

Závislosť náhodných udalostí a premenných

Formálna definícia závislosti udalostí:

Náhodné udalosti A_1, \dots, A_n sú združene nezávislé, ak pre akýkoľvek výber z nich A_{i_1}, \dots, A_{i_m} platí

$$P[A_{i_1} \& \dots \& A_{i_m}] = P[A_{i_1}] \dots P[A_{i_m}].$$

Intuitívne:

A_1, \dots, A_n sú združene nezávislé, ak informácia, že nastala udalosť vytvorená z A_{k_1}, \dots, A_{k_m} , nezmení naše očakávanie, že nastala udalosť vytvorená z iných udalostí A_{j_1}, \dots, A_{j_l} („z iných udalostí“ = $\{k_1, \dots, k_m\}$ a $\{j_1, \dots, j_l\}$ sú disjunktné).

Formálna definícia závislosti náhodných premenných:

Náhodné premenné X_1, \dots, X_n sú združene nezávislé, ak pre akékoľvek B_1, \dots, B_n sú $[X_1 \text{ patrí do } B_1], \dots, [X_n \text{ patrí do } B_n]$ združene nezávislé udalosti.

Intuitívne:

X_1, \dots, X_n sú združene nezávislé, ak informácia o tom, aké hodnoty nadobudli X_{k_1}, \dots, X_{k_m} , nezmení naše očakávanie ohľadom hodnôt, ktoré nadobudli iné premenné X_{j_1}, \dots, X_{j_l} („iné premenné“ = $\{k_1, \dots, k_m\}$ a $\{j_1, \dots, j_l\}$ sú disjunktné).

Nezávislosť sa typicky vyskytuje v:

Merania (pozorovania/pokusy/výbery), o ktorých výsledku nám neposkytujú žiadnu informáciu výsledky iných meraní (pozorovaní/pokusov/výberov).

Závislosť sa typicky vyskytuje v:

Merania (pozorovania/pokusy/výbery), o ktorých výsledku nám poskytujú informáciu výsledky iných (napríklad časovo predchádzajúcich) meraní (pozorovaní/pokusov/výberov)

Štatistické testovanie nezávislosti:

To, či sú dáta realizáciou nezávislých premenných vieme (nedokonale!) testovať štatistickými testami, napríklad:

- Testy nulovosti rôznych mier „konkordancie“ dvojice náhodných premenných (napr. Pearsonovho korelačného koeficientu a iných),
- Testy rovnosti združeného rozdelenia a súčinu marginálnych rozdelení (napr. chíkvadrát test nezávislosti a mnohé iné).

Kovariancia náhodných premenných

Formálna definícia kovariancie a vlastnosti:

Pre náhodné premenné X, Y definujeme

$$\text{cov}(X, Y) = E[(X - EX)(Y - EY)].$$

$\text{cov}(X, Y) = \text{cov}(Y, X)$, $\text{cov}(X+a, Y+b) = \text{cov}(X, Y)$, $\text{cov}(aX, bY) = a \cdot b \cdot \text{cov}(X, Y)$ pre všetky a, b . Ak sú X, Y nezávislé, $\text{cov}(X, Y) = 0$.

Intuitívne:

Kovariancia je symetrická miera „lineárnej“ závislosti premenných X a Y , na ktorú nemá vplyv posun hodnôt premenných X a Y , avšak má vplyv zmena jednotiek X a Y .

Interpretácia:

- $\text{cov}(X, Y) > 0$: s vyšším X sa dajú čakať vyššie Y ,
- $\text{cov}(X, Y) < 0$: s vyšším X sa dajú čakať nižšie Y ,
- $\text{cov}(X, Y) = 0$: zvyšovanie X neovplyvňuje stredný rast/pokles Y ,
- Zo samotnej hodnoty $\text{cov}(X, Y)$ je ťažké niečo vyčítať...

Výberový odhad kovariancie z n nezávislých dvojíc realizácií:

$$\widehat{\text{cov}}(X, Y) = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

V tomto vzťahu sú realizácie $(x_1, y_1), \dots, (x_n, y_n)$ a \bar{x}, \bar{y} znamenajú aritmetický priemer.

(Pearsonova) korelácia náhodných premenných

Formálna definícia:

Pre náhodné premenné X, Y definujeme

$$\text{cor}(X, Y) = \frac{\text{cov}(X, Y)}{\sqrt{DX}\sqrt{DY}}$$

$\text{cor}(X, Y) = \text{cor}(Y, X)$, $\text{cor}(X+a, Y+b) = \text{cor}(X, Y)$ pre všetky a, b ,
 $\text{cor}(aX, bY) = \text{cor}(X, Y)$ pre všetky $a > 0, b > 0$. Ak sú X, Y nezávislé,
 $\text{cor}(X, Y) = 0$. Tiež

$$-1 \leq \text{cor}(X, Y) \leq 1$$

Intuitívne:

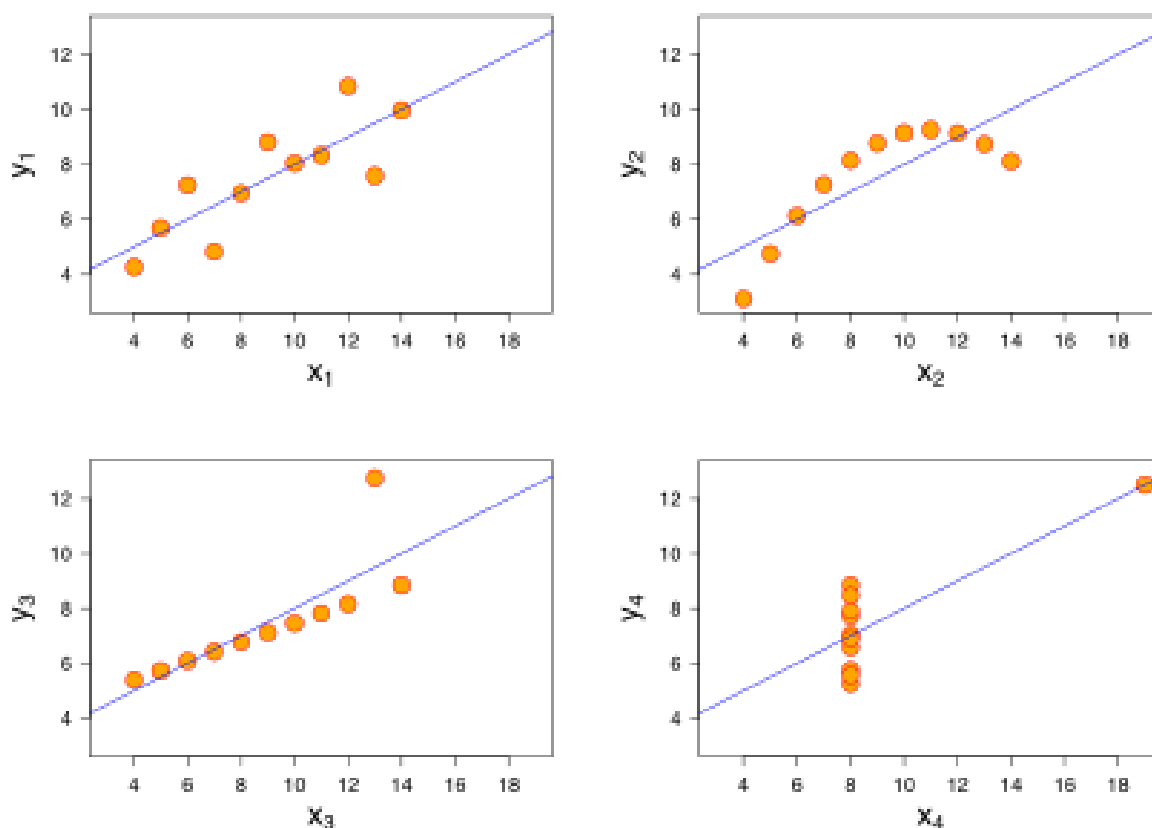
Korelácia je symetrická, normalizovaná miera „lineárnej“ závislosti premenných X a Y , na ktorú nemá vplyv posun hodnôt premenných X a Y , ani zmena jednotiek X a Y .

Interpretácia:

- $\text{cor}(X, Y) > 0$: s vyšším X sa dajú čakať vyššie Y ,
- $\text{cor}(X, Y) < 0$: s vyšším X sa dajú čakať nižšie Y ,
- $\text{cor}(X, Y) = 0$: zvyšovanie X nevplyva na stredný rast/pokles Y ,
- Ak je $\text{cor}(X, Y)$ blízke -1 , máme skoro presnú klesajúcu lineárnu závislosť premenných X a Y ,
- Ak je $\text{cor}(X, Y)$ blízke 1 , máme skoro presnú rastúcu lineárnu závislosť premenných X a Y .

"Nedokonalosti" $cor(X,Y)$ z praktického hľadiska:

- Rôzne "formy" závislosti môžu viesť na ten istý korelačný koeficient, vid' napr. *Anscombe's quartet* (obrázok je z wiki)



- Pozri tiež *Datasaurus dozen* (napr. <https://cran.r-project.org/web/packages/datasauRus/vignettes/Datasaurus.html>)

- Závislé náhodné premenné môžu mať nulový korelačný koeficient, ako je vidieť z *Datasaurus dozen*

Výberový odhad korelácie z n nezávislých dvojíc realizácií:

$$\widehat{cov}(X, Y) = \frac{cov(X, Y)}{S_X S_Y},$$

kde S_X je odmocnina výberového rozptylu realizácií x_1, \dots, x_n náhodnej premennej X a analogicky S_Y pre náhodnú premennú Y.

Súvis výberovej korelácie s regresnou priamkou:

Ak $\hat{a}x + \hat{b}$ je regresná priamka vypočítaná z dvojíc $(x_1, y_1), \dots, (x_n, y_n)$ a $\widehat{cov}(X, Y)$ je výberový korelačný koeficient vypočítaný z tých istých dvojíc, tak

$$\hat{a} = \widehat{cov}(X, Y) \frac{S_Y}{S_X}.$$

Test nulovosti korelačného koeficientu:

$$t = \widehat{cov}(X, Y) \sqrt{\frac{n-2}{1 - \widehat{cov}(X, Y)^2}}$$

má tazvané Studentovo rozdelenie s n-2 stupňami voľnosti, ak sú $(x_1, y_1), \dots, (x_n, y_n)$ realizácie dvojrozmerného normálneho rozdelenia a X, Y majú nulový korelačný koeficient.

Iné typy mier závislosti:

Rôzne druhy korelačných koeficientov (Spearmanovo rho, Kendallovo tau), alebo úplne odlišné miery závislosti

Kauzalita

Definícia (Collins dictionary): *Causality is the relationship of cause and effect. The cause of an event is the thing that makes it happen. The effect of one thing on another is the change that the first thing causes in the second thing. :)*

Otázky týkajúce sa kauzality, ktoré tu nemôžeme riešiť:

- Čo presne je kauzalita?

Jedna zo základných otázok filozofie spadajúca do oblasti ontológie a epistemológie. Táto otázka je zásadná (a nevyjasnená) pre mnohé vedecké disciplíny.

- Ako matematicky formalizovať uvažovanie o kauzalite?

Množstvo prístupov: pravdepodobnostné teórie kauzality, kauzálny kalkulus, teória kontrafaktuálov a iné. Ide o zložité teórie, ktoré si vyžadujú dlhodobé štúdium.

- Ako zistiť, či nejaká premenná kauzálne ovplyvňuje inú premennú?

Viacero prístupov, jeden z nich sa opiera o kontrolované experimenty, resp. o "intervenciu" (o "vynútenie" hodnôt alebo rozdelenia nejakej premennej).

Korelácia nie je kauzalita:

Spresnenie: Nenulová korelácia premenných X a Y neimplikuje kauzálnu súvislosť* medzi X a Y (čiže ani to, že X je príčinou Y, ani to, že Y je príčinou X).

*Nech by sme kauzalitu (pojmy príčiny a následku) formalizovali akýmkoľvek „rozumným“ spôsobom.

To, že premenné X a Y sú korelované je často dôsledkom nie kauzálnej závislosti X od Y alebo Y od X, ale tým, že existuje latentná premenná Z, ktorá kauzálny ovplyvňuje X aj Y.

Príklad: Nech X je hmotnosť žiaka základnej školy a Y je veľkosť jeho slovnej zásoby. X a Y sú korelované, avšak je zjavné*, že nemôžeme tvrdiť, že X je príčinou Y alebo naopak. Existuje však iná premenná, V, napríklad vek (alebo vyspelosť) žiaka, ktorá je kladne korelovaná s X aj s Y, a o ktorej môžeme tvrdiť, že kauzálny ovplyvňuje aj X aj Y.

Úloha: Štúdia v časopise Nature

(<https://www.nature.com/articles/20094>) zistila významnú kladnú koreláciu medzi krátkozrakosťou detí a tým, koľko boli vystavené nočnému svetlu ako novorodenci. Dá sa z toho usúdiť, že nočné svetlo spôsobuje krátkozrakosť?