr proměnných samotnou analýzou

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



	Wilks' Lambda	Partial Eta Squared	F (Sig.)	Toler.	V. Toler.	
FW-YG	0.55562	0.841704	11.39024	0.067342	0.084970	0.115320
RH-AS	0.670534	0.838063	11.78516	0.001078	0.933334	0.066566
ZM-ZM	0.530334	0.964823	2.22404	0.141031	0.919489	0.089511

V našem případě (dopředná analýza) vybrány pouze tři proměnné

## Kontingenční tabulky

**Test dobré shody** (Pearsonův chí-kvadrát test)

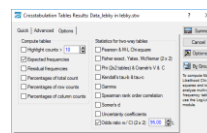
Testuje shodu reálné distribuce hodnot do n skupin s teoretickou distribucí.

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

je A nezávislé na B a naopak

	A	B	Σ
+	a	b	
-	c	d	
Σ			suma sum

Statistics > Basic statistics > Tables and banners... > Options > Expected frequencies

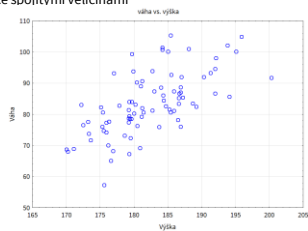
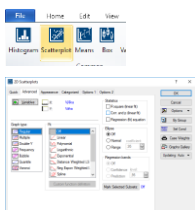


zobrazení výsledků  
 Advanced > Detailed Two-way Tables

## Korelační analýza

Hodnocení vztahu mezi dvěma a více spojitými veličinami

**Bodový graf**

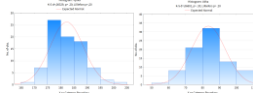


Vztah existuje, ale vztah není přesně lineární

## Korelační analýza

**Korelační koeficienty**

Předpokladem parametrických je normalita rozložení



**Pearsonův korelační koeficient**

$$R(X,Y) = \frac{E[(X - E(X)) \cdot (Y - E(Y))]}{\sqrt{D(X) \cdot D(Y)}}$$

Basic statistics > Correlation matrices



**Spearmanův korelační koeficient**

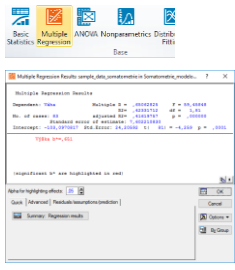
Non-parametrics > Correlations



## Regresní analýza

Vysvětluje vztah dvou a více proměnných. Jak **vysvětlovaná proměnná** závisí na jiných **proměnných (prediktorech)**. **Model musí odpovídat typu vztahu** – pokud je přímkový, můžeme použít lineární model.

$$Y = b_0 + b_1X + E$$



Dependent – závislá (vysvětlovaná proměnná)  
 Multiple R – koeficient vícerozměrné korelace  
 R2 – koeficient determinace – podíl vysvětlované variability  
 Adjusted R2 – podobný, ale bere v úvahu počet regresorů  
 F, df a p – F test vztahu mezi závislou proměnnou a množinou nezávislých proměnných  
 – F regreseční průměr čtverců/reziduální průměr čtverců  
 Standard error of estimate – směrodatná chyba odhadu – rozptýlení 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 porovnat vliv jednotlivých proměnných

## 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)

	Model	Sum of Squares	df	Mean Square	F	Sig.	Partial	eta Squared
Intercept	1	103.097	1	103.097	4.2017	0.00003		
Výška	2	0.85262	1	0.85262	0.10276	0.00002		

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

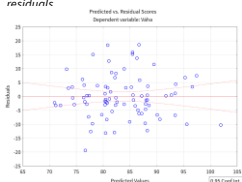
v našem případě – hmotnost = -103,097 + 1,024\*výška+E

## Regresní analýza

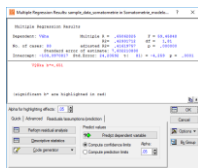
Ověření předpokladů:

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

Perform residual analysis > Scatterplots > Predicted vs. residual



Rezidua konstantně rozptýlena kolem nulové střední hodnoty

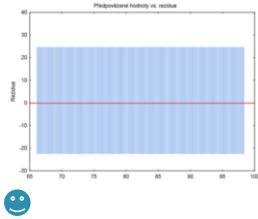


Nepřímková závislost

## Regresní analýza

Ověření předpokladů:

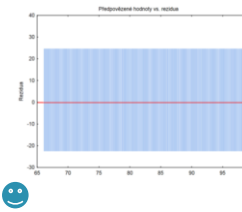
- 1) **Správně specifikovaný model**
- 2) **Střední hodnota chybové složky je rovna 0**
- 3) **Chybová složka má konstantní rozptyl**
- 4) **Jednotlivé složky chybového vektoru jsou nekorelované**
- 5) **Reziduální složka má normální rozdělení**



## Regresní analýza

Ověření předpokladů:

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



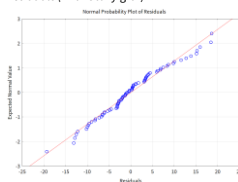
## Regresní analýza

Ověření předpokladů:

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

Nemusíme ověřovat – jde o nezávislé jedince

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 reziduí) – odhady parametrů modelu a regr. rovnice jsou v pořádku, ale **ne významnost regr. parametrů a konfidenční intervaly**

## Regresní analýza

### Predikce

*Predict dependent variable*

*Compute confidence limits*

interval spolehlivosti pro průměrnou hodnotu odezvy

udává rozmezí, ve kterém se s 95% spolehlivostí nachází „true best fit“ dané populace

*Compute prediction limits*

interval spolehlivosti pro individuální hodnotu odezvy

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

