

5. Porovnání empirického a teoretického rozložení

5.1. Motivace: Možnost použití statistických testů je podmíněna nějakými předpoklady o datech. Velmi často je to předpoklad o typu rozložení, z něhož získaná data pocházejí. Mnoho testů je založeno na předpokladu normality. Opomíjení předpokladů o typu rozložení může v praxi vést i ke zcela zavádějícím výsledkům, proto je nutné věnovat tomuto problému patřičnou pozornost.

5.2. Popis Kolmogorovova – Smirnovova testu a jeho Lilieforsovy varianty

Testujeme hypotézu, která tvrdí, že náhodný výběr X_1, \dots, X_n pochází z rozložení s distribuční funkcí $\Phi(x)$. Necht' $F_n(x) = \frac{1}{n} \text{card}\{i; X_i \leq x\}$ je výběrová distribuční funkce. Testovou statistikou je statistika $D_n = \sup_{-\infty < x < \infty} |F_n(x) - \Phi(x)|$. Nulovou hypotézu zamítáme na hladině významnosti α , když $D_n \geq D_n(\alpha)$, kde $D_n(\alpha)$ je tabelovaná kritická hodnota.

(Pro $n \geq 30$ lze $D_n(\alpha)$ aproximovat výrazem $\sqrt{\frac{1}{2n} \ln \frac{2}{\alpha}}$.)

Upozornění: Nulová hypotéza musí specifikovat distribuční funkci zcela přesně, včetně všech jejích případných parametrů. Např. K-S test lze použít pro testování hypotézy, že náhodný výběr X_1, \dots, X_n pochází z rozložení $Rs(0,1)$, což se využívá při testování generátorů náhodných čísel.

Lilieforsova modifikace Kolmogorovova – Smirnovova testu

Necht' nulová hypotéza tvrdí, že náhodný výběr pochází z normálního rozložení, jehož parametry μ a σ^2 neznáme. Tyto parametry musíme odhadnout z dat. Tím se změní rozložení testové statistiky D_n . V takovém případě jde o **Lilieforsovu modifikaci** Kolmogorovova – Smirnovova testu. Příslušné modifikované kvantily byly určeny pomocí simulačních studií.

Poznámka ke K-S testu ve STATISTICE: Test normality poskytuje hodnotu testové statistiky (ozn. d) a dvě p-hodnoty. První se vztahuje k případu, kdy μ a σ^2 známe předem, druhá (ozn. Liliefors p) se vztahuje k případu, kdy μ a σ^2 neznáme. Objeví-li se ve výstupu $p = \text{n.s.}$ (tj. non significant), pak hypotézu o normalitě nezamítáme na hladině významnosti 0,05.

5.3. Příklad: Jsou dány hodnoty 10, 12, 8, 9, 16. Pomocí Lilieforsovy varianty K-S testu ověřte na hladině významnosti 0,05, zda tato data pocházejí z normálního rozložení.

Řešení: Odhadem střední hodnoty je výběrový průměr $m = 11$, odhadem rozptylu je výběrový rozptyl $s^2 = 10$. Uspořádaný náhodný výběr je (8, 9, 10, 12, 16). Vypočteme hodnoty výběrové distribuční funkce:

$$x < 8 : F_5(x) = 0, 8 \leq x < 9 : F_5(x) = \frac{1}{5} = 0,2, 9 \leq x < 10 : F_5(x) = \frac{2}{5} = 0,4,$$

$$10 \leq x < 12 : F_5(x) = \frac{3}{5} = 0,6, 12 \leq x < 16 : F_5(x) = \frac{4}{5} = 0,8, x \geq 16 : F_5(x) = 1$$

Hodnoty teoretické distribuční funkce $\Phi_T(x)$ v bodech 8, 9, 10, 12, 16:

$$\Phi_T(8) = \Phi\left(\frac{8-11}{\sqrt{10}}\right) = \Phi(-0,95) = 1 - \Phi(0,95) = 1 - 0,82894 = 0,17106$$

$$\Phi_T(9) = \Phi\left(\frac{9-11}{\sqrt{10}}\right) = \Phi(-0,63) = 1 - \Phi(0,63) = 1 - 0,73565 = 0,26435$$

$$\Phi_T(10) = \Phi\left(\frac{10-11}{\sqrt{10}}\right) = \Phi(-0,32) = 1 - \Phi(0,32) = 1 - 0,62552 = 0,37448$$

$$\Phi_T(12) = \Phi\left(\frac{12-11}{\sqrt{10}}\right) = \Phi(0,32) = 0,62552$$

$$\Phi_T(16) = \Phi\left(\frac{16-11}{\sqrt{10}}\right) = \Phi(1,58) = 0,94295$$

(Φ je distribuční funkce rozložení $N(0,1)$.)

Rozdíly mezi výběrovou distribuční funkcí $F_5(x)$ a teoretickou distribuční funkcí $\Phi_T(x)$:

$$d_1 = 0,2 - 0,17106 = 0,02894; d_2 = 0,4 - 0,26435 = 0,13565; d_3 = 0,6 - 0,37448 = 0,22552;$$

$$d_4 = 0,8 - 0,62552 = 0,17448; d_5 = 1 - 0,94295 = 0,05705.$$

Testová statistika: $D_5 = 0,22552$, modifikovaná kritická hodnota pro $n = 5$, $\alpha = 0,05$ je 0,343. Protože $0,22552 < 0,343$, hypotézu o normalitě nezamítáme na hladině významnosti 0,05.

5.4. Popis Shapirova – Wilkova testu

Testujeme hypotézu, která tvrdí, že náhodný výběr X_1, \dots, X_n pochází z rozložení $N(\mu, \sigma^2)$.

Testová statistika má tvar:

$$W = \frac{\left[\sum_{i=1}^m a_i^{(n)} (X_{(n-i+1)} - X_{(i)}) \right]^2}{\sum_{i=1}^n (X_i - M)^2},$$

kde $m = n/2$ pro n sudé a $m = (n-1)/2$ pro n liché. Koeficienty $a_i^{(n)}$ jsou tabelovány.

Na testovou statistiku W lze pohlížet jako na korelační koeficient mezi uspořádanými pozorováními a jim odpovídajícími kvantily standardizovaného normálního rozložení. V případě, že data vykazují perfektní shodu s normálním rozložením, bude mít W hodnotu 1. Hypotézu o normalitě tedy zamítneme na hladině významnosti α , když se na této hladině neprokáže korelace mezi daty a jim odpovídajícími kvantily rozložení $N(0,1)$.

Lze také říci, že $S - W$ test je založen na zjištění, zda body v Q-Q grafu jsou významně odlišné od regresní přímky proložené těmito body.

(S-W test se používá především pro výběry menších rozsahů, $n < 50$, ale v systému STATISTICA je implementováno jeho rozšíření i na výběry velkých rozsahů, kolem 2000.)

Výpočet pomocí systému STATISTICA:

V sedmi náhodně vybraných prodejnách byly zjištěny následující ceny určitého druhu zboží (v Kč): 35, 29, 30, 33, 45, 33, 36. Rozhodněte pomocí Lilieforsovy varianty K-S testu a S-W testu na hladině významnosti 0,05, zda lze tyto ceny považovat za realizace náhodného výběru z normálního rozložení.

Řešení:

Otevřeme nový datový soubor o jedné proměnné a 7 případech. Do proměnné X jsou zapíšeme zjištěné ceny.

Statistiky – Základní statistiky a tabulky – Tabulky četností - OK – Proměnné X , OK – Normalita – zaškrtneme Lilieforsův test a Shapiro - Wilksův W test – Testy normality

Proměnná	Testy normality (Tabulka22)				
	N	max D	Lilliefors p	W	p
x	7	0,240290	p > .20	0,868661	0,180679

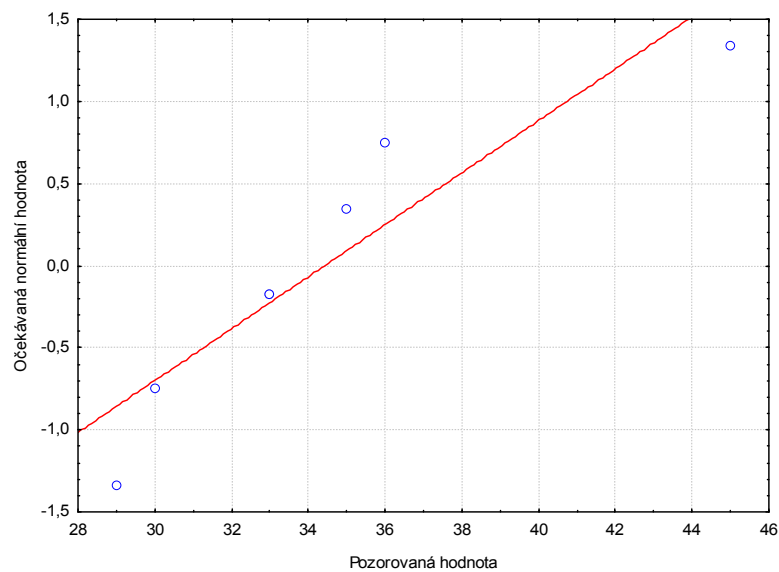
V tabulce je uvedena hodnota testové statistiky pro Lilieforsův test ($d = 0,24029$) a pro S-W test ($W = 0,86866$) a odpovídající p-hodnoty. Lilieforsovo p

je počítáno na základě parametrů odhadnutých z dat. V našem případě $p > 0,2$ a pro S-W test $p = 0,18068$. Ani jeden z testů nezamítá nulovou hypotézu o normalitě.

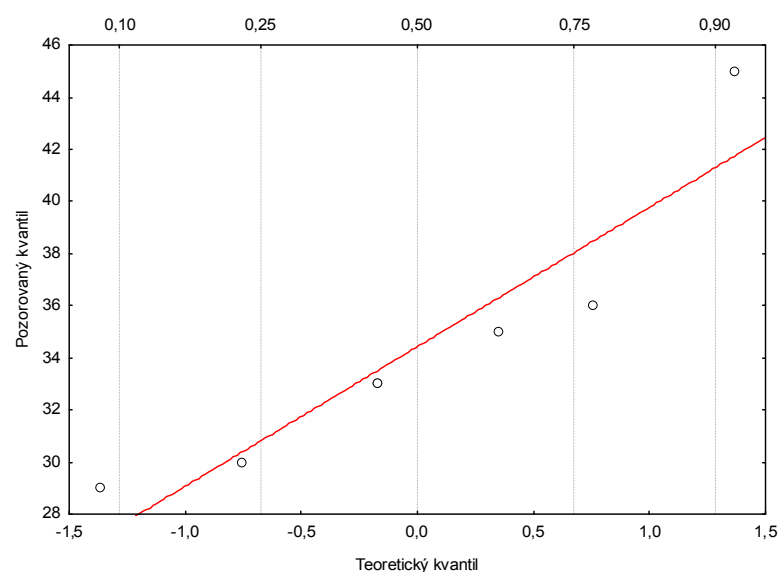
Výpočet doplníme normálním pravděpodobnostním grafem a kvantil – kvantilovým grafem:

Graphs – 2D Graphs - Normal Probability Plots (resp. Quantile- Quantile plot)- Variables X – OK.

N-P plot:



Q-Q plot:



5.5. Poznámka: Další testy normality

Existují testy normality založené na výběrové šikmosti a špičatosti. Pro náhodnou veličinu s normálním rozložením platí, že její šikmost i špičatost jsou nulové. Pro výběr z normálního rozložení by tedy výběrová šikmost a špičatost měly být blízké 0.

Nechť X_1, \dots, X_n je náhodný výběr.

$$\text{Výběrová šikmost: } A_3 = \frac{\frac{1}{n} \sum_{i=1}^n (X_i - M)^3}{\left[\sqrt{\frac{1}{n} \sum_{i=1}^n (X_i - M)^2} \right]^3}$$

$$\text{Výběrová špičatost: } A_4 = \frac{\frac{1}{n} \sum_{i=1}^n (X_i - M)^4}{\left[\sqrt{\frac{1}{n} \sum_{i=1}^n (X_i - M)^2} \right]^4} - 3$$

Lze dokázat, že pro výběr z normálního rozložení platí:

$$E(A_3) = 0, \quad D(A_3) = \frac{6(n-2)}{(n+1)(n+3)}, \quad E(A_4) = -\frac{6}{n+1}, \quad D(A_4) = \frac{24n(n-2)(n-3)}{(n+1)^2(n+3)(n+5)}.$$

Pro $n \rightarrow \infty$ se statistiky $A_3 \sqrt{n}$ a $A_4 \sqrt{n}$ asymptoticky řídí normálním rozložením.

Test založený na šikmosti zamítne hypotézu o normalitě na asymptotické hladině významnosti α , když

$$U_3 = \frac{|A_3|}{\sqrt{D(A_3)}} \geq u_{1-\alpha/2}.$$

D'Agostinův test: zavedeme pomocné veličiny

$$b = \frac{3(n^2 + 27n - 70)(n+1)(n+3)}{(n-2)(n+5)(n+7)(n+9)},$$

$$W^2 = \sqrt{2(b-1)} - 1,$$

$$d = \frac{1}{\sqrt{\ln W}}, \quad a = \sqrt{\frac{2}{W^2 - 1}}.$$

Testová statistika má tvar $Z_3 = d \cdot \ln \left[\frac{U_3}{a} + \sqrt{\left(\frac{U_3}{a}\right)^2 + 1} \right]$ a platí, že má přibližně rozložení $N(0,1)$. Pro $n > 8$ zamítáme hypotézu o normalitě pokud $|Z_3| \geq u_{1-\alpha/2}$.

Test založený na špičatosti zamítne hypotézu o normalitě na asymptotické hladině významnosti α , když

$$U_4 = \frac{|A_4 - E(A_4)|}{\sqrt{D(A_4)}} \geq u_{1-\alpha/2}.$$

Také v tomto případě existuje D'Agostinova modifikace testu, nebudeme ji ale uvádět. Z dalších testů normality lze jmenovat např. **Andersonův-Darlingův** nebo **Jarque-Beraův** test.

5.6. Popis testu dobré shody v diskrétním a spojitém případě

Testujeme hypotézu, která tvrdí, že náhodný výběr X_1, \dots, X_n pochází z rozložení s distribuční funkcí $\Phi(x)$.

- Je-li distribuční funkce spojitá, pak data rozdělíme do r třídících intervalů (u_j, u_{j+1}) , $j = 1, \dots, r$. Zjistíme absolutní četnost n_j j -tého třídícího intervalu a vypočteme pravděpodobnost p_j , že náhodná veličina X s distribuční funkcí $\Phi(x)$ se bude realizovat v j -tém třídícím intervalu. Platí-li nulová hypotéza, pak $p_j = \Phi(u_{j+1}) - \Phi(u_j)$.
- Má-li distribuční funkce nejvýše spočetně mnoho bodů nespojitosti, pak místo třídících intervalů použijeme varianty $x_{[j]}$, $j = 1, \dots, r$. Pro variantu $x_{[j]}$ zjistíme absolutní četnost n_j a vypočteme pravděpodobnost p_j , že náhodná veličina X s distribuční funkcí $\Phi(x)$ se bude realizovat variantou $x_{[j]}$. Platí-li nulová hypotéza, pak $p_j = \Phi(x_{[j]}) - \lim_{x \rightarrow x_{[j]}^-} \Phi(x) = P(X = x_{[j]})$.

Testová statistika:
$$K = \sum_{j=1}^r \frac{(n_j - np_j)^2}{np_j}.$$

Platí-li nulová hypotéza, pak $K \approx \chi^2(r-1-p)$, kde p je počet odhadovaných parametrů daného rozložení. (Např. pro normální rozložení $p = 2$, protože z dat odhadujeme střední hodnotu a rozptyl.) Nulovou hypotézu zamítáme na asymptotické hladině významnosti α , když $K \geq \chi^2_{1-\alpha}(r-1-p)$. Aproximace se považuje za vyhovující, když tzv. teoretické četnosti $np_j \geq 5$, $j = 1, \dots, r$.

Upozornění: Hodnota testové statistiky K je silně závislá na volbě třídících intervalů. Navíc při nesplnění podmínky $np_j \geq 5$, $j = 1, \dots, r$ je třeba některé intervaly resp. varianty slučovat, což vede ke ztrátě informace.

5.7. Příklad (test dobré shody pro diskrétní rozložení): Byl zjišťován počet poruch určitého zařízení za 100 hodin provozu ve 150 disjunktních 100 h intervalech. Výsledky měření:

Počet poruch za 100 hodin provozu 0 1 2 3 4 a víc
 Absolutní četnost 52 48 36 10 4

Na asymptotické hladině významnosti 0,05 testujte hypotézu, že náhodný výběr X_1, \dots, X_{150} pochází z rozložení $Po(1,2)$.

Řešení:

Pravděpodobnost, že náhodná veličina s rozložením $Po(\lambda)$, kde $\lambda = 1,2$ bude nabývat hodnot p_0, \dots, p_4 a víc je

$$p_j = \frac{\lambda^j}{j!} e^{-\lambda} = \frac{1,2^j}{j!} e^{-1,2}, j=0,1,2,3, p_4 = 1 - (p_0 + p_1 + p_2 + p_3).$$

Výpočty potřebné pro stanovení testové statistiky K uspořádáme do tabulky.

j	n_j	p_j	np_j	$\frac{(n_j - np_j)^2}{np_j}$
0	52	0,301	150.0,301=45,15	1,039
1	48	0,361	150.0,361=54,15	0,698
2	36	0,217	150.0,217=32,55	0,366
3	10	0,087	150.0,087=13,05	0,713
4	4	0,034	150.0,034=5,1	0,237

$K = 1,039 + 0,698 + 0,713 + 0,237 = 3,053$, $r = 5$, $\chi^2_{0,95}(4) = 9,488$. Protože $3,053 < 9,488$, nulovou hypotézu nezamítáme na asymptotické hladině významnosti 0,05.

Výpočet pomocí systému STATISTICA:

Vytvoříme datový soubor o dvou proměnných (POČET a ČETNOST) a pěti případech a zapíšeme do něj hodnoty 0 1 2 3 4 a 52 48 36 10 4.

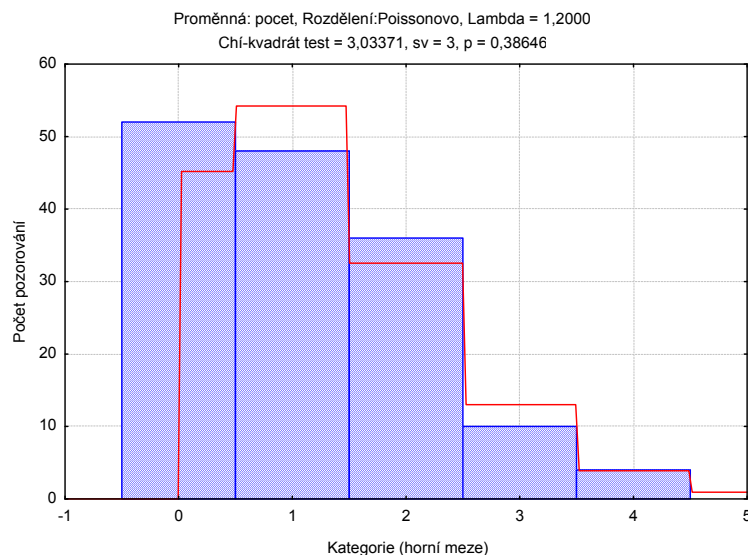
Statistiky – Prokládání rozdělení – Diskrétní rozdělení – Poissonovo – OK – Proměnná POČET – Proměnná vah ČETNOST – Stav zapnuto – OK – Parametry Lambda 1,2, OK.

Proměnná: pocet, Rozdělení:Poissonovo, Lambda = 1,2000 (Tabulka4) Chi-kvadrát = 3,03371, sv = 3, p = 0,38646								
Kategorie	Pozorované Četnosti	Kumulativ. Pozorované	Procent Pozorované	Kumul. % Pozorované	Očekáv. Četnosti	Kumulativ. Očekáv.	Procent Očekáv.	Kumul. % Očekáv.
<= 0,00000	52	52	34,66667	34,6667	45,17914	45,1791	30,11943	30,1194
1,00000	48	100	32,00000	66,6667	54,21495	99,3941	36,14330	66,2627
2,00000	36	136	24,00000	90,6667	32,52897	131,9231	21,68598	87,9487
3,00000	10	146	6,66667	97,3333	13,01159	144,9347	8,67439	96,6231
< Nekonečno	4	150	2,66667	100,0000	5,06535	150,0000	3,37690	100,0000

Ve výstupní tabulce je uvedena hodnota testového kritéria (3,03371) a odpovídající p-hodnota (0,38646). Nulová hypotéza se tedy nezamítá na asymptotické hladině významnosti 0,05.

(Podmínky dobré aproximace jsou splněny, všechny teoretické četnosti - uvedené ve sloupci Očekávané četnosti – jsou větší než 5.)

Pro vytvoření grafu se vrátíme do Proložení diskretních rozložení – Základní výsledky – Graf pozorovaného a očekávaného rozdělení.



5.8. Příklad (test dobré shody pro spojité rozložení): Byl pořízen náhodný výběr rozsahu $n = 100$. Jeho číselné realizace byly rozříděny do 5 ekvidistantních třídících intervalů o délce 0,04, přičemž dolní mez prvního třídícího intervalu je 3,92. Absolutní četnosti jednotlivých třídících intervalů jsou: 11, 20, 44, 19, 6.

Výběrový průměr se realizoval hodnotou $m = 4,02$ a výběrová směrodatná odchylka hodnotou $s = 0,04$.

Na asymptotické hladině významnosti 0,05 testujte hypotézu, že náhodný výběr pochází z normálního rozložení.

Řešení:

Výpočty potřebné pro stanovení testové statistiky K uspořádáme do tabulky. Přitom symbolem Φ značíme distribuční funkci rozložení $N(\mu, \sigma^2)$, kde $\mu = 4,02$ a $\sigma = 0,04$.

(u_j, u_{j+1})	n_j	$p_j = \Phi(u_{j+1}) - \Phi(u_j)$	np_j	$(n_j - np_j)^2$	$\frac{(n_j - np_j)^2}{np_j}$
$(3,92, 3,96)$	11	0,060598	6,0598	24,4060	4,0276
$(3,96, 4,00)$	20	0,241730	24,1730	17,4142	0,7204
$(4,00, 4,04)$	44	0,382925	38,2925	32,5756	0,8507
$(4,04, 4,08)$	19	0,241730	24,1730	26,7608	1,1070
$(4,08, 4,12)$	6	0,060598	6,0598	0,0036	0,0006

$$K = 4,0276 + 0,7204 + 0,8507 + 1,1070 + 0,0006 = 6,7063$$

$$\text{Kritický obor: } W = \langle \chi^2_{1-\alpha}(r-1-p), \infty \rangle = \langle \chi^2_{0,95}(5-1-2), \infty \rangle = \langle 5,9915, \infty \rangle$$

Protože testová statistika se realizuje v kritickém oboru, hypotézu o normalitě zamítáme na asymptotické hladině významnosti 0,05.

Výpočet pomocí systému STATISTICA:

Protože nemáme k dispozici původní data, ale jenom třídící intervaly a jejich četnosti, do nového datového souboru o dvou proměnných x_j a n_j zadáme středy třídících intervalů a jejich absolutní četnosti:

	1	2
	x_j	n_j
1	3,94	11
2	3,98	20
3	4,02	44
4	4,06	19
5	4,1	6

Statistiky – Prokládání rozdělení – ponecháme implicitní nastavení pro Normální rozdělení – OK – Proměnná x_j – klikneme na ikonu se závažím – Proměnná vah n_j – Stav Zapnuto – OK – Parametry – Počet kategorií 5, Průměr 4,02, Rozptyl 0,0016, OK.

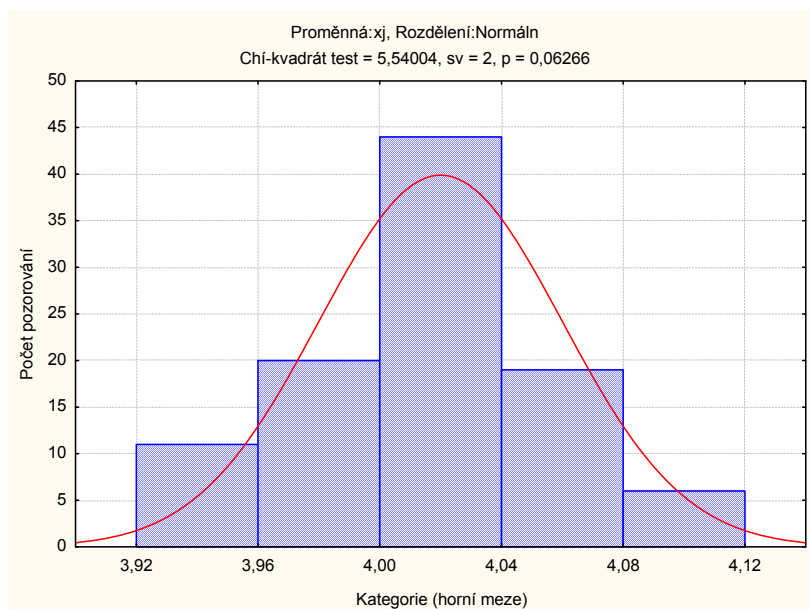
Dostaneme výstupní tabulku:

Horní hranice	Proměnná: xj, Rozdělení: Normální (Tabulka10) Chi-kvadrát = 5,54004, sv = 2, p = 0,06266							
	Pozorované Četnosti	Kumulativ. Pozorované	Procent Pozorované	Kumul. % Pozorované	Očekáv. Četnosti	Kumulativ. Očekáv.	Procent Očekáv.	Kumul. % Očekáv.
<= 3,96000	11	11	11,00000	11,0000	6,68072	6,6807	6,68072	6,6807
4,00000	20	31	20,00000	31,0000	24,17303	30,8538	24,17303	30,8538
4,04000	44	75	44,00000	75,0000	38,29249	69,1462	38,29249	69,1462
4,08000	19	94	19,00000	94,0000	24,17303	93,3193	24,17303	93,3193
< Nekonečno	6	100	6,00000	100,0000	6,68072	100,0000	6,68072	100,0000

V záhlaví výstupní tabulky je uvedena hodnota testového kritéria (5,54004), počet stupňů volnosti = 2 a p-hodnota (0,06266). Nulová hypotéza se tedy nezamítá na asymptotické hladině významnosti 0,05.

Rozdíl oproti ručnímu výpočtu je způsoben tím, že systém STATISTICA uvažuje první interval $(-\infty, 3,96)$ a poslední interval $(4,08, \infty)$.

Pro vytvoření grafu se vrátíme do Proložení spojitých rozdělení – Základní výsledky – Graf pozorovaného a očekávaného rozdělení.



5.9. Poznámka: Test dobré shody může být použit i v těch případech, kdy rozložení, z něhož daný náhodný výběr pochází, neodpovídá nějakému známému rozložení (např. exponenciálnímu, normálnímu, Poissonovu, ...), ale je určeno intuitivně nebo na základě zkušenosti.

5.10. Příklad: Ve svých pokusech pozoroval J.G. Mendel 10 rostlin hrachu a na každé z nich počet žlutých a zelených semen. Výsledky pokusu:

č.rostliny	1	2	3	4	5	6	7	8	9	10
počet žlutých semen	25	32	14	70	24	20	32	44	50	44
počet zelených semen	11	7	5	27	13	6	13	9	14	18
celkem	36	39	19	97	37	26	45	53	64	62

Z genetických modelů vyplývá, že pravděpodobnost výskytu žlutého semene by měla být 0,75 a zeleného 0,25. Na asymptotické hladině významnosti 0,05 testujte hypotézu, že výsledky Mendelových pokusů se shodují s modelem.

Řešení:

Výpočty potřebné pro stanovení testové statistiky K uspořádáme do tabulky.

j	n_j	p_j	np_j	$\frac{(n_j - np_j)^2}{np_j}$
1	25	0,75	$36 \cdot 0,75 = 27$	0,148148
2	32	0,75	$39 \cdot 0,75 = 29,25$	0,258547
⋮	⋮	⋮	⋮	⋮
10	44	0,75	$62 \cdot 0,75 = 46,5$	0,134409

$K = 0,148148 + 0,258547 + \dots + 0,134409 = 1,797495$, $r = 10$, $\chi^2_{0,95}(9) = 16,9$.
Protože $1,797495 < 16,9$, nulovou hypotézu nezamítáme na asymptotické hladině významnosti 0,05.

Výpočet pomocí systému STATISTICA:

Vytvoříme datový soubor se třemi proměnnými Celkem, X a Y a 10 případy. Do proměnné Celkem zapíšeme celkový počet žlutých a zelených semen, do X zapíšeme pozorované absolutní četnosti žlutých semen, do proměnné Y vypočítané teoretické četnosti (v našem případě $\text{Celkem} \cdot 0,75$).

Statistiky – Neparametrická statistika – Pozorované vs. očekávané χ^2 – Proměnné Pozorované četnosti X, Očekávané četnosti Y, OK – Výpočet.

Pozorované vs. očekávané četnosti (Mendel hrach Chi-Kvadr. = 1,797495 sv = 9 p = ,994280 POZN.: Nestejné součty pozor. a oček. četností)				
Případ	pozorov. X	očekáv. Y	P - O	(P-O)^2 /O
C: 1	25,0000	27,0000	-2,00000	0,148148
C: 2	32,0000	29,2500	2,75000	0,258547
C: 3	14,0000	14,2500	-0,25000	0,004386
C: 4	70,0000	72,7500	-2,75000	0,103952
C: 5	24,0000	27,7500	-3,75000	0,506757
C: 6	20,0000	19,5000	0,50000	0,012821
C: 7	32,0000	33,7500	-1,75000	0,090741
C: 8	44,0000	39,7500	4,25000	0,454403
C: 9	50,0000	48,0000	2,00000	0,083333
C: 10	44,0000	46,5000	-2,50000	0,134409
Sčt	355,0000	358,5000	-3,50000	1,797495

Ve výstupní tabulce najdeme hodnotu testové statistiky (Chi-Kvadr. = 1,797495) a odpovídající p-hodnotu, kterou porovnáme se zvolenou hladinou významnosti. V našem případě je p-hodnota 0,99428, takže nulová hypotéza se nezamítá na asymptotické hladině významnosti 0,05.