

# REGRESNÍ ANALÝZA

**Příklad:**

**Z dat PISA zjišťují vztah mezi:**

**matematickým skóre žáků**

**a**

**úrovní příjmů domácnosti v níž žijí**

**/pro přehlednost je vše v tomto příkladu  
vymyšleno a souvislost mezi příjmem a  
výkony žáků není zdaleka tak silná**

**OSA Y – HODNOTY DRUHÉ PROMĚNNÉ**

**Pro regresní modelování je typická představa dat v souřadné soustavě...**

**OSA X – HODNOTY JEDNÉ PROMĚNNÉ**

SKÓRE MATEMATICKÝCH DOVEDNOSTÍ  
0-100

80 bodů

Zde máme Artura: jeho rodiče jsou poměrně bohatí a má dobré výsledky

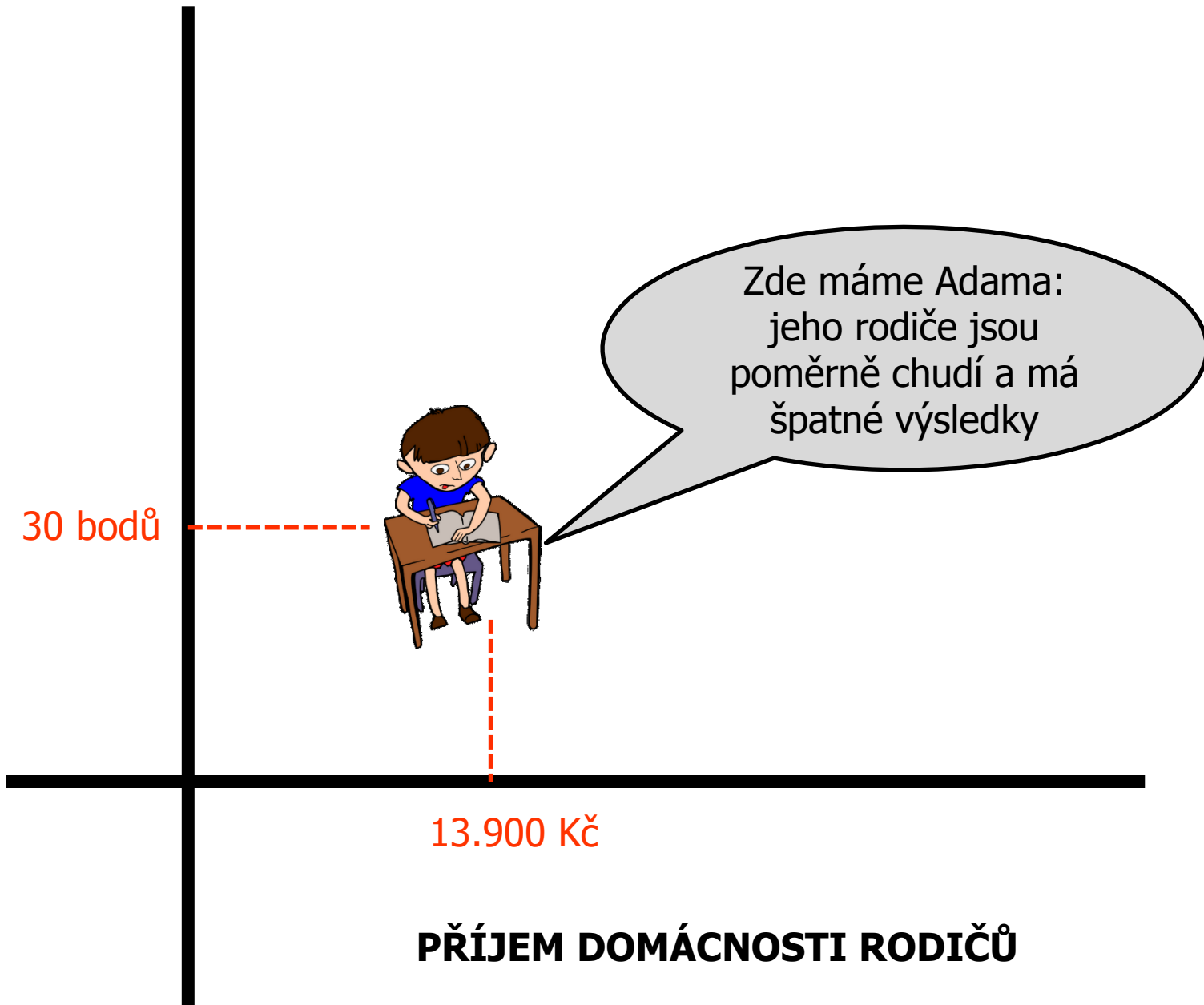


Představme si tedy, že vynášíme jednotlivé případy do bodového grafu...

68.500 Kč

PŘÍJEM DOMÁCNOSTI RODIČŮ

**SKÓRE MATEMATICKÝCH DOVEDNOSTÍ**  
**0-100**



**SKÓRE MATEMATICKÝCH DOVEDNOSTÍ**  
**0-100**

70 bodů



Zde máme Andreu: její rodiče jsou na tom průměrně a Andrea má celkem dobré výsledky

28.300 Kč

**PŘÍJEM DOMÁCNOSTI RODIČŮ**

**SKÓRE MATEMATICKÝCH DOVEDNOSTÍ**  
**0-100**



**Takto vypadají uspořádaná data za tři žáky...**

**PŘÍJEM DOMÁCNOSTI RODIČŮ**

**SKÓRE MATEMATICKÝCH DOVEDNOSTÍ**  
**0-100**

**Přidáváme další...**

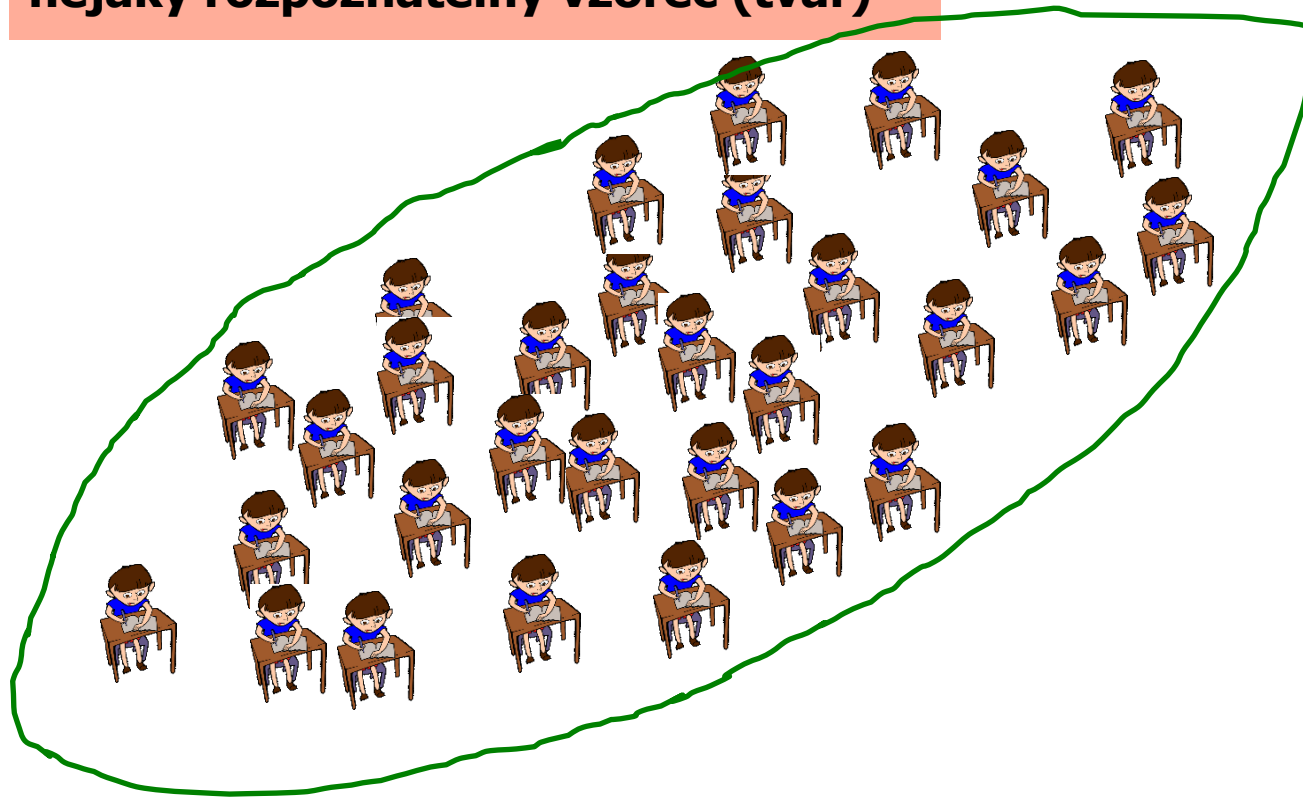


**PŘÍJEM DOMÁCNOSTI RODIČŮ**



**SKÓRE MATEMATICKÝCH DOVEDNOSTÍ**  
**0-100**

a další...  
A sledujeme, jestli v grafu vidíme  
nějaký rozpoznatelný vzorec (tvar)

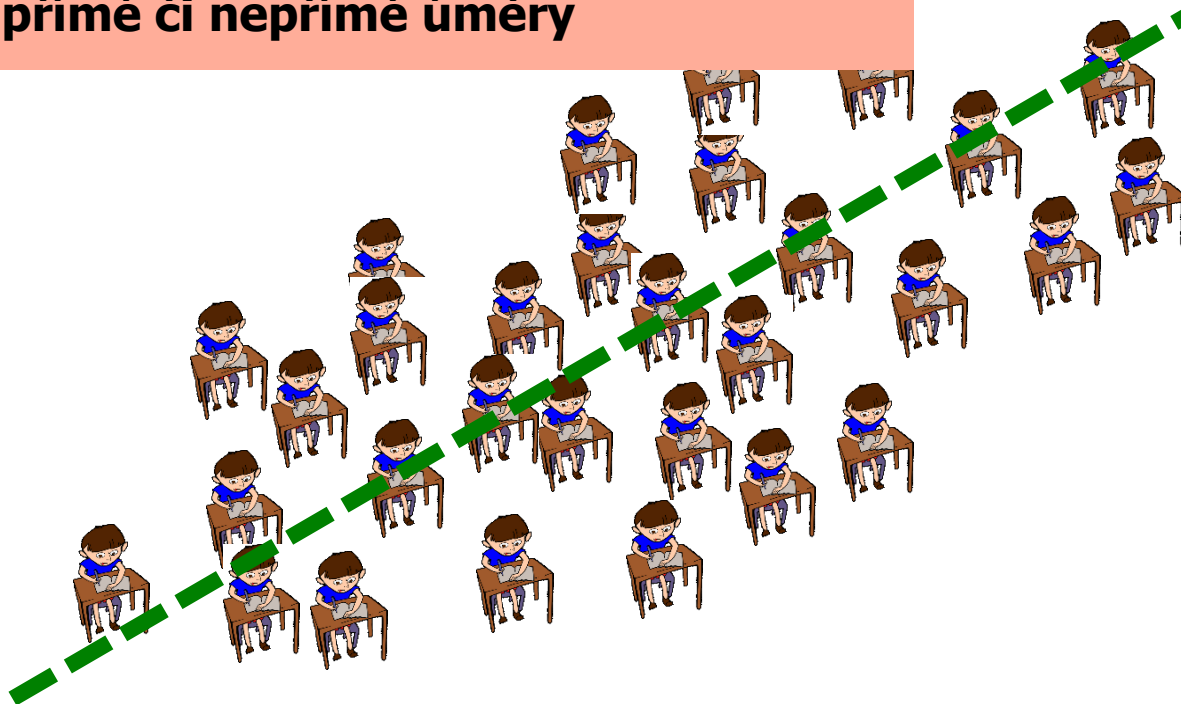


**PŘÍJEM DOMÁCNOSTI RODIČŮ**

SKÓRE MATEMATICKÝCH DOVEDNOSTÍ

0-100

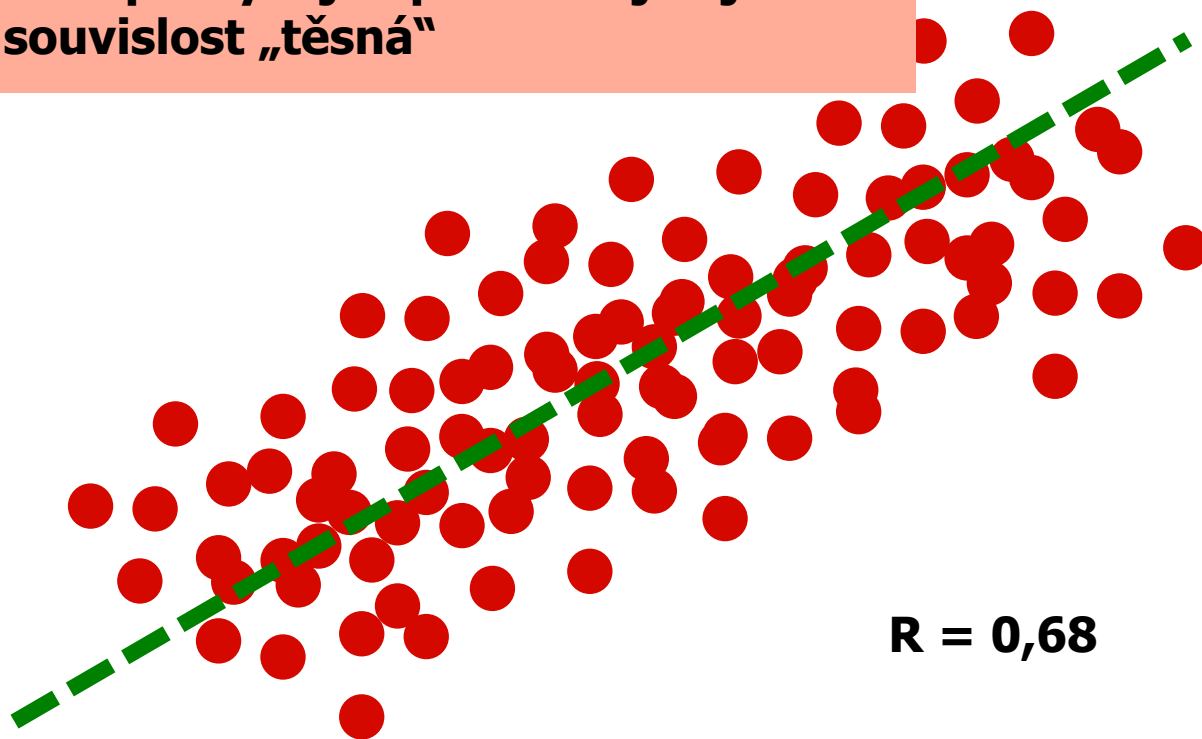
Úkolem regrese je proložit daty  
přímku – tu považujeme za vhodný a  
nejjednodušší model souvislosti  
přímé či nepřímé úměry



PŘÍJEM DOMÁCNOSTI RODIČŮ

SKÓRE MATEMATICKÝCH DOVEDNOSTÍ  
0-100

Tohle už známe z korelací –  
koeficient korelace vyjadřuje jak se  
data přimykají k přímce – jak je  
souvislost „těsná“

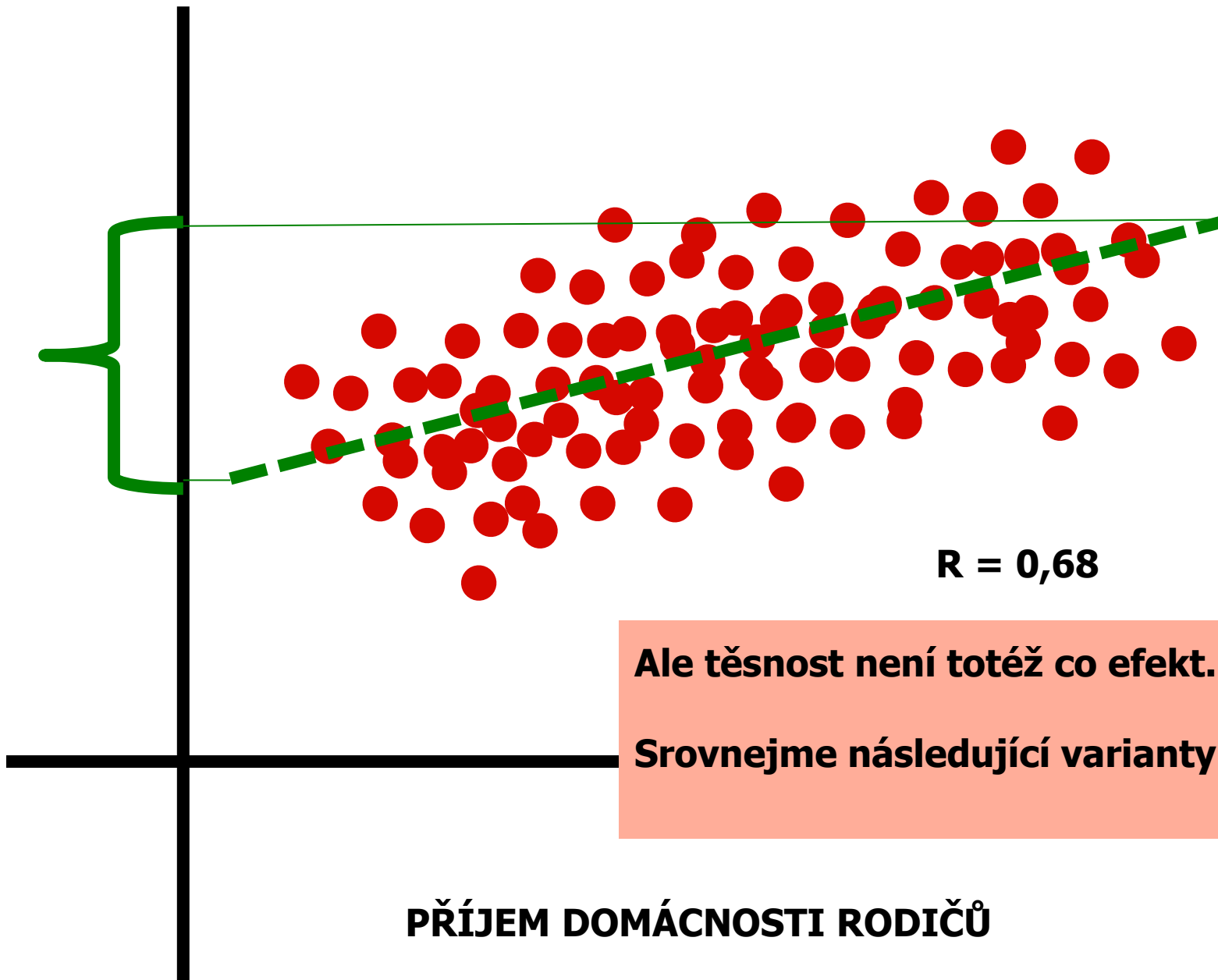


$R = 0,68$

PŘÍJEM DOMÁCNOSTI RODIČŮ

SKÓRE MATEMATICKÝCH DOVEDNOSTÍ

0-100

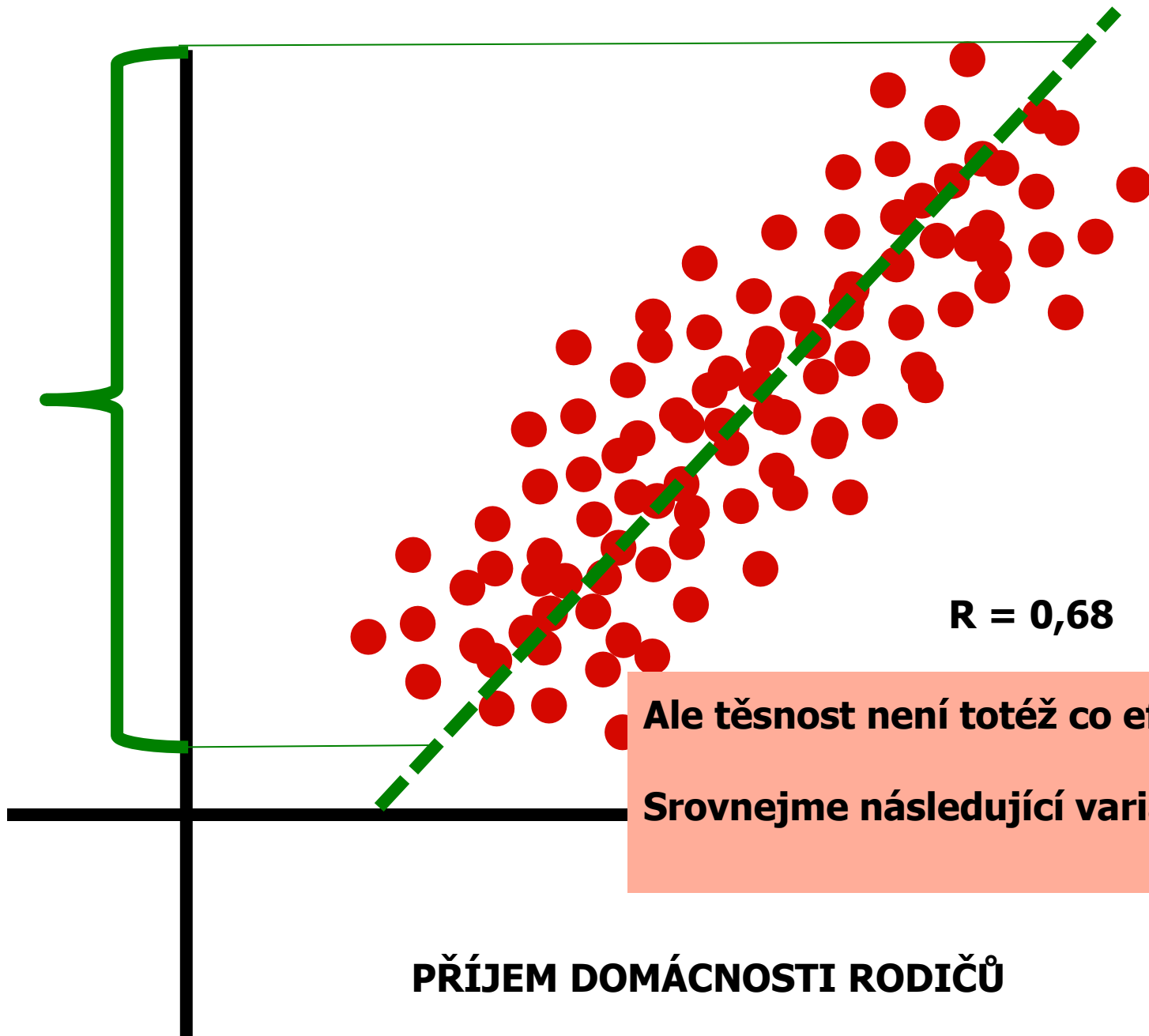


Ale těsnost není totéž co efekt...

Srovnejme následující varianty:

PŘÍJEM DOMÁCNOSTI RODIČŮ

SKÓRE MATEMATICKÝCH DOVEDNOSTÍ  
0-100



$R = 0,68$

Ale těsnost není totéž co efekt...

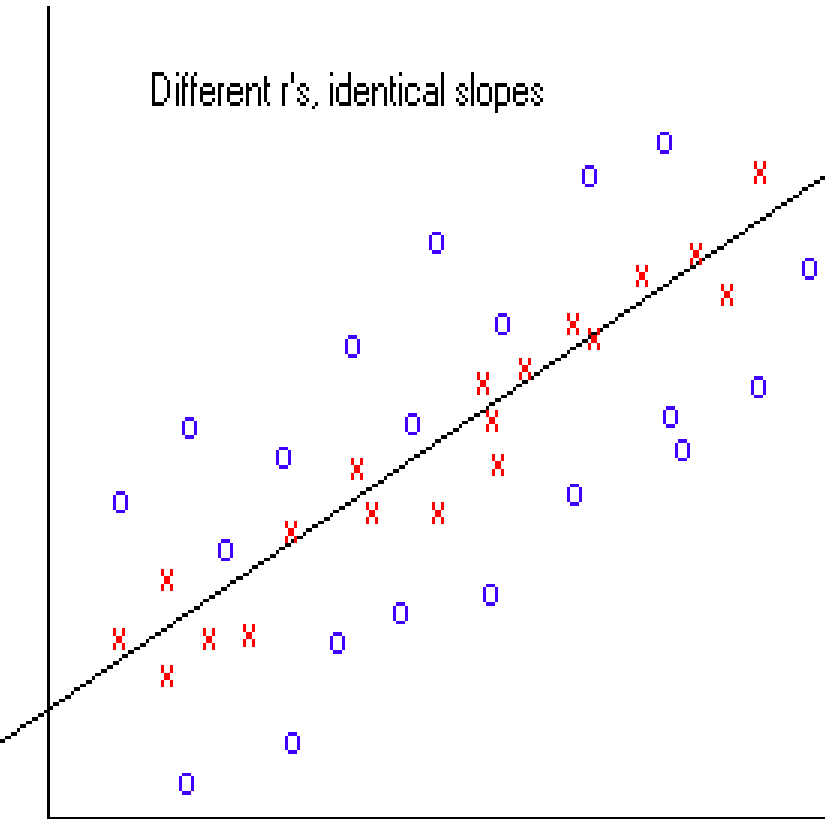
Srovnejme následující varianty:

PŘÍJEM DOMÁCNOSTI RODIČŮ

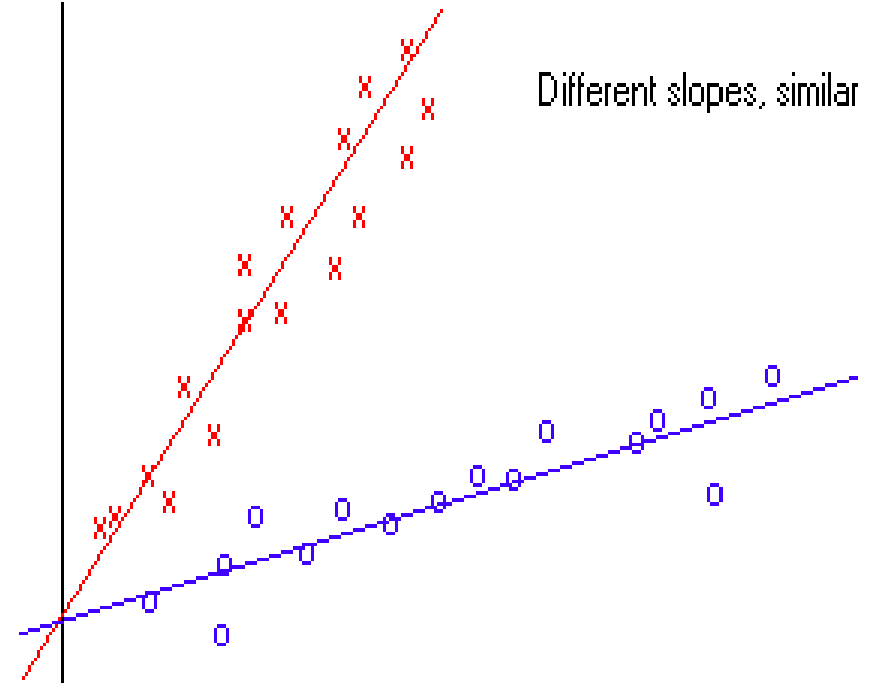
**TĚSNOST** = jak často je změna v příjmu rodičů asociována se změnou skóre žáka

**EFEKT** = jak moc se mění skóre žáka v závislosti na příjmu rodičů

Different r's, identical slopes



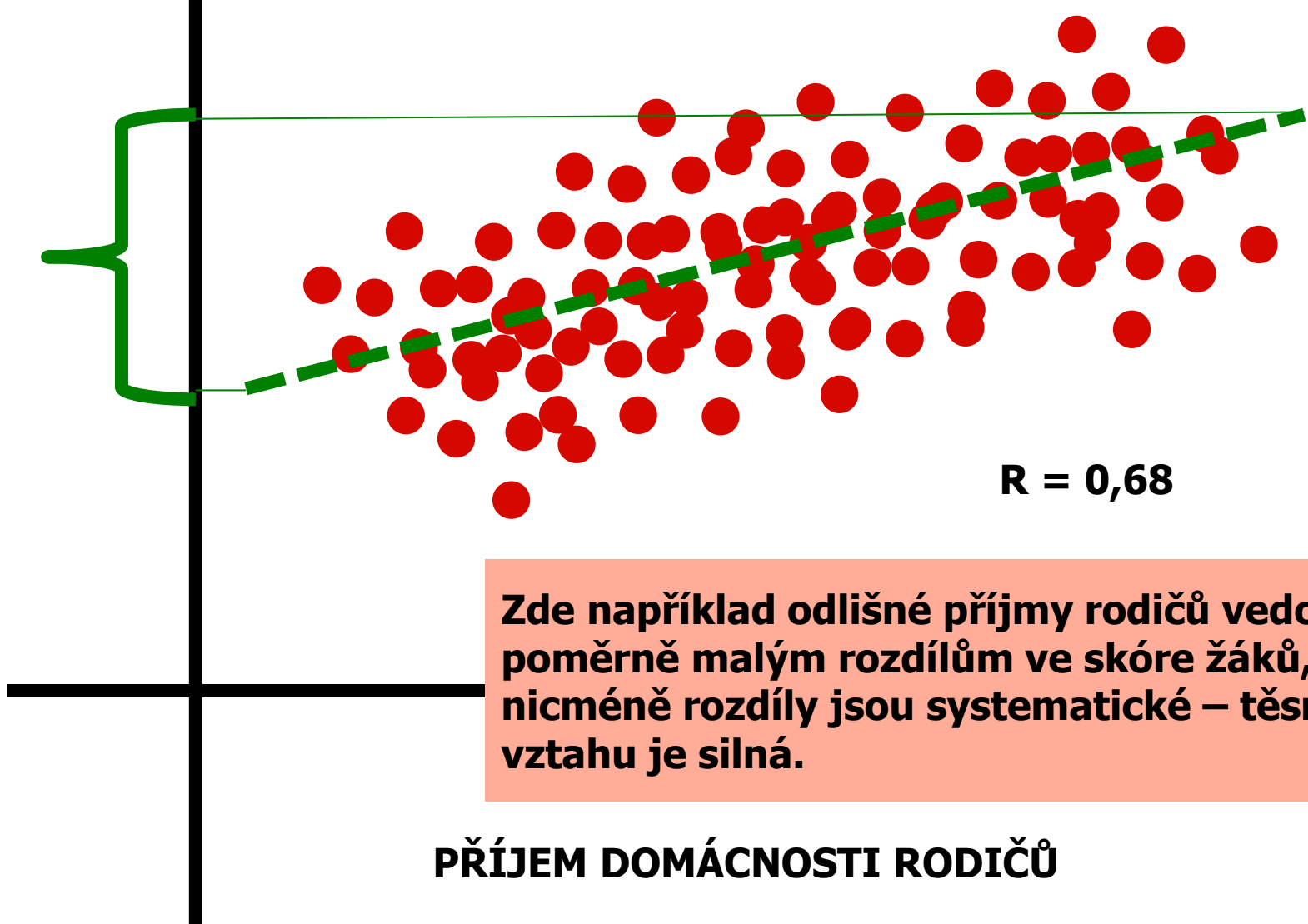
Different slopes, similar r's



SKÓRE MATEMATICKÝCH DOVEDNOSTÍ

0-100

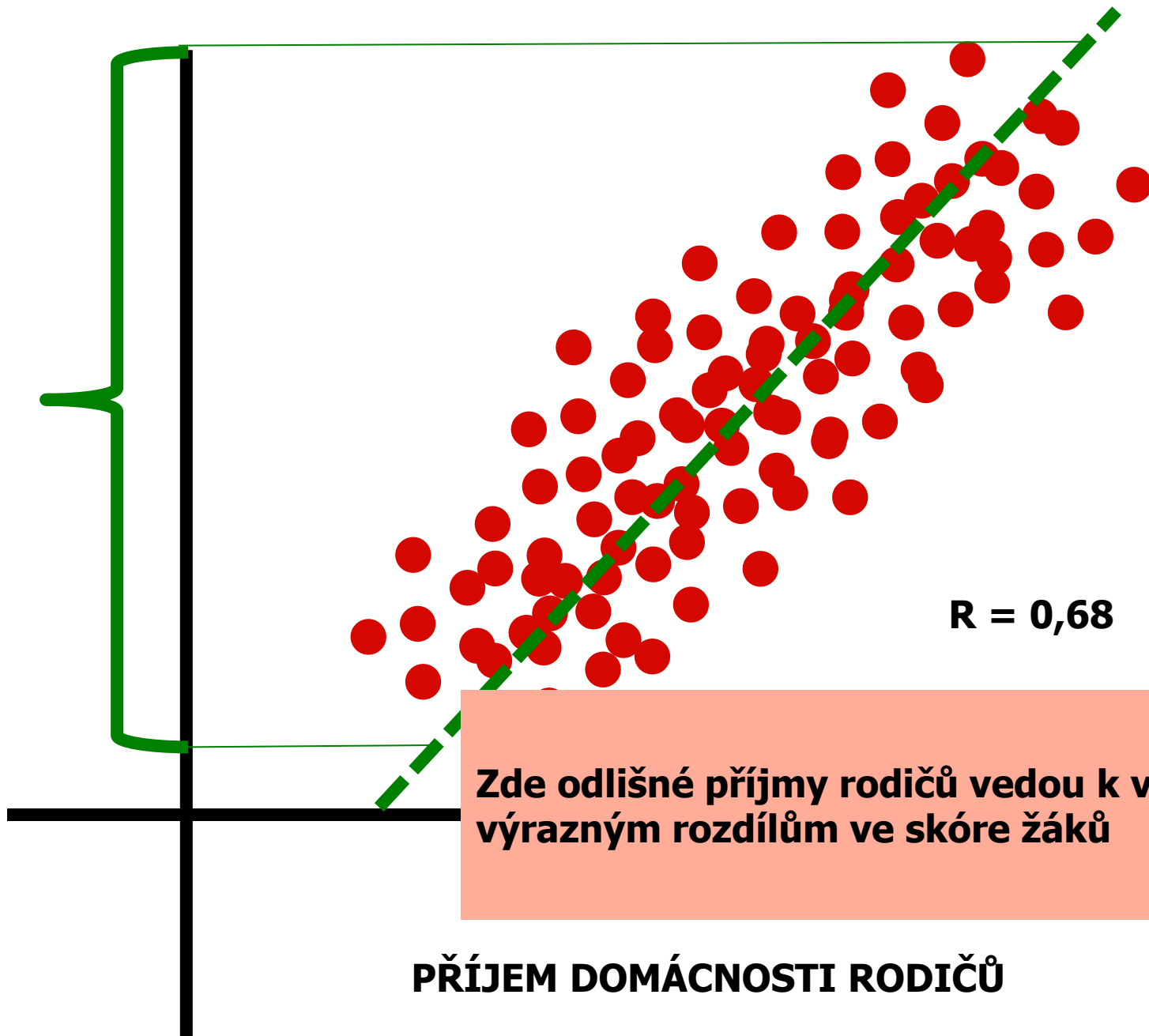
REGRESE nám tedy na rozdíl od korelace prozradí, jak velký vliv má příjem rodičů na skóre žáků



Zde například odlišné příjmy rodičů vedou k poměrně malým rozdílům ve skóre žáků, nicméně rozdíly jsou systematické – těsnost vztahu je silná.

PŘÍJEM DOMÁCNOSTI RODIČŮ

**SKÓRE MATEMATICKÝCH DOVEDNOSTÍ**  
**0-100**





SKÓRE MATEMATICKÝCH DOVEDNOSTÍ

0-100

Zároveň regrese umožňuje  
predikovat...



AHOJ, JSEM  
BEDŘICH A JSEM TU  
NOVÝ...

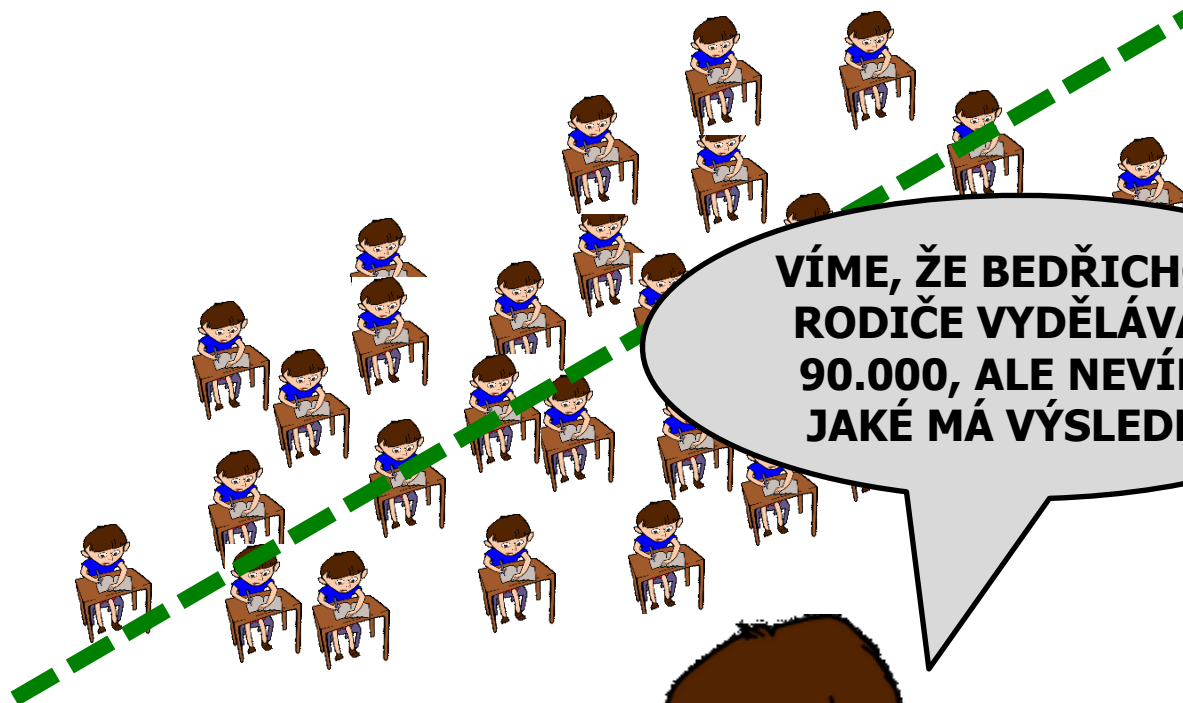
PŘÍJEM DOMÁC



SKÓRE MATEMATICKÝCH DOVEDNOSTÍ

0-100

Zároveň regrese umožňuje predikovat...



VÍME, ŽE BEDŘICHOVI  
RODIČE VYDĚLÁVÁJÍ  
90.000, ALE NEVÍME  
JAKÉ MÁ VÝSLEDKY

PŘÍJEM DOMÁCÍ



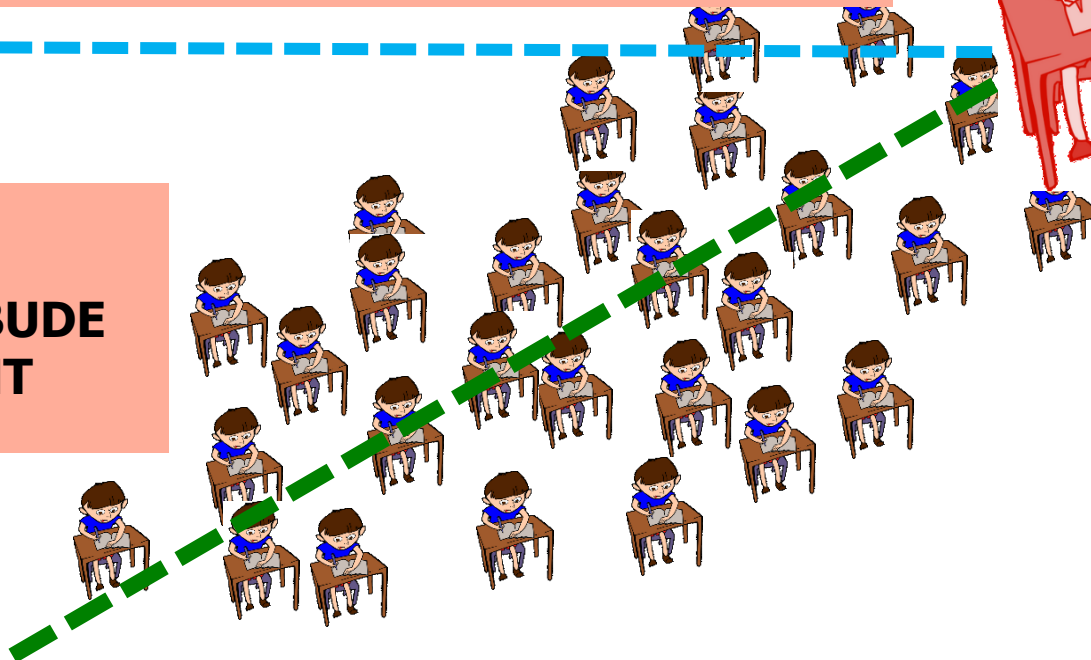
HOVĚDNOSTÍ

Podle hodnoty nezávisle proměnné predikujeme pozici na přímce a z ní odečteme hodnotu závisle proměnné...

Odhad:  
85 bodů

BEDŘICH  
PRAVDĚPODOBĚNĚ BUDE  
NÁŠ NOVÝ PREMIANT

SKÓRE MATĚ



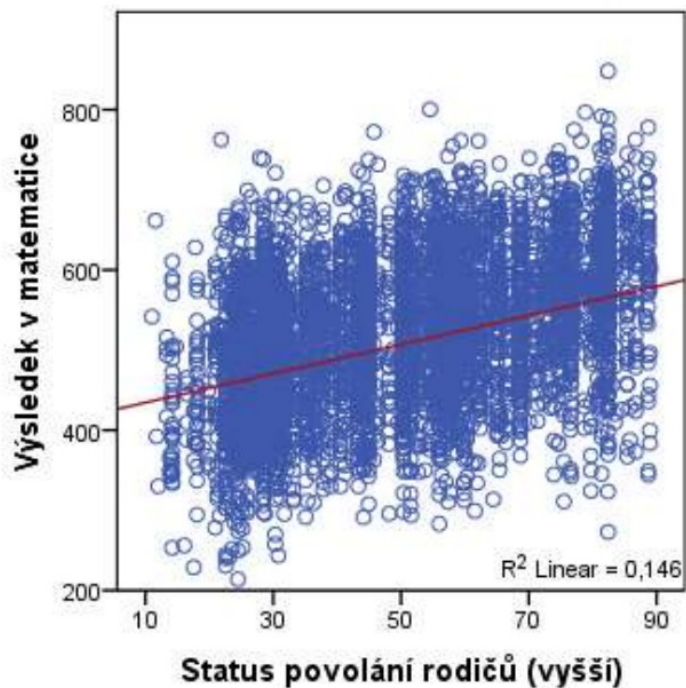
90.000 Kč

PŘÍJEM DOMÁCNOSTI RODIČŮ

## VAROVÁNÍ K PŘÍKLADU:

Omluvte prosím stigmatizující a potenciálně stereotypizující charakter přechozího příkladu – děti chudých se neučí o tolik hůře jako děti bohatých

Graf č. 4 Vztah úrovně matematické způsobilosti a společenského postavení rodin žáků



## Výsledky z reálných dat PISA 2012

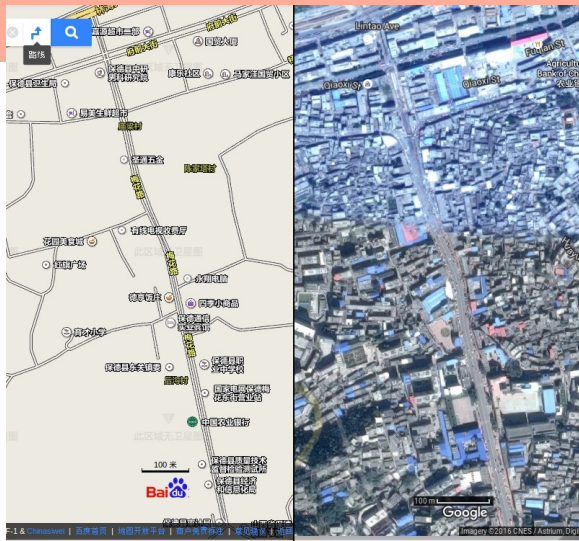
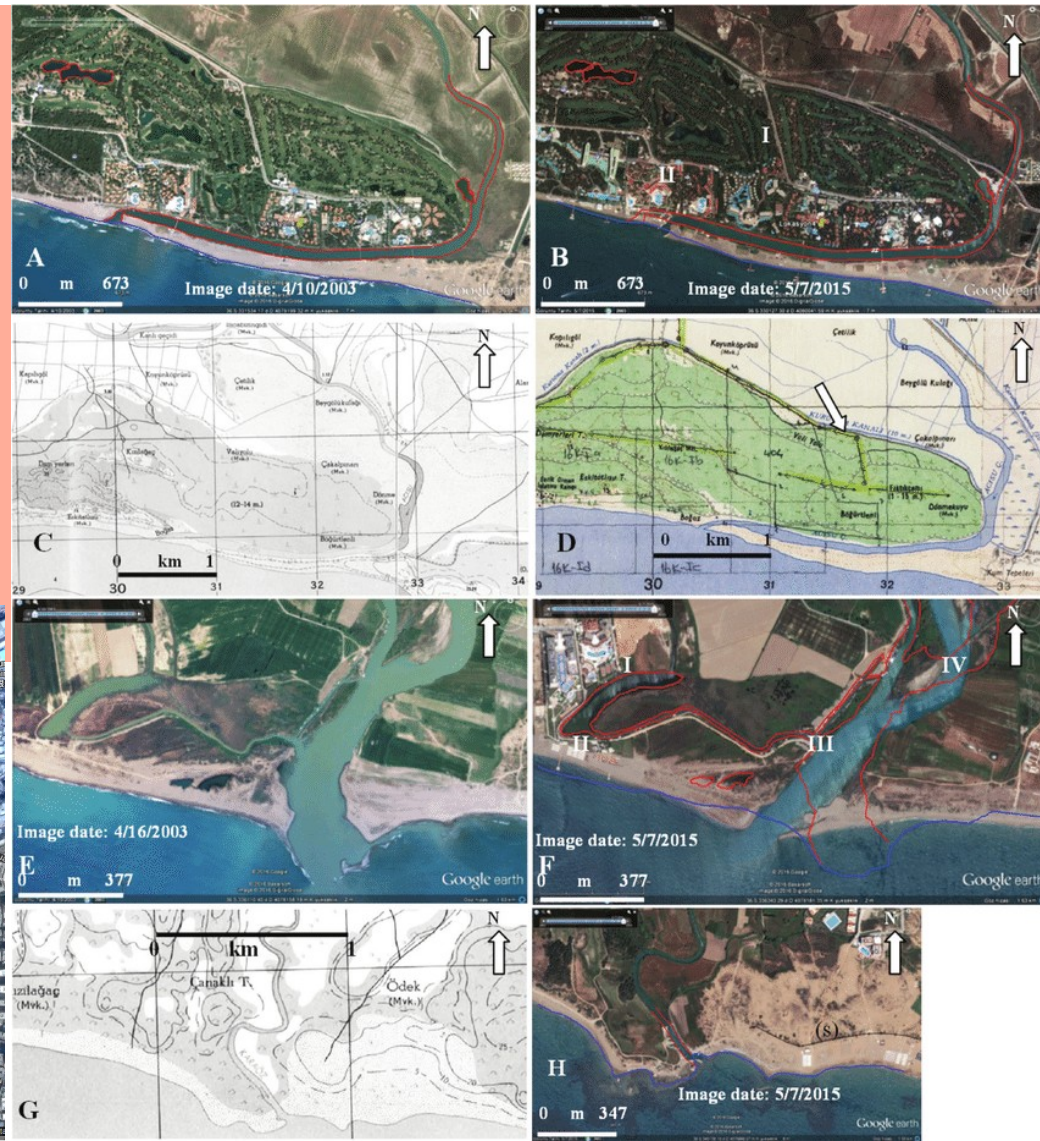
[https://www.csicr.cz/Csicr/media/Prilohy/PDF\\_el.\\_publikace/Mezinárodní%20šetření/PISA\\_2012\\_S\\_A.pdf](https://www.csicr.cz/Csicr/media/Prilohy/PDF_el._publikace/Mezinárodní%20šetření/PISA_2012_S_A.pdf)

Zdroj: PISA 2012

# S regresní analýzou opouštíme oblast popisných statistik a vstupujeme na pole modelování.

Model je vždy abstrakce - odhlížíme od detailů ve prospěch zachycení důležitého vzorce

Srov.  
Letecký snímek vs.  
mapa

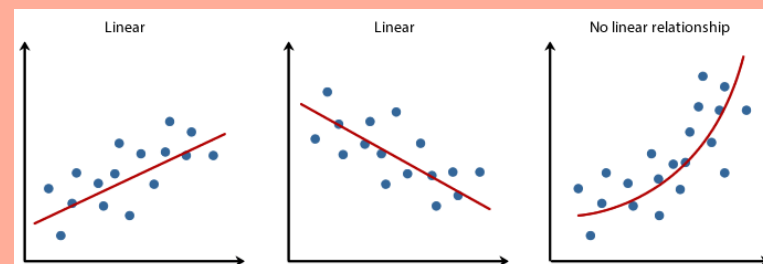


# U MODELŮ LZE OBECNĚ UVAŽOVAT O DVOU ASPEKTECH:

## PARAMETRY MODELU: co model říká

- jaký tvar má souvislost,
- jaký je sklon přímky...

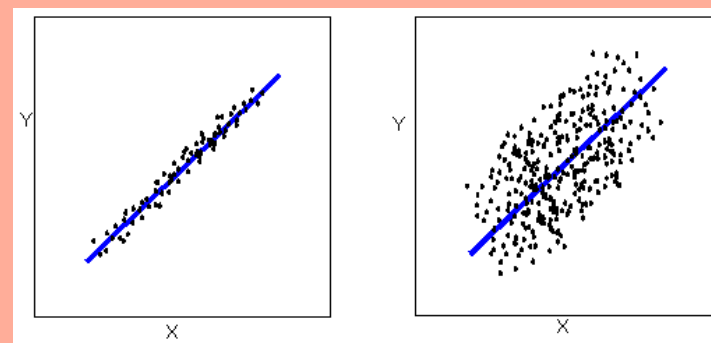
## Regresní rovnice a její členy



## VHODNOST MODELU: jak dobře model reprezentuje data, jak "sedí,"

### Míry vhodnosti modelu

(nový význam  $R^2$  – zde jako míra vhodnosti modelu, podobně také ANOVA – analýza rozptylu)



# Cíl lineární regrese

Sumarizovat vztah mezi dvěma proměnnými ve formě **rovnice přímky** (neboť předpokládáme lineární, tj. přímkový vztah) prostřednictvím výpočtu **regresního koeficientu**:

$$y = a + b \cdot x \quad (Y = b_0 + b_1 \cdot X_1)$$

y .... hodnota závisle proměnné, tu chceme predikovat  
(*outcome*)

a .... parametr, který říká, v jakém bodě přímka protíná  
vertikální osu Y (hodnota Y, když X = 0),

b .... regresní koeficient -- určuje směr přímky,  
(*predictor*)

x.... hodnota nezávisle proměnné, slouží k predikci  
hodnoty y

S regresní rovnicí se můžeme setkat v různých formách:

$$Y_i = \beta_0 + \beta_1 X_i + \epsilon_i$$

Labels: Dependent Variable ( $Y_i$ ), Population Y intercept ( $\beta_0$ ), Population Slope Coefficient ( $\beta_1$ ), Independent Variable ( $X_i$ ), Random Error term ( $\epsilon_i$ ).  
Components: Linear component ( $\beta_0 + \beta_1 X_i$ ), Random Error component ( $\epsilon_i$ ).

nejčastěji

$$y = a + b * x$$

$$y_i = b_0 + b_1 x + e$$

Labels: Estimated (or predicted) y value ( $y_i$ ), Estimate of the regression intercept ( $b_0$ ), Estimate of the regression slope ( $b_1$ ), Independent variable ( $x$ ), Error term ( $e$ ).

$$Y' = A + B * X$$

SIMPLE REGRESSION EQUATION

Labels:  $X$ : predictor (present in data),  $B$ : coefficient (estimated by regression),  $A$ : intercept (estimated by regression),  $Y'$ : predicted value (calculated from A, B and X).

**Vždy ale musí obsahovat dva členy (které se však taky nazývají různě):**

Něco, co určuje, **kde přímka protíná osu** Y – to je nutné pro její umístění v rovině – názvy: *konstanta, posunutí, intercept*

Něco, co určuje **sklon přímky**: *směrnice, koeficient, slope*



**Y**

**b = Regresní koeficient**  
(*slope*)

**b**

1

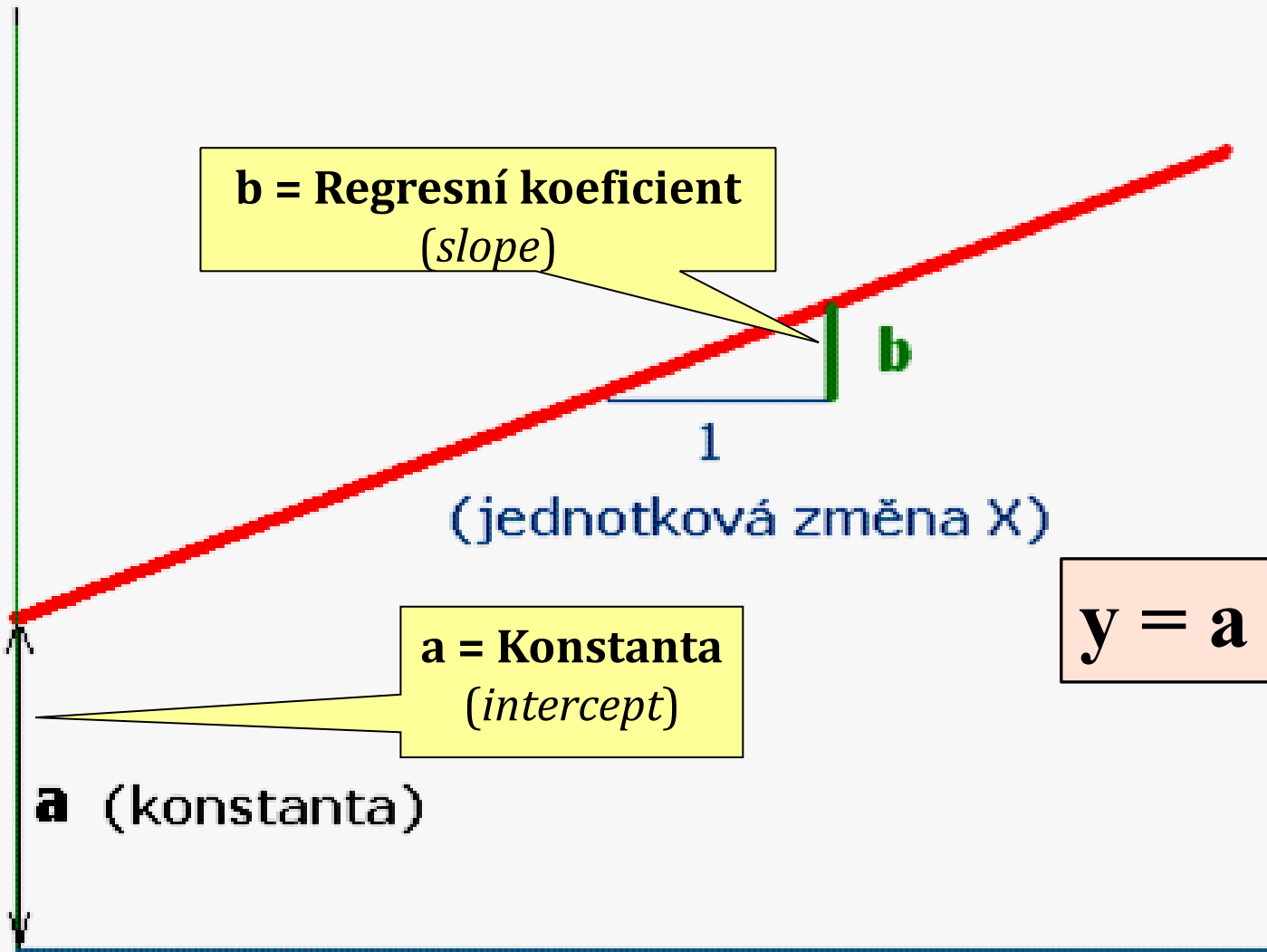
(jednotková změna X)

**a = Konstanta**  
(*intercept*)

**a** (konstanta)

$$y = a + b * x$$

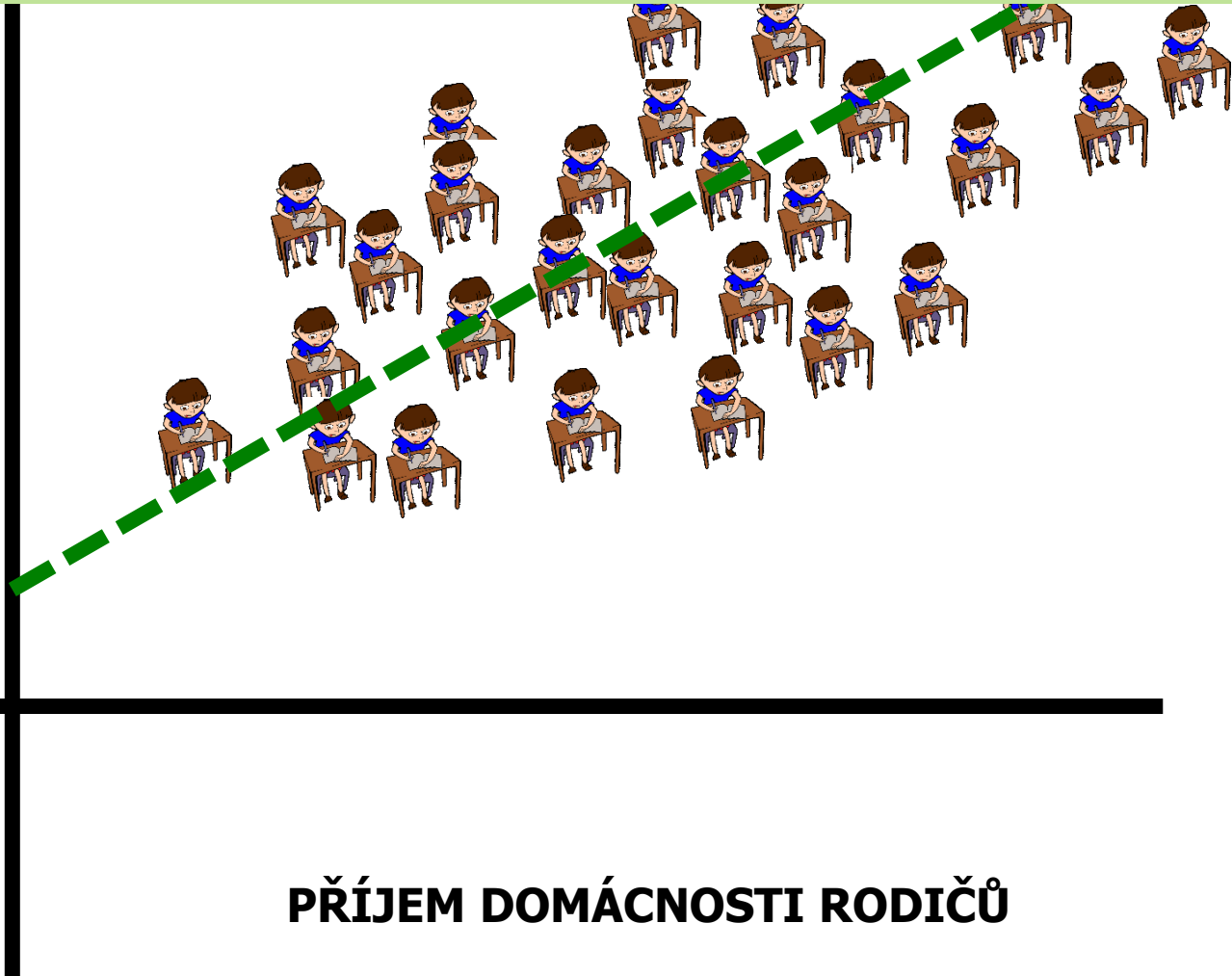
**X**



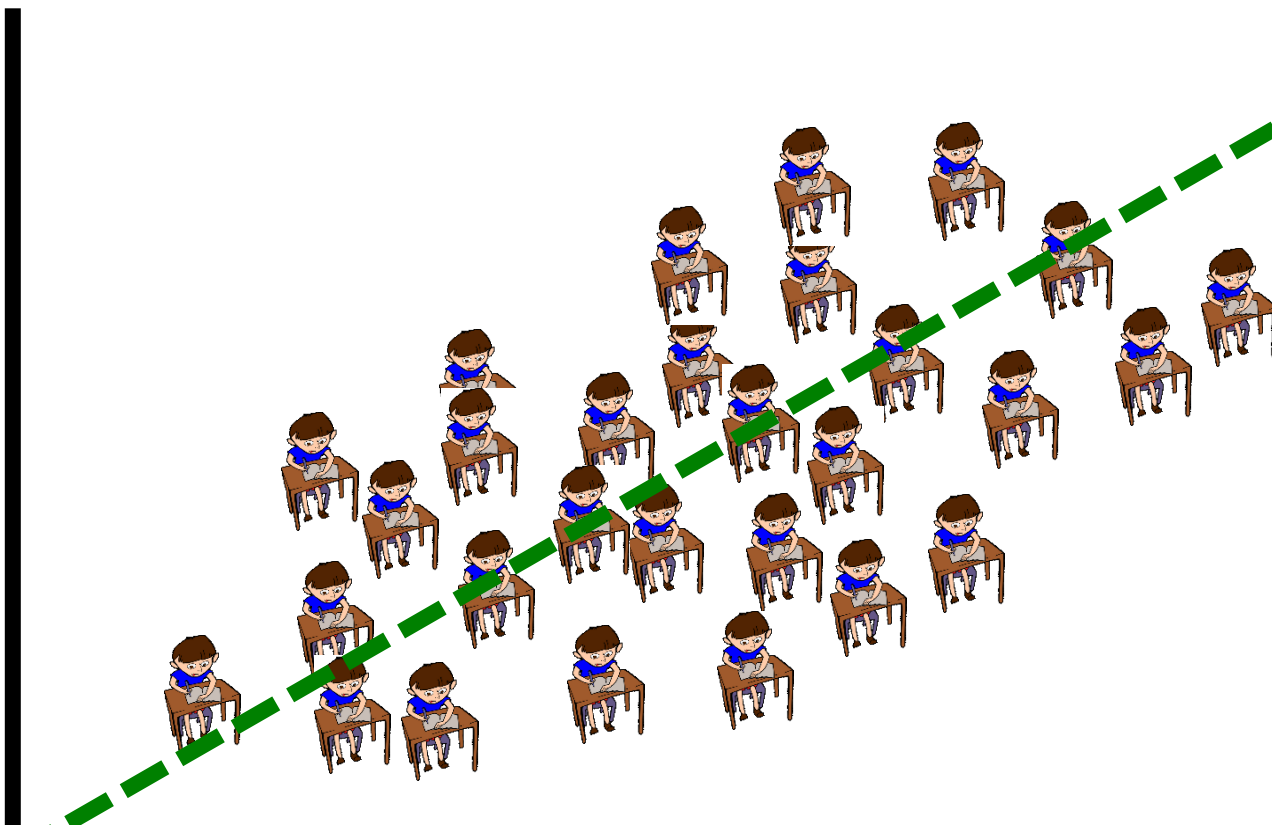
SKÓRE MATEMATICKÝCH DOVEDNOSTÍ  
0-100

Takže například:

skóre =  $15 + 0,005 * \text{příjem rodičů}$



RE MATEMATICKÝCH DOVEDNOSTÍ  
0-100



**Takže například:**

**Jestliže se příjem domácnosti zvýší o korunu, zvýší se skóre žáka o 0,005 bodu**

**... jestliže se zvýší o 1000 Kč, zvýší se skóre o 5 bodů**

# Podmínky pro užití lineární regresní analýzy:

1. Vztah mezi analyzovanými proměnnými musí být lineární,
2. závisle proměnná  $Y$  je měřena na intervalové úrovni a nezávisle proměnná  $X$  je buď intervalová, nebo dichotomická,
3. obě proměnné by měly být přibližně normálně rozloženy
4. V datech není přítomna heteroskedasticita