

# Úvod do praktické fyziky, cvičení 12

## Metoda nejmenších čtverců potřetí, test kvality fitu

Jan Matoušek

21. 12. 2020



UNIVERZITA KARLOVA  
Matematicko-fyzikální  
fakulta

## Metoda nejmenších čtverců (cv. 10):

- Necht' mezi veličinou  $x$  a  $y$  platí vztah

$$y = \lambda(x|\boldsymbol{\theta}).$$

- Modelová funkce  $\lambda$  závisí na parametrech

$$\boldsymbol{\theta} = (\theta_1, \theta_2, \dots, \theta_m).$$

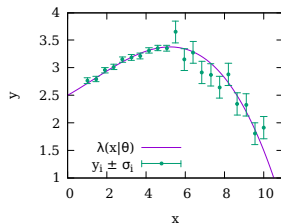
- V bodech  $x_1, x_2, \dots, x_N$  jsme naměřili hodnoty  $y_1, y_2, \dots, y_N$  se standardními odchylkami  $\sigma_1, \sigma_2, \dots, \sigma_N$ .

- $y_i \in N(\lambda(x_i|\boldsymbol{\theta}), \sigma_i)$

- Nejlepší odhad  $\hat{\boldsymbol{\theta}}$  je ten, který minimalizuje  $\chi^2$

$$\chi^2(\boldsymbol{\theta}|\mathbf{x}, \mathbf{y}, \boldsymbol{\sigma}) = \sum_{i=1}^N \frac{[y_i - \lambda(x_i|\boldsymbol{\theta})]^2}{\sigma_i^2}$$

(je to ten, který maximalizuje  $\ln L$ ).



## Lineární model ( $\lambda(x|\boldsymbol{\theta})$ lineární v $\theta_i$ )

$$\lambda(x|\boldsymbol{\theta}) = \theta_1 f_1(x) + \theta_2 f_2(x) + \dots + \theta_M f_M(x),$$

$$\frac{\partial \chi^2}{\partial \theta_i} = 0 \quad \text{soustava lineárních rovnic.}$$

→ lze řešit analyticky (viz cvičení 10 a 11).

## Nelineární model

$$\frac{\partial \chi^2}{\partial \theta_i} = 0 \quad \text{soustava nelineárních rovnic.}$$

→ obecně nutno řešit numericky.

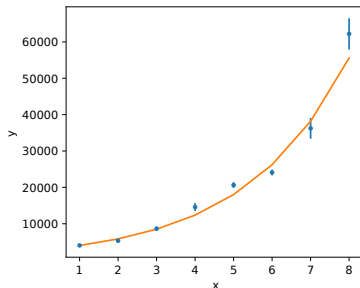
# Nelineární metoda nejmenších čtverců

## Příklad (upf\_cv12\_covid.py)

- Data: `upf_cv12_covid.txt` – aktivní případy COVID-19 od 24. 8. do 18. 10. [MZ ČR].
- Sloupce: číslo týdne  $x_i = i$  a týdenní průměr  $y_i = \bar{n}_i$ ,  $\sigma_i^2 = \frac{1}{6} \sum_{j=1}^7 (n_{ij} - \bar{n}_i)^2$ .
- Model: exponenciální růst  $\lambda(x|a, b) = ae^{bx}$ .

### Nelineární fit

- Numerická minimizace  $\chi^2$ , např. pomocí `scipy.optimize.curve_fit()`
- Někdy se vyplatí použít logaritmickou škálu.



### Linearizace fitu

$$y_i \rightarrow \ln y_i,$$

$$\lambda(x|a, b) \rightarrow \ln \lambda(x|a, b) = \ln a + bx,$$

$$\sigma_i \rightarrow \sigma_{\ln y_i} = \frac{\partial \ln y_i}{\partial y_i} \sigma_i = \frac{\sigma_i}{y_i}.$$

- Lineární regrese, např. `numpy.polyfit()`
- Riziko: I když  $y_i \in N(\lambda(x_i), \sigma_i)$ , po transformaci už to platit nemusí.

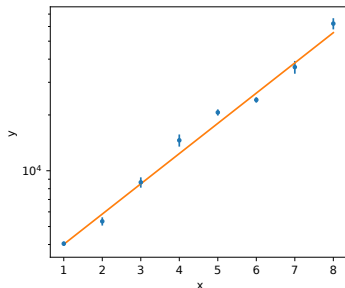
# Nelineární metoda nejmenších čtverců

## Příklad (upf\_cv12\_covid.py)

- Data: `upf_cv12_covid.txt` – aktivní případy COVID-19 od 24. 8. do 18. 10. [MZ ČR].
- Sloupce: číslo týdne  $x_i = i$  a týdenní průměr  $y_i = \bar{n}_i$ ,  $\sigma_i^2 = \frac{1}{6} \sum_{j=1}^7 (n_{ij} - \bar{n}_i)^2$ .
- Model: exponenciální růst  $\lambda(x|a, b) = ae^{bx}$ .

### Nelineární fit

- Numerická minimizace  $\chi^2$ , např. pomocí `scipy.optimize.curve_fit()`
- Někdy se vyplatí použít logaritmickou škálu.



### Linearizace fitu

$$y_i \rightarrow \ln y_i,$$

$$\lambda(x|a, b) \rightarrow \ln \lambda(x|a, b) = \ln a + bx,$$

$$\sigma_i \rightarrow \sigma_{\ln y_i} = \frac{\partial \ln y_i}{\partial y_i} \sigma_i = \frac{\sigma_i}{y_i}.$$

- Lineární regrese, např. `numpy.polyfit()`
- Riziko: I když  $y_i \in N(\lambda(x_i), \sigma_i)$ , po transformaci už to platit nemusí.

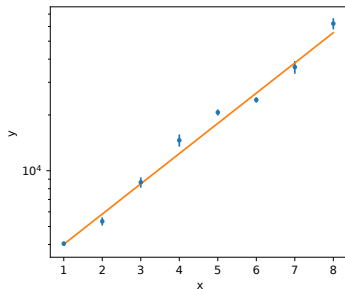
# Nelineární metoda nejmenších čtverců

## Příklad (upf\_cv12\_covid.py)

- Data: `upf_cv12_covid.txt` – aktivní případy COVID-19 od 24. 8. do 18. 10. [MZ ČR].
- Sloupce: číslo týdne  $x_i = i$  a týdenní průměr  $y_i = \bar{n}_i$ ,  $\sigma_i^2 = \frac{1}{6} \sum_{j=1}^7 (n_{ij} - \bar{n}_i)^2$ .
- Model: exponenciální růst  $\lambda(x|a, b) = ae^{bx}$ .

### Nelineární fit

- Numerická minimizace  $\chi^2$ , např. pomocí `scipy.optimize.curve_fit()`
- Někdy se vyplatí použít logaritmickou škálu.



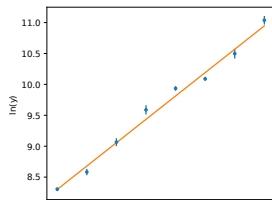
### Linearizace fitu

$$y_i \rightarrow \ln y_i,$$

$$\lambda(x|a, b) \rightarrow \ln \lambda(x|a, b) = \ln a + bx,$$

$$\sigma_i \rightarrow \sigma_{\ln y_i} = \frac{\partial \ln y_i}{\partial y_i} \sigma_i = \frac{\sigma_i}{y_i}.$$

- Lineární regrese, např. `numpy.polyfit()`
- Riziko: I když  $y_i \in N(\lambda(x_i), \sigma_i)$ , po transformaci už to platit nemusí.



# Metoda nejmenších čtverců – chyby obou proměnných

- $x$  i  $y$  jsou náhodné proměnné.
- Vzdálenost bodu  $[x_i, y_i]$  od modelové funkce, vážená standardními odchylkami:

$$d_i^2 = \frac{(x_i - \tilde{x})^2}{\sigma_{x_i}^2} + \frac{(y_i - f(\tilde{x}))^2}{\sigma_{y_i}^2}$$

- Vzdálenost bude minimální, pokud

$$\frac{\partial d_i^2}{\partial \tilde{x}} = 0$$

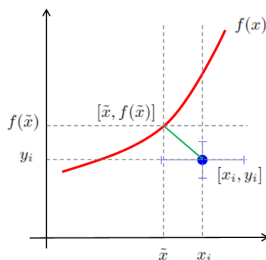
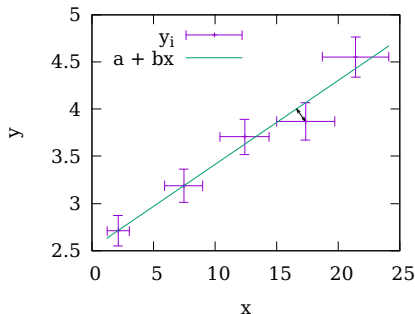
- Odtud můžeme odvodit

$$d_i^2 = \frac{(y_i - f(x_i))^2}{\sigma_{y_i}^2 + \sigma_{x_i}^2 f'^2(x_i)} = \frac{(y_i - f(x_i))^2}{\sigma_{\text{tot.}}^2}$$

- Nakonec  $\chi^2$  upravený pro chyby v  $y$  i  $x$ :

$$\chi^2 = \sum_{i=1}^N d_i^2.$$

- Prakticky např.:  
Root (TGraphErrors::Fit() automaticky),  
Gnuplot (fit... u 1:2:3:4 xyerrors...)



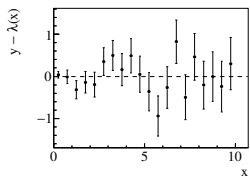
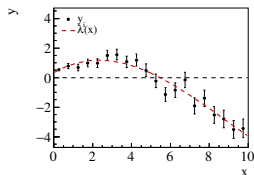
# Kvalita fitu pomocí reziduí

- Fitujeme body  $[x_i, y_i]$  funkcí  $\lambda(x|\theta)$ .
- Rezidua:  $r_i = y_i - \lambda(x_i|\theta)$ .
- Měla by být rovnoměrně rozdělená kolem nuly.
- Hledáme systematické odchyly.
- Můžeme odhalit oblasti, kde fit nepopisuje data dobře.
- Měly by mít normální rozdělení (standardizovaná rezidua  $\frac{r_i}{\sigma_{r_i}}$  normální standardní).
- Standardní odchylnka reziduí:

$$\begin{aligned}\sigma_{r_i}^2 &= \left(\frac{\partial r_i}{\partial y_i} \sigma_i\right)^2 + \left(\frac{\partial r_i}{\partial \lambda} \sigma_\lambda\right)^2 + 2 \frac{\partial r_i}{\partial y_i} \frac{\partial r_i}{\partial \lambda} \text{cov}(y_i, \lambda(x_i|\theta)) \\ &= \sigma_i^2 + \sigma_\lambda^2 - 2 \text{cov}(y_i, \lambda(x_i|\theta))\end{aligned}$$

- $\sigma_{r_i} \approx \sigma_i$ , pokud  $\sigma_i \gg \sigma_\lambda$  (např. je-li bodů hodně).
- Pro vážený průměr  $\lambda(x_i|\theta) = \bar{y}$ :

$$\text{cov}(y_i, \bar{y}) = \sigma_{\bar{y}}^2 \quad \rightarrow \quad \sigma_{r_i}^2 = \sigma_i^2 - \sigma_{\bar{y}}^2.$$

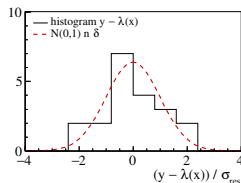
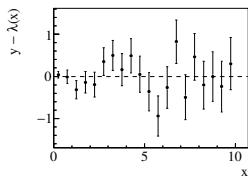
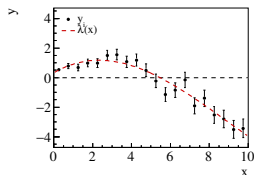


- Fitujeme body  $[x_i, y_i]$  funkcí  $\lambda(x|\boldsymbol{\theta})$ .
- Rezidua:  $r_i = y_i - \lambda(x_i|\boldsymbol{\theta})$ .
- Měla by být rovnoměrně rozdělená kolem nuly.
- Hledáme systematické odchyly.
- Můžeme odhalit oblasti, kde fit nepopisuje data dobře.
- Měly by mít normální rozdělení (standardizovaná rezidua  $\frac{r_i}{\sigma_{r_i}}$  normální standardní).
- Standardní odchylna reziduí:

$$\begin{aligned}\sigma_{r_i}^2 &= \left( \frac{\partial r_i}{\partial y_i} \sigma_i \right)^2 + \left( \frac{\partial r_i}{\partial \lambda} \sigma_\lambda \right)^2 + 2 \frac{\partial r_i}{\partial y_i} \frac{\partial r_i}{\partial \lambda} \text{cov}(y_i, \lambda(x_i|\boldsymbol{\theta})) \\ &= \sigma_i^2 + \sigma_\lambda^2 - 2 \text{cov}(y_i, \lambda(x_i|\boldsymbol{\theta}))\end{aligned}$$

- $\sigma_{r_i} \approx \sigma_i$ , pokud  $\sigma_i \gg \sigma_\lambda$  (např. je-li bodů hodně).
- Pro vážený průměr  $\lambda(x_i|\boldsymbol{\theta}) = \bar{y}$ :

$$\text{cov}(y_i, \bar{y}) = \sigma_{\bar{y}}^2 \quad \rightarrow \quad \sigma_{r_i}^2 = \sigma_i^2 - \sigma_{\bar{y}}^2.$$





- Náhodné proměnné  $z_i \in N(0, 1)$ .
- Náhodná proměnná  $y = \sum_{i=1}^N z_i^2$ .
- Pak  $y \in \chi^2(N)$ , tj. má rozdělení  $\chi^2$  s  $N$  stupni volnosti.
- Hustota pravděpodobnosti:

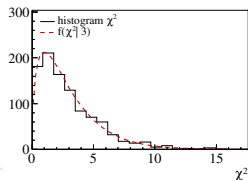
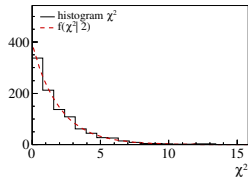
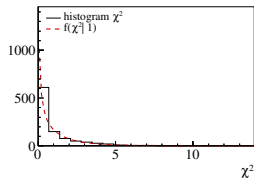
$$f_{\chi^2}(y|N) = \frac{1}{2^{\frac{N}{2}} \Gamma(\frac{N}{2})} y^{\frac{N}{2}-1} e^{-\frac{y}{2}}, \quad y \in [0, \infty), \quad N = 1, 2, \dots$$

- Gamma funkce: definovaná integrálem

$$\Gamma(x) = \int_0^{\infty} e^{-t} t^{x-1} dt.$$

$$\Gamma\left(\frac{1}{2}\right) = \sqrt{\pi}, \quad \Gamma(x+1) = x\Gamma(x), \quad \Gamma(n) = (n-1)!$$

- Momenty rozdělení:  $E[y] = N$ ,  $V[y] = 2N$ .



- Náhodné proměnné  $z_i \in N(0, 1)$ .
- Náhodná proměnná  $y = \sum_{i=1}^N z_i^2$ .
- Pak  $y \in \chi^2(N)$ , tj. má rozdělení  $\chi^2$  s  $N$  stupni volnosti.
- Hustota pravděpodobnosti:

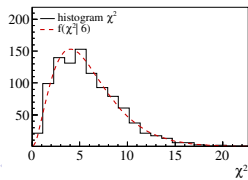
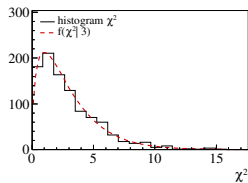
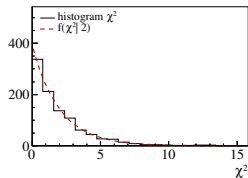
$$f_{\chi^2}(y|N) = \frac{1}{2^{\frac{N}{2}} \Gamma(\frac{N}{2})} y^{\frac{N}{2}-1} e^{-\frac{y}{2}}, \quad y \in [0, \infty), \quad N = 1, 2, \dots$$

- Gamma funkce: definovaná integrálem

$$\Gamma(x) = \int_0^{\infty} e^{-t} t^{x-1} dt.$$

$$\Gamma\left(\frac{1}{2}\right) = \sqrt{\pi}, \quad \Gamma(x+1) = x\Gamma(x), \quad \Gamma(n) = (n-1)!$$

- Momenty rozdělení:  $E[y] = N$ ,  $V[y] = 2N$ .



- Náhodné proměnné  $z_i \in N(0, 1)$ .
- Náhodná proměnná  $y = \sum_{i=1}^N z_i^2$ .
- Pak  $y \in \chi^2(N)$ , tj. má rozdělení  $\chi^2$  s  $N$  stupni volnosti.
- Hustota pravděpodobnosti:

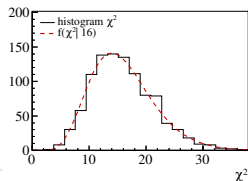
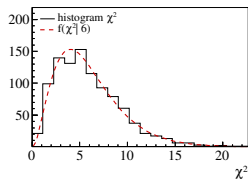
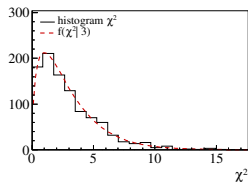
$$f_{\chi^2}(y|N) = \frac{1}{2^{\frac{N}{2}} \Gamma(\frac{N}{2})} y^{\frac{N}{2}-1} e^{-\frac{y}{2}}, \quad y \in [0, \infty), \quad N = 1, 2, \dots$$

- Gamma funkce: definovaná integrálem

$$\Gamma(x) = \int_0^{\infty} e^{-t} t^{x-1} dt.$$

$$\Gamma\left(\frac{1}{2}\right) = \sqrt{\pi}, \quad \Gamma(x+1) = x\Gamma(x), \quad \Gamma(n) = (n-1)!$$

- Momenty rozdělení:  $E[y] = N$ ,  $V[y] = 2N$ .



- Náhodné proměnné  $z_i \in N(0, 1)$ .
- Náhodná proměnná  $y = \sum_{i=1}^N z_i^2$ .
- Pak  $y \in \chi^2(N)$ , tj. má rozdělení  $\chi^2$  s  $N$  stupni volnosti.
- Hustota pravděpodobnosti:

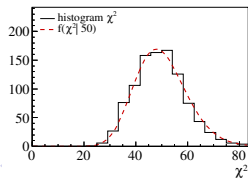
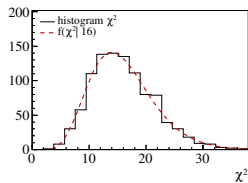
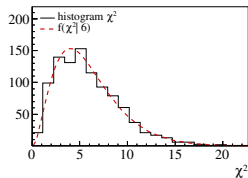
$$f_{\chi^2}(y|N) = \frac{1}{2^{\frac{N}{2}} \Gamma(\frac{N}{2})} y^{\frac{N}{2}-1} e^{-\frac{y}{2}}, \quad y \in [0, \infty), \quad N = 1, 2, \dots$$

- Gamma funkce: definovaná integrálem

$$\Gamma(x) = \int_0^{\infty} e^{-t} t^{x-1} dt.$$

$$\Gamma\left(\frac{1}{2}\right) = \sqrt{\pi}, \quad \Gamma(x+1) = x\Gamma(x), \quad \Gamma(n) = (n-1)!$$

- Momenty rozdělení:  $E[y] = N$ ,  $V[y] = 2N$ .



# Test kvality fitu pomocí $\chi^2$

- Nezávislá měření  $y_1, y_2, \dots, y_N$ .
- $y_i \in N(\lambda(x_i|\boldsymbol{\theta}), \sigma_i)$
- $\lambda$  modelová funkce s parametry  $\theta_1, \theta_2, \dots, \theta_M$ .
- Potom

$$\chi^2(\boldsymbol{\theta}|\mathbf{x}, \mathbf{y}, \boldsymbol{\sigma}) = \sum_{i=1}^N \frac{[y_i - \lambda(x_i|\boldsymbol{\theta})]^2}{\sigma_i^2}$$

má rozdělení  $\chi^2(N - M)$ .

- Počet stupňů volnosti je  $N - M$ , tedy počet bodů mínus počet fitovaných parametrů.

$$E[\chi^2] = N - M, \quad V[\chi^2] = 2(N - M).$$

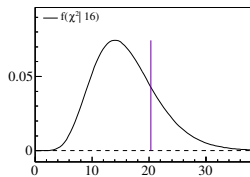
- $\chi^2$  na počet stupňů volnosti

$$E\left[\frac{\chi^2}{N - M}\right] = 1.$$

- Pro jeden fit spočítáme jeho  $\chi_0^2$  a jeho p-hodnotu:

$$P(y \geq \chi_0^2) = \int_{\chi_0^2}^{\infty} f_{\chi^2}(y|N - M)dy = 1 - F_{\chi^2}(\chi_0^2|N - M)$$

- Je-li  $P(y \geq \chi_0^2) < \alpha$ , zamítneme nulovou hypotézu, tj. že data lze popsat funkcí  $\lambda(x, \boldsymbol{\theta})$ .
- $\alpha$  je hladina signifikance, obvykle 0.05 nebo 0.01.



$\chi_0^2$  spočítaný z jedné sady nafitovaných dat, porovnaný s hustotou pravděpodobnosti  $\chi^2$ . P-hodnota odpovídá integrálu od vertikální čáry doprava.

# Test kvality fitu pomocí $\chi^2$

- Nezávislá měření  $y_1, y_2, \dots, y_N$ .
- $y_i \in N(\lambda(x_i|\boldsymbol{\theta}), \sigma_i)$
- $\lambda$  modelová funkce s parametry  $\theta_1, \theta_2, \dots, \theta_M$ .
- Potom

$$\chi^2(\boldsymbol{\theta}|\mathbf{x}, \mathbf{y}, \boldsymbol{\sigma}) = \sum_{i=1}^N \frac{[y_i - \lambda(x_i|\boldsymbol{\theta})]^2}{\sigma_i^2}$$

má rozdělení  $\chi^2(N - M)$ .

- Počet stupňů volnosti je  $N - M$ , tedy počet bodů mínus počet fitovaných parametrů.

$$E[\chi^2] = N - M, \quad V[\chi^2] = 2(N - M).$$

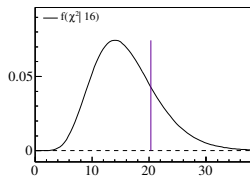
- $\chi^2$  na počet stupňů volnosti

$$E\left[\frac{\chi^2}{N - M}\right] = 1.$$

- Pro jeden fit spočítáme jeho  $\chi_0^2$  a jeho p-hodnotu:

$$P(y \geq \chi_0^2) = \int_{\chi_0^2}^{\infty} f_{\chi^2}(y|N - M)dy = 1 - F_{\chi^2}(\chi_0^2|N - M)$$

- Je-li  $P(y \geq \chi_0^2) < \alpha$ , zamítneme nulovou hypotézu, tj. že data lze popsat funkcí  $\lambda(x, \boldsymbol{\theta})$ .
- $\alpha$  je hladina signifikance, obvykle 0.05 nebo 0.01.



$\chi_0^2$  spočítaný z jedné sady nafitovaných dat, porovnaný s hustotou pravděpodobnosti  $\chi^2$ . P-hodnota odpovídá integrálu od vertikální čáry doprava.

# Test kvality fitu pomocí $\chi^2$

- Programy pro zpracování dat poskytují  $F_{\chi^2}$  nebo přímo p-hodnotu:
  - gnuplot automaticky vypisuje p-hodnotu při fitu.
  - $F_{\chi^2, r} = \text{Excel: CHISQ.DIST}(x, r, \text{true})$ , Root: `ROOT::Math::chisquared_cdf(x, r)`, Python: `stats.chi2.cdf(x, r)`...
- Lze zjistit i jaký  $\chi^2$  odpovídá určité hladině signifikance:
  - $\chi^2 = F^{-1}(P, r) = \text{Excel: CHISQ.INV}(P, r)$ , Root: `ROOT::Math::chisquared_quantile(P, r)`...
- Tabulka p-hodnot pro počet stupňů volnosti 1–10:

Počet stupňů volnosti	$\chi^2$										
	hladina signifikance $\alpha = 5\%$ $\alpha = 1\%$										
1	0.004	0.02	0.06	0.15	0.46	1.07	1.64	2.71	3.84	6.64	10.83
2	0.10	0.21	0.45	0.71	1.39	2.41	3.22	4.60	5.99	9.21	13.82
3	0.35	0.58	1.01	1.42	2.37	3.66	4.64	6.25	7.82	11.34	16.27
4	0.71	1.06	1.65	2.20	3.36	4.88	5.99	7.78	9.49	13.28	18.47
5	1.14	1.61	2.34	3.00	4.35	6.06	7.29	9.24	11.07	15.09	20.52
6	1.63	2.20	3.07	3.83	5.35	7.23	8.56	10.64	12.59	16.81	22.46
7	2.17	2.83	3.82	4.67	6.35	8.38	9.80	12.02	14.07	18.48	24.32
8	2.73	3.49	4.59	5.53	7.34	9.52	11.03	13.36	15.51	20.09	26.12
9	3.32	4.17	5.38	6.39	8.34	10.66	12.24	14.68	16.92	21.67	27.88
10	3.94	4.87	6.18	7.27	9.34	11.78	13.44	15.99	18.31	23.21	29.59
$P[y \geq \chi^2]$	0.95	0.90	0.80	0.70	0.50	0.30	0.20	0.10	0.05	0.01	0.001

- Pro počet stupňů volnosti  $> 10$  konverguje k  $N(k, \sqrt{2k})$ .