

## Výfučtení: O kráse regrese

### Co znamená lineární?

Následující text se vás pokusí nenásilně seznámit s procesem lineární regrese, který je pro zpracování výsledků fyzikálního měření mnohdy velice praktický. Mnoho fyzikálních veličin spolu souvisí přes nějakou konstantu a jejich závislost je lineární. Slovo lineární znamená, že danou závislost lze zapsat jako polynom prvního stupně, tedy

$$y = ax + b,$$

kde  $a$  a  $b$  jsou nějaké pevné konstanty a  $x$  je proměnná. Grafem takovéto funkce je přímka, linie, odtud slovo lineární. Konstanta  $b$  určuje, kde bude graf funkce protínat osu  $y$  (dosadíme-li za  $x = 0$ , dostaneme  $y = b$ ) a konstanta  $a$  zase určuje, jak strmý bude růst veličiny  $y$  (graf  $y = 3x + 1$  poroste strměji než  $y = 2x + 1$ ). Pokud konstanta  $b$  bude nulová, graf funkce bude procházet počátkem, pokud bude nulová konstanta  $a$ , závislost  $y = 0x + b = b$  je konstantní a už o lineární závislosti nemluvíme.

Příkladem lineárních závislostí v přírodě může být například závislost napětí na rezistoru na proudu, který jím prochází,

$$U = RI.$$

Konstanta  $b$  je nulová (při nulovém proudu je nulové napětí) a konstanta  $a$  je rovna odporu rezistoru  $a = R$ .

Další lineární závislostí může být třeba závislost objemu vody, která je v nádobě s kolmými stěnami na výšce hladiny v ní. Konstanta úměrnosti  $a$  je v tomto případě plocha dna nádoby.

$$V = Sh.$$

### Co znamená regrese?

Regrese pochází z latinského slova *regredi*, což znamená navracet se, ustupovat. Do statistiky toto slovo zavedl Francis Galton a označil tím „návrat k průměru“, fakt, že vysocí otcové mívají často syny menší, než jsou oni sami a synové malých otců bývají zase často vyšší, než jsou oni sami. Pojem se rozšířil na jakékoliv zkoumání závislostí náhodných veličin.

My budeme pojmem „lineární regrese“ označovat prokládání přímky naměřenými daty tak, aby přímka co nejlépe vystihovala jejich závislost. Body reprezentující jednotlivá měření přitom nemusí na přímce úplně ležet, ale rádi bychom, aby jim přímka odpovídala co nejlépe.

### Elektrický příklad

Lineární regrese se například využije v případě měření odporu drátu. Do měření mohou vstoupit nejrůznější chyby. Ručička přístrojů se například může na určitých místech stupnice trochu zadržávat, zdroj napětí může při určitých napětích více kolísat než při jiných. Abychom minimalizovali vliv těchto chyb, nebudeme měřit proud a napětí na drátu pouze jednou, ale několikrát a to při různých hodnotách. Výsledky měření naleznete v tabulce 1.

Nyní bychom chtěli data proložit lineární funkcí  $f(I) = a \cdot I + b$  tak, aby byla hodnota  $f(I)$  vždy naměřeným hodnotám co nejbliže. Toto prokládání se děje metodou nejmenších čtverců. To znamená, že vyzkoušíme spoustu různých lineárních funkcí  $f(I)$  a pak pro každý naměřený

Tabulka 1: Závislost napětí na proudu

$I/A$	1	2	3	4	5	6	7
$U/V$	2,1	3,9	5,7	8,2	10,2	11,9	13,8

proud spočítáme „čtverec“  $c$ , což je druhá mocnina rozdílu naměřeného napětí  $U$  při proudu  $I$  a hodnoty lineární funkce v bodě  $I - f(I)$ . Zapsáno vzorcem

$$c(I) = (U(I) - f(I))^2.$$

Všechny tyto čtverce pro všech sedm naměřených hodnot bychom sečetli a pak bychom je mezi sebou porovnali pro různé zkoušené lineární funkce. U té funkce  $f(I)$ , kde by byl součet čtverců nejmenší, je jasné, že naměřeným hodnotám vyhovuje nejlépe. Takto bychom třeba zjistili, že  $f_1(I) = 2 \cdot I + 0$  je lepší než  $f_2(I) = 3 \cdot I + 1$ . Protože pro funkci  $f_1$  je součet čtverců 0,24 a pro funkci  $f_2$  je tento součet 206,04.

Je  $f_1$  nejlepší? Co třeba  $f_3(I) = 1,9 \cdot I + 0$ ? Nebo  $f_2(I) = 1,09 \cdot I + 0,1$ ? Lineárních funkcí je nekonečně mnoho (za  $a$  a  $b$  si můžeme zvolit jakákoli čísla). Proto je celé toto počítání metody nejmenších čtverců vhodné vhodně svěřit počítači. Ten sice nevyzkouší všechny lineární funkce, ale jistě jich i pro velké počty naměřených dat zvládne víc než vy. Funkce lineární regrese má zabudovaně například tabulkový procesor, nebo můžete použít program na kreslení grafů *gnuplot*, jako jsme to udělali my. Program *gnuplot* určil, že nejlepší lineární funkce (viz obrázek 1), jaká sedí na naše data je

$$f(I) = 1,99 \cdot I + 0,1.$$

Dokonce zvládne vypočítat, jak přesné je určení koeficientů  $a$ ,  $b$ :

$$a = 1,99 \pm 0,04,$$

$$b = 0,1 \pm 0,1.$$

Konstanta  $a$  pak hraje roli námi hledaného odporu a proto můžeme zpracování slavnostně zakončit

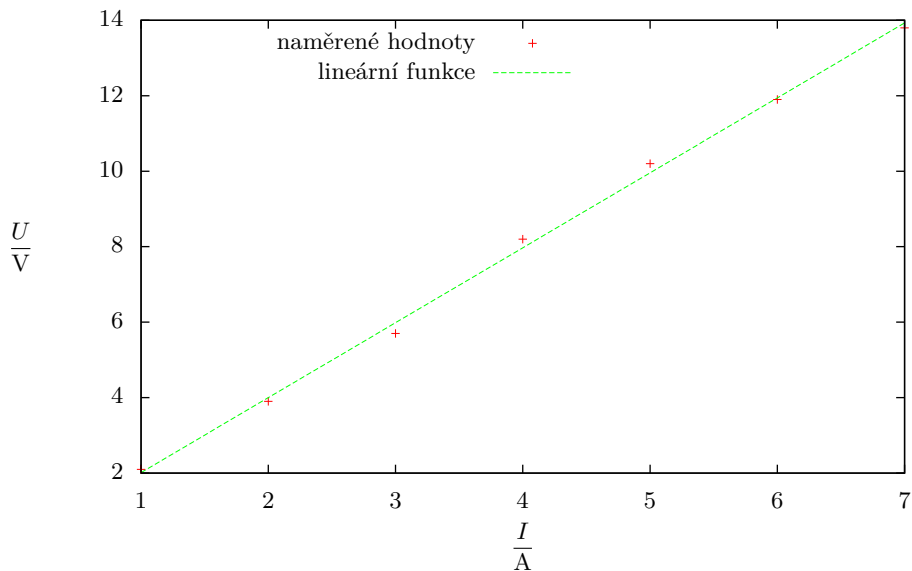
$$a = 1,99 \pm 0,04 \Omega.$$

### Co když závislost není lineární?

Metodu nejmenších čtverců můžeme samozřejmě použít i na jiné funkce, než lineární. Nejdříve musíme odhadnout, jak by závislost mohla vypadat. Můžeme zkoušet například různé polynomy, tam už budeme potřebovat více parametrů, například u polynomu třetího stupně budeme hledat čtyři konstanty

$$f(x) = ax^3 + bx^2 + cx + d.$$

Idea je stejná, vyzkoušet různé funkce a porovnat součty čtverců pro všechna naměřená data. Obecně je to ale výpočetně náročnější a ne vždy se to povede. Pokud bude totiž chyba určení



Obr. 1: Data proložená přímkou dle lineární regrese

parametrů  $a, b, c, \dots$  stejně velká nebo dokonce větší než parametry samy, nemá takový výsledek valného významu a program, který použijeme, může zahlásit, že se mu vhodné parametry nepodařilo nalézt. Potom je třeba použít jiný typ funkce.

---

Fyzikální korespondenční seminář je organizován studenty MFF UK. Je zastřešen Oddělením pro vnější vztahy a propagaci MFF UK a podporován Ústavem teoretické fyziky MFF UK, jeho zaměstnanci a Jednotou českých matematiků a fyziků.

Toto dílo je šířeno pod licencí Creative Commons Attribution-Share Alike 3.0 Unported. Pro zobrazení kopie této licence, navštivte <http://creativecommons.org/licenses/by-sa/3.0/>.