

3. cvičení z PaSti – 2022-03-01

Nezávislé jevy

- Jevy A a B jsou nezávislé, pokud $P(A \cap B) = P(A)P(B)$. To je totéž jako $P(A | B) = P(A)$ nebo $P(B) = 0$.
 - Jevy A_1, \dots, A_n jsou nezávislé, pokud $P(\cap_{i \in I} A_i) = \prod_{i \in I} P(A_i)$, kdykoliv $\emptyset \neq I \subseteq \{1, \dots, n\}$.
1. Jsou-li jevy A, B nezávislé, pak jsou nezávislé i A, B^c a také A^c, B^c .
 2. Mohou být dva jevy nezávislé a zároveň disjunktní?
 3. Najděte jevy A, B, C takové, že jsou po dvou nezávislé, ale $P(A \cap B \cap C) \neq P(A)P(B)P(C)$.
 4. Najděte jevy A, B, C takové, že $P(A \cap B \cap C) = P(A)P(B)P(C)$, ale jevy nejsou po dvou nezávislé.

Testování vlastností

Motivace: Nechtě některé elementární jevy mají určitou vlastnost a máme k dispozici test, který ji ne úplně spolehlivě odhaluje – příkladem může být třeba výskyt nemoci v populaci a test na tuto nemoc. Uvažme jevy $N =$ „je nemocný“ a $T =$ „test vyšel pozitivně“.

Představme si, že test danému jedinci vyšel pozitivně. To může znamenat, že nemoc má (*pravý pozitivní výsledek*), nebo že ji nemá, ale test se spletl (*falešný pozitivní výsledek*). Podobně pokud test vyšel negativní, může to být výsledek *pravý negativní* nebo *falešný negativní*. (Anglicky true/false positive/negative.)

Pravděpodobnost pravého pozitivního výsledku je zjevně $P(N | T)$, falešného pozitivního $P(N^c | T)$, podobně pravého negativního $P(N^c | T^c)$ a falešného negativního $P(N | T^c)$. Přitom zjevně $P(N | T) = 1 - P(N^c | T)$ apod.

Při výrobě testu obvykle statisticky testujeme úspěšnost testu na lidech, u kterých víme, zda jsou nemocní. Tím získáme opačně podmíněné pravděpodobnosti $P(T | N)$ a $P(T | N^c)$. Bayseova věta nám je umožňuje obrátit a spočítat pravděpodobnosti z předchozího odstavce:

$$P(B_i | A) = \frac{P(A | B_i)P(B_i)}{\sum_i P(A | B_i)P(B_i)}, \quad P(N | T) = \frac{P(T | N)P(N)}{P(T | N)P(N) + P(T | N^c)P(N^c)}.$$

Ještě k terminologii: U testů se také mluví o *senzitivitě* a *specifitě*. To první je $P(T|N)$ – jaký zlomek nemocných se nám podařilo identifikovat – to druhé pak $P(T^c|N^c)$, tedy totéž pro zlomek zdravých.

Podobnou situaci najdeme i v kontextu vyhledávání informací: představte si třeba internetový vyhledávač, kterému položíte dotaz. Nemocní jsou správné odpovědi na dotaz, pozitivně testovaní jsou výsledky vyhledávače. Tentokrát se parametrům říká

precision a *recall*. Precision je zlomek správných odpovědí mezi nalezenými, tedy $P(N | T)$, čili pravděpodobnost pravých pozitiv. Recall říká, jaký zlomek správných odpovědí jsme našli – to je $P(T | N)$, čili totéž, co senzitivita. Ať žije jednotnost terminologie ;)

Bayesova věta

5. Petr dostává hodně emailů, ale 80 % z nich jsou spamy. Jeho spamový filtr 90 % spamů správně označí, ale také 5 % řádných emailů označí jako spam.

- Kolik procent emailů bude označeno jako spamy?
- Kolik procent řádných emailů je mezi těmi, co jsou označené jako spamy?
- Kolik procent spamů je mezi emaily, které testem prošly?

6. Kouřovými signály přenášíme binární soubor. Je proto poměrně vysoká pravděpodobnost chyby u každého bitu: 0 se jako 0 přeneso jen s pravděpodobností 0.9, 1 jako 1 jen s pravděpodobností 0.8. Předpokládejme (trochu neseriózně), že obsahem souboru jsou náhodné bity, které se přenášejí nezávisle.

- Pokud jsme dostali signál 0, jaká je pravděpodobnost, že byl opravdu vyslán?
- Dostali jsme zprávu 0010. Jaká je pravděpodobnost, že byla opravdu vyslána?
- Jak se výpočet změní, pokud budeme pro kontrolu vysílat každý symbol třikrát (a pak vezmeme častější z těch tří pokusů)?

Úlohu si můžete zjednodušit předpokladem, že 0 a 1 se mají stejnou pravděpodobnost správného přenosu.

K procvičení

7. Pro plánování výletu do Krkonoš používáme českou a polskou předpověď počasí. Každá nám dá binární výsledek *hezky/ošklivo* a má pravděpodobnost úspěchu $p \in [0, 1]$; obě předpovědi jsou nezávislé. Používáme je takto: pokud se shodují, věříme jim, pokud ne, hodíme si korunou. Jaká je pravděpodobnost, že se rozhodneme správně?

8. Na chorobu C máme dva testy, A a B . Test A má sensitivitu i specificitu $p = 0.95$. Test B vždy řekne, že pacient je zdravý. Předpokládejte, že $P(C) = 0.01$.

- Spočtete pro oba testy pravděpodobnost úspěchu (tj. správné odpovědi), použijeme-li je na náhodného pacienta. Co to říká o užitečnosti obou testů?
- Pro jaké p je pravděpodobnost úspěchu obou testů stejná?

9. Ve volbách hlasují lidé pro dva kandidáty, A a B . Při odchodu z volební místnosti jsou voliči náhodně požádáni o účast v exit-poll. Předpokládejme, že kdo odpoví, odpoví popravdě koho volil, ale ne všichni se zúčastní. Označíme-li E množinu voličů, kteří se exit-pollu zúčastní, tak předpokládejme $P(E | A) = 0.7$ a $P(E | A^c) = 0.4$. Výsledky exit-pollu jsou 60 % pro A . Jaký je skutečný podíl lidí, kteří hlasovali pro A ?