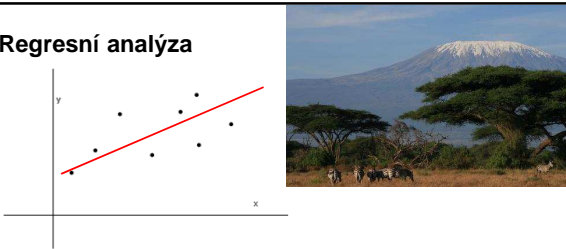



**MASARYKOVA UNIVERZITA**  
**Z1069 Statistické metody a zpracování dat**  
**VII. Regresní počít**

  
 INVESTICE DO ROZVOJE VZDĚLÁVÁNÍ

### Regresní analýza



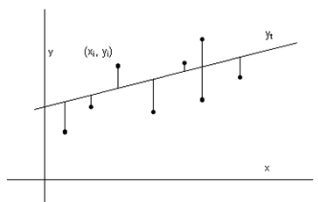
Úkolem regresní analýzy je sestavit **vztah (model)** závislosti mezi závisle a nezávisle proměnnou.

**Regresní analýza řeší :**

- odhady neznámých parametrů regresní funkce
- testování hypotéz o těchto parametrech
- ověřování předpokladů regresního modelu

### Určení lineární regresní závislosti

Nejjednodušším případem regresní závislosti je případ, kdy regresní funkce je přímkou. Rovnice regresní přímky má tvar:

$$y' = a + bx$$


Symbol  $y'$  se používá pro označení **nejpravděpodobnější teoretické hodnoty**  $y$  odpovídající danému  $x$ , která leží na regresní přímce a která se odlišuje od konkrétních hodnot  $y_i$ , které se nacházejí mimo ni.

### Metoda nejmenších čtverců

Průběh regresní přímky je určen tzv. **metodou nejmenších čtverců**, kdy musí být splněna podmínka takového průběhu přímky, při kterém je součet čtverců vzdáleností všech bodů pole od přímky minimální, tedy platí:

$$\sum (y_i - y'_i)^2 = \min$$

Výpočet vertikální vzdálenosti bodů korelačního pole od regresní přímky se provádí podle uvedeného obrázku. Z něho je zřejmé, že pro vzdálenost konkrétní hodnoty závisle proměnné  $y_i$  od bodu regresní přímky  $y'_i$  musí platit:

$$y_i - y'_i = y_i - (a + bx_i) = y_i - a - bx_i$$

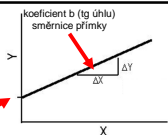
Součet čtverců svislých vzdáleností  $y_i$  od regresní přímky je potom:

$$\sum (y_i - y'_i)^2 = \sum (y_i - a - bx_i)^2 = A$$

Pro MNČ musí platit

$$A = \sum (y_i - a - bx_i)^2 = \min$$

### Výpočet koeficientů regresní přímky $y' = a + bx$



Z výše uvedených vztahů lze následnými úpravami odžít výrazy pro výpočet koeficientů regresní přímky  $a, b$

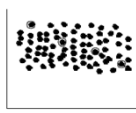
$$b = \frac{\sum x_i y_i - n \bar{x} \bar{y}}{\sum x_i^2 - n \bar{x}^2} \quad \text{případně} \quad b = \frac{s_{xy}}{s_x^2} \quad a = \bar{y} - b \bar{x}$$

Kovariance  $s_{xy}$  lomeno druhá mocnina směrodatné odchylky  $s_x$

**Koeficient b** (angl. slope) se označuje jako **koeficient regrese** a je směrnicí regresní přímky (tangentou úhlu, který přímka a svírá s osou  $x$ ). Je-li  $b > 0$ , mluvíme o regresi pozitivní, je-li  $b < 0$  o regresi negativní.

**Koeficient a** (angl. intercept) představuje  $y$ -ovou souřadnici průsečíku regresní přímky s osou  $y$  (tedy při  $x=0$ ).

### Intervaly a pásy spolehlivosti lineární regresní závislosti



Koeficienty (parametry) regresní přímky jsou **bodovými odhady** !

- Konstrukci regresní přímky provádíme na základě **výběrových souborů**.
- Proto se její rovnice může u různých výběrů ze stejných základních souborů lišit.
- Z tohoto důvodu je vhodné doplnit průběh regresní přímky také tzv. **intervaly spolehlivosti**.
- Výpočtem intervalů spolehlivosti určujeme pro vybraný  $x$  interval, v němž se mohou s určitou pravděpodobností vyskytovat hodnoty  $y$  s tím, že jejich nejrepresentativnější hodnota je  $y'$ .

## Intervaly a pásy spolehlivosti

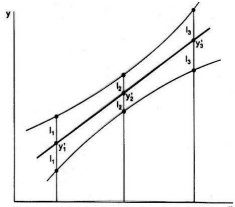
Nejprve je zapotřebí zvolit hladinu spolehlivosti – tedy pravděpodobnost, s níž očekáváme výskyt hodnot  $y$  v určených mezích  $1-p$  ( $p=0,05$  či  $0,01$ ). Poloviční šířka intervalu spolehlivosti  $l$  je dána výrazem:

$$l = t_{1-p} \cdot \frac{h\sqrt{A}}{\sqrt{n-2}} \quad h = \sqrt{\frac{1}{n} + \frac{(x-\bar{x})^2}{(n-1)s_x^2}}$$

Hodnota  $t_p$  je kritická hodnota rozdělení pro  $n-2$  stupňů volnosti a hladinu významnosti  $p$ . Meze intervalů spolehlivosti určíme pomocí hodnot  $y'$  z rovnice  $y' - \bar{y} = b(x - \bar{x})$

horní mez:  $y' + l$

dolní mez:  $y' - l$

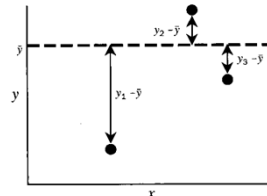


Pásky spolehlivosti vzniknou spojením krajních bodů intervalů spolehlivosti.

## Testování vhodnosti regresní závislosti

• Nejčastěji se k testování používá **analýza rozptylu (ANOVA)**.

• **Princip:** Zjistíme celkovou proměnlivost hodnot  $y$  a následně vypočteme, z jaké části je tato celková variabilita objasněna proměnlivostí v hodnotách  $x$ .

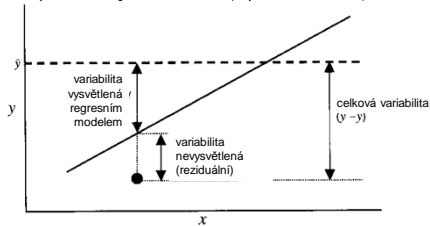


**Celková variabilita:** celková suma čtverců: od každé hodnoty  $y$  odečteme průměr, výsledek povýšíme na druhou a sečteme pro všechna  $y$ .

## Testování významnosti regresní závislosti

Celkovou variabilitu lze rozdělit na dvě části:

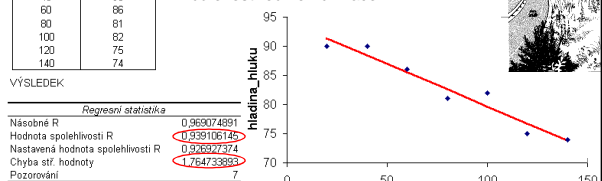
- variabilitu **vysvětlenou** regresní čarou
- variabilita **nevysvětlená** regresním modelem (zbytkovou, reziduální)



- testování je založeno na porovnání množství vysvětlené a nevysvětlené variability
- je tedy **obdobou F-testu**: testujeme, zda variabilita vysvětlená modelem se významně liší od variability reziduální, tedy: **H0: nelíší se**
- Vypočte se testovací kritérium ( $F$ ) a pokud jemu příslušející  $p$ -hodnota je menší než alfa ( $0,05$ ) potom **zamítáme nulovou hypotézu** a konstatujeme, že **regresní model je vhodný**

## Příklad regresní analýzy v EXCELU

Zjistěte, jak souvisí hladina hluku se vzdáleností od komunikace.

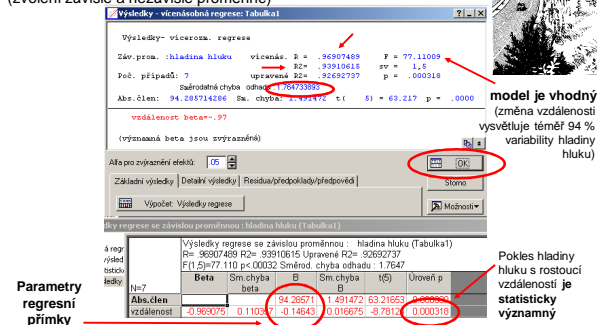


ANOVA		model je vhodný		vzdálenost	
Regrese	Rezidua	SS	MS	F	Významnost F
1	5	240,1428571	240,1429	77,11009	0,000317688
5	6	15,57142857	3,114286		
Celkem		6	255,7142857		

**Existuje signifikantní pokles hladiny hluku se vzdáleností od komunikace. Lineární regresní model vysvětluje 93,9% variability hodnot hladiny hluku.**

## Řešení v programu Statistica

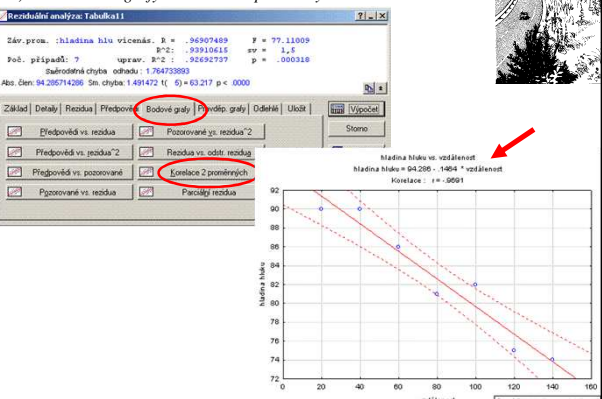
1) Statistika – Vícezměrná regrese (zvolení závisle a nezávisle proměnné)



U regresního modelu se často testuje, zda se **směrnice přímky významně liší od nuly**. Používá se **t-testu**, kterým lze prokázat statisticky významný růst či pokles hodnot závisle proměnné při změně hodnot nezávisle proměnné.

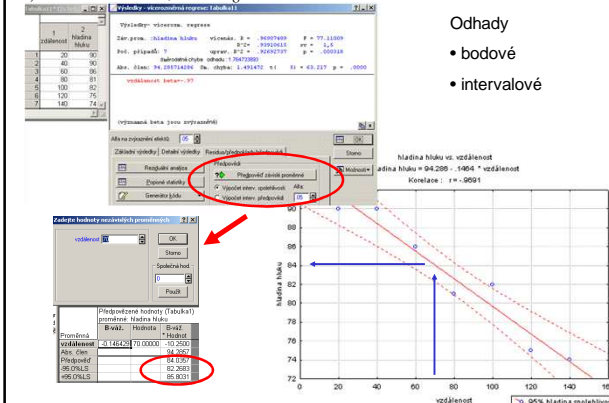
## Řešení v programu Statistica

2) OK – Bodové grafy – Korelace 2 proměnných



### Výpočet neznámé hodnoty - předpovědi

#### 3) Statistika – Vícerozměrná regrese

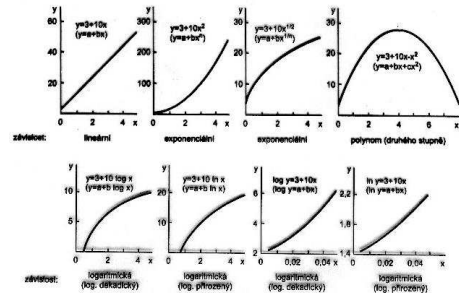


- Odhady
- bodové
  - intervalové

### Další typy regresních funkcí

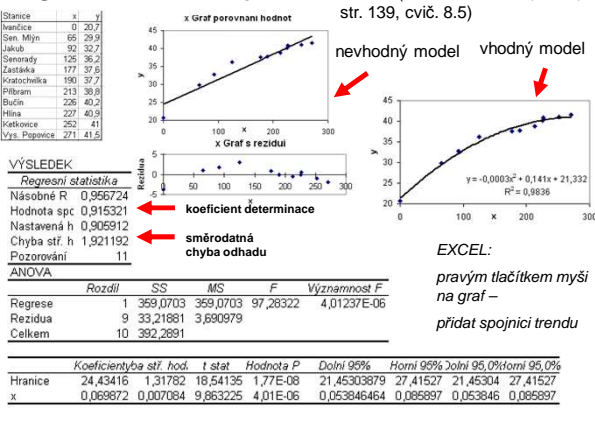
Regresní vztah dvou proměnných často nelze vhodně vyjádřit přímkou – jiné typy funkcí.

Může mít tvar např. logaritmických či exponenciálních funkcí a nebo je vztah vyjádřen rovnicí polynomu m-tého stupně.



### Regresní závislost není přímka

Příklad (viz. Brázdil a kol., 1995, str. 139, cvič. 8.5)



### Hledání vhodného regresního modelu

Lze postupovat dvěma způsoby:

1. Volba vhodného modelu na základě praktické zkušenosti či teoretických předpokladů
2. Posouzením bodového grafu a interpretací nástrojů regresní analýzy

Způsoby hodnocení vhodnosti regresního modelu

- analýza reziduálních hodnot
- výpočet směrodatné chyby odhadu ( $s_e$ )
- výpočet koeficientu determinace ( $r^2_{xy}$ ).

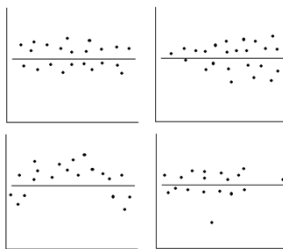
### Hledání vhodného regresního modelu

#### Analýza reziduálních hodnot

Rezidua jsou vzdálenosti skutečných hodnot  $y_i$  od modelem odhadnutých hodnot  $y_i^{\wedge}$

Zvolený regresní model považujeme za vhodný, pokud reziduální hodnoty splňují všechny následující podmínky:

- rezidua jsou **náhodná** a **nezávislá**
- mají **normální rozdělení** s nulovým průměrem a konstantním rozptylem
- rozptyl reziduí je **konstantní**.



### Hledání vhodného regresního modelu

**Směrodatná chyba odhadu** – je vyjádřením směrodatné odchylky resp. rozptýlu reziduálních hodnot a vhodnou mírou pro posouzení vhodnosti použité regresní závislosti

$$s_e = \sqrt{\frac{\sum_{j=1}^n (y_j - y_j^{\wedge})^2}{n-2}}$$

Čím je hodnota reziduálního rozptýlu nižší, tím je model vhodnější.

**Koeficient determinace** ( $r^2_{xy}$ ) – viz. Korelační počet

$$r^2 = \frac{SS_{regres}}{SS_{total}}$$

Čím je hodnota koeficientu determinace větší, tím je model vhodnější.

### Hledání vhodného regresního modelu

Grafy – Bodové grafy

volba modelu

### Vícerozměrná regrese

Popisuje závislost více proměnných z nichž více je příčinami (vysvětlující proměnné) a jen jedna je důsledek (vysvětlovaná proměnná).

Jsou-li dvě vysvětlující proměnné regresní model je rovina

Odhad parametrů se provádí MNČ

Výstupy a interpretace jsou „obdobné“ jako u modelu jednorozměrné regrese

$$y' = a + b_1x_1 + b_2x_2 + \dots$$

Např.:  
 $\text{Úhrn}_\text{sřázek} = 345,6 + 0,45 \cdot \text{zem}_\text{délka} + 1,23 \cdot \text{nadm}_\text{výška}$

### Vícerozměrná regrese

Statistika – vícerozměrná regrese (data viz. Brázdil a kol., 1995, str. 129, cvič. 8.3)

Punkva – pod Sk. Mlýnem (y) model odtoku:  
 $y = 11,0289 + 1,2487x_1 + 1,0497x_2$

↑ standardizované regresní koeficienty      ↑ regresní koeficienty

### Měření závislosti kvalitativních znaků

- Kvalitativní znaky mají slovní charakter a získáváme je v sociologických průzkumech, při terénním šetření apod.
- K charakterizování závislosti kvalitativních znaků slouží tzv. kontingenční tabulky

- Z kontingenční tabulky lze určit intenzitu závislosti ve dvojici slovních znaků.
- Máme-li dva alternativní znaky dostaneme tzv. čtyřpolní tabulku.

### Měření závislosti kvalitativních znaků

Obecně může mít každý kvalitativní znak A r tříd a znak B s tříd. Výsledky šetření potom sestavujeme do kontingenční tabulky r x s.

Pozorované četnosti v jednotlivých buňkách označujeme dvěma indexy – obecně  $n_{ij}$ .

Také marginální četnosti mají dva indexy.

Ten, přes který je sčítáno je označen hvězdičkou – tedy  $n_{2*}$  značí součet četností v druhé řádce,  $n_{*1}$  značí součet četností v prvním sloupci.

Tabulka bývá doplněna hodnotami procentuálních (relativních) četností. Častým požadavkem je konstantní délka intervalů tvořících třídy.

Stejně jako v případě kvantitativních znaků ověřujeme i zde existenci vztahu testy významnosti a hodnotíme ho vhodnou mírou závislosti.

### Kontingenční tabulka typu r x s

Tříděný znak	Znak B					Součet		
	$b_1$	$b_2$	...	$b_j$	...		$b_s$	
Znak A	$a_1$	$n_{11}$	$n_{12}$	...	$n_{1j}$	...	$n_{1s}$	$n_{1*}$
	$a_2$	$n_{21}$						$n_{2*}$
	⋮							⋮
	$a_i$				$n_{ij}$			$n_{i*}$
	⋮							⋮
	$a_r$	$n_{r1}$					$n_{rs}$	$n_{r*}$
	Součet	$n_{*1}$	$n_{*2}$	...	$n_{*j}$	...	$n_{*s}$	$n_{**} = n$

### Posuzování závislosti v kontingenčních tabulkách

Podmíněné četnosti uvnitř kontingenční tabulky mají podobný význam jako body korelačního diagramu — jejich rozmístění umožňuje usuzovat na charakter závislosti tříděných znaků.

Pro posouzení nezávislosti obou znaků můžeme vedle pozorovaných četností stanovit pro jednotlivá pole také očekávané (teoretické) četnosti :

$$n_{ij} = \frac{n_{i*} \cdot n_{*j}}{n}$$

tedy jako součin okrajových četností příslušného řádku a sloupce dělený rozsahem souboru.

Pro každé pole kontingenční tabulky existuje dvojice četností - četnost pozorovaná a četnost vypočtená.

### Hypotéza nezávislosti

Ukazatel, který pro tabulku jako celek měří rozdílnost pozorovaných a vypočtených četností v jednotlivých polích tabulky se nazývá čtvercová kontingence  $\chi^2$

$$\chi^2 = \sum_{i=1}^r \sum_{j=1}^s \frac{(n_{ij} - n_{ij}')^2}{n_{ij}'}$$

Je to bezrozměrná hodnota a platí:  $\chi^2 \geq 0$

Hodnoty nula nabývá pouze v případě, že znaky v kontingenční tabulce jsou nezávislé.

Vypočtená hodnota  $\chi^2$  se porovnává na zvolené hladině významnosti  $\alpha$  s kritickou hodnotou  $\chi^2$  rozdělení pro  $(r-1)(s-1)$  stupňů volnosti.

Hypotézu (H0) o nezávislosti dvou studovaných znaků zamítáme, jestliže vypočtená hodnota  $\chi^2$  je větší než tabulková; případně, když jí příslušející *p-hodnota* je menší než zvolená hladina významnosti.

### Příklad analýzy závislosti v tabulce r x s

Pro výběr 234 studentů zjišťujeme, zda existuje vztah mezi sportem, který provozují a sportovními pořady, které sledují v televizi.

Sestavíme tabulku typu 4 x 4:

Obľíbenost při sledování televize	Obľíbenost při sportování				Řádkové součty
	hry	atletika	gymnastika	plavání	
hry	133	6	2	4	145
atletika	15	10	4	3	32
gymnastika	4	1	25	0	30
plavání	9	0	1	17	27
<b>Sloupcové součty</b>	161	17	32	24	234

Hypotéza nezávislosti H<sub>0</sub>: Neexistuje vztah mezi provozovaným sportem a sportem sledovaným v TV.

Vypočtená hodnota testovacího kritéria  $\chi^2 = 273,3$

Kritická hodnota z tabulek pro  $p=0,05$  a  $(4-1)(4-1)=9$  stupňů volnosti:

$$\chi^2 = 16,9$$

**Závěr: H<sub>0</sub> zamítáme, existuje významný vztah.**

### Testování nezávislosti v tabulce 2 x 2

	Zájem o statistiku		řádkové součty
	ano	ne	
chlapci	a	b	a + b
dívky	c	d	c + d
<b>sloupcové součty</b>	a + c	b + d	n

Pro výpočet testovacího kritéria  $\chi^2$  v tabulce 2 x 2 můžeme využít zjednodušený vzorec:

$$\chi^2 = \frac{n(ad - bc)^2}{(a + b)(c + d)(a + c)(b + d)}$$

Protože v 2x2 tabulce můžeme uvažovat i směr poruchy nulové hypotézy – proto musíme rozhodnout, zda použijeme test jednostranný či dvoustranný.

Kritické hodnoty jsou uvedeny v tabulce  $\chi^2$  rozdělení o jednom stupni volnosti.

### Příklad analýzy závislosti v tabulce 2 x 2

	Zájem o statistiku		řádkové součty
	ano	ne	
chlapci	30	36	66
dívky	11	63	74
<b>sloupcové součty</b>	41	99	140

Hypotéza nezávislosti H<sub>0</sub>: Relativní četnost studentů se zájmem o statistiku je nezávislá na pohlaví.

Vypočtená hodnota testovacího kritéria:  $\chi^2 = \frac{140(30 \times 63 - 11 \times 36)^2}{41 \times 99 \times 66 \times 74} = 15,8$

Kritická hodnota  $\chi^2$ -rozdělení z tabulek pro  $\alpha=0,05$ : 3,84

**Závěr: H<sub>0</sub> zamítáme, existuje významný rozdíl.**

Zájem u chlapců:  $30/66 = 0,45$

Zájem u dívek:  $11/74 = 0,14$

**Chlapci mají zhruba 3x větší zájem o statistiku než dívky.**

### Čyřpolní tabulka - řešení v programu Statistica

Statistiky – Neparametrická statistika – Tabulka 2 x 2

The screenshot shows the Statistica interface with a 2x2 contingency table and its analysis results. The table data is as follows:

	Sloupec1	Sloupec2	Řádek celkem
Počet, řádek 1	30	36	66
Procent z celku	21,429%	25,714%	47,143%
Počet, řádek 2	11	63	74
Procent z celku	7,857%	45,000%	52,857%
Sloupec celkem	41	99	140
Procent z celku	29,286%	70,714%	

The analysis results shown in the software are:

- N-kvadrát (sym): 15,76 p= 0,001
- Yatesův kongovaný chi-kv: 14,32 p= 0,002
- Fikvadrát: 11,259
- Fisherovo p, jednodr.: p= 0,001
- oboustr.: p= 0,001
- McNemarův chi-kvadrát: 11,01 p= 0,009
- Chi-kvadrát: 12,26 p= 0,005