

Skupinové testování

ANTONÍN JANČAŘÍK – TOMÁŠ KEPKA

Pedagogická fakulta UK, Praha – Matematicko-fyzikální fakulta UK, Praha

Během pandemie covid-19 byla často diskutována kapacita testovacích zařízení a možnosti, jak ji navýšit. Jednou z metod, který byla v praxi skutečně realizována, bylo i poolování. Jedná se o metodu, kdy je smícháno a současně testováno více vzorků. V praxi bylo poolování použito tak, že test probíhal ve dvou krocích. V prvním kroku byly testovány malé skupiny vzorků (dle doporučení MZČR o velikosti 6–10 vzorků). Pokud některá skupina byla pozitivní, byl ve druhém kroku testován každý vzorek dané skupiny samostatně. V tomto článku představíme metody skupinového testování (Combinatorial group testing), které stojí na pomezí mezi matematikou a informatikou a které umožňují přímo v prvním kroku testu detekovat konkrétní infikované vzorky. Prezentované algoritmy jsou dostupné žákům středních škol a využívají především vlastnosti zápisu čísla v soustavě o jiném základu.

1. Skupinové testování

Metody skupinového testování (Group testing) byly poprvé představeny v roce 1943 jako nástroj pro ekonomicky výhodnější testování vzorků krve vojáků (Dorfman, 1943). Cílem použití matematických metod je snížení počtu testů, které je nutné realizovat, oproti stavu, kdyby byl každý vzorek analyzován samostatně. Je zjevné, že metody pro testování krve lze přímo přenést na PCR testování vzorků covid-19. Ukazuje se však, že metody, které byly původně připraveny pro potřeby zdravotnictví, jsou uplatnitelné v mnoha dalších oblastech, jako je kontrola kvality výrobků či prohledávání souborových uložišť.

V tomto textu se budeme nejprve zabývat pouze případy, kdy všechny testy musí probíhat v jednom okamžiku a počet „infikovaných“ vzorků je velmi malý. Představíme známé algoritmy na nalezení jednoho, dvou a tří „infikovaných“ vzorků. Přičemž u testu směřujícím k nalezení tří infikovaných vzorků představíme test v drobné modifikaci, která v některých případech přináší zlepšení oproti doposud známým testům. V navržených metodách, na rozdíl od současných požadavků praxe, neuvažujeme žádné omezení na počet vzorků, které lze smíchat a testovat současně.

2. Nalezení jednoho vadného vzorku

Na případu jednoho infikovaného vzorku si představíme postupy, které jsou při skupinovém testování používány a otázky, které si můžeme v souvislosti s testováním klást.

Cílem je nalézt jeden „infikovaný“, resp. „vadný“ vzorek mezi n vzorky. Úlohu lze rozdělit na dva případy, na případ, kdy víme, že je právě jeden vzorek infikován a případ, kdy může předpokládat, že je nejvýše jeden vzorek infikován. Další důležitou otázkou je, zda test dokáže odhalit případ, kdyby mezi vzorky bylo více infikovaných a upozornit na tuto skutečnost.

Úloha tohoto typu je často řešena, čtenářům doporučujeme především knihu Dennisa E. Shashy, kde téma testování krve představuje poměrně atraktivní formou (Shasha, 2005, s. 30–35).

Úloha 1

Předpokládejme, že máme $n = 2^q$ vzorků, přičemž právě jeden z těchto vzorků je infikovaný. Navrhněte postup, jak pomocí co nejmenšího počtu testů zjistit, o který vzorek se jedná, za předpokladu, že v průběhu jednoho testu můžete smíchat více vzorků dohromady a test vyjde pozitivní právě tehdy, když jeden ze vzorků byl infikován.

Řešení. Každému ze vzorků přiřadíme index $X = X_{q-1} \dots X_0$, který odpovídá zápisu jeho pořadového čísla v binární soustavě, tedy $X_p \in \{0, 1\}$. Nyní stačí provést právě q testů, přičemž v testu i , $0 \leq i < q$, otestujeme všechny vzorky, pro jejichž index platí $X_i = 1$.

Protože ze zadání víme, že právě jeden ze vzorků je pozitivní, můžeme snadno sestavit jeho index. Pokud $X_i = 1$ v případě, že test i byl pozitivní, v opačném případě $X_i = 0$.

Diskuze řešení. Takto navržený algoritmus je velmi efektivní, dokáže nalézt konkrétní infikovaný vzorek v čase $O(\log_2 n)$, tedy například mezi tisícem vzorků je schopen infikovaný identifikovat na základě pouhých 10 testů. Na druhou stranu v každém testu musí být vždy zahrnuta polovina sledovaných vzorků, což může převyšovat technické možnosti realizace.

Test v předložené podobě není schopen rozlišit případ, kdy mezi nakaženými vzorky není žádný infikovaný od případu, kdy má infikovaný vzorek pořadové číslo 0. Tento problém by bylo možné odstranit přidáním jediného testu.

Stejně tak není test schopen identifikovat, zda mezi testovanými vzorky není více než jeden infikovaný. Tento problém by šlo řešit pomocí zdvojnásobení počtu testů, kdy pro každou pozici X_i by jeden test ověřoval,

vzorky pro které $X_i = 1$ a druhý $X_i = 0$. Takto upravený test by byl schopen určit, zda jsou všechny vzorky negativní, rozpoznat infikovaný vzorek za předpokladu, že je právě jeden a určit, že v testované sadě je více než jeden infikovaný vzorek. Tato situace nastává v případech, kdy pro některé i budou oba testy spojené s touto pozicí pozitivní. Ve speciální případě, kdy tato situace nastane právě pro jedno i , je test schopen přesně určit i oba infikované vzorky.

3. Nalezení dvou vadných vzorků

V rozšířeném testu z předchozího příkladu může ale nastat situace, kdy oba testy vyjdou pozitivní pro dvě (a více) různých pozic. Předpokládejme, že oba testy vychází pozitivní pro pozice i a j a testovaná sada obsahuje právě 2 infikované vzorky. Za této situace nemáme žádný indikátor, dle kterého by bylo možné rozhodnout mezi možnostmi, že jeden z infikovaných vzorků má na pozicích i a j kombinaci (0,0) a druhý (1,1) a možnostmi, že první ze vzorků má na pozicích i a j kombinaci (1,0) a druhý (0,1).

V nejhorsím možném případě může nastat situace, kdy všechny testy vyjdou i při pouhých dvou infikovaných vzorcích pozitivně, a to v případě, kdy první vzorek má index $X_{q-1} \dots X_0$ a druhý $X'_{q-1} \dots X'_0$, kde $\{X_i X'_i\} = \{0, 1\}$ pro všechna i .

Tuto situaci lze samozřejmě řešit tak, že budeme pro každou dvojici indexů testovat, která z možností nastává. Ve skutečnosti však není potřeba testovat všechny případy, které pro danou dvojici mohou nastat, ale pouze jeden z nich, tedy by stačilo přidat $\binom{q}{2}$ testů, kde pro dvojici (i, j) by byly testovány všechny vzorky, kde $X_i = 1 = X_j$. Například pro otestování již výše zmiňovaného tisíce vzorků by stačilo přidat 65 testů. Takto navržený test však jde ještě vylepšit.

Úloha 2

Předpokládejme, že máme $n = 3^q$ vzorků, přičemž nejvýše dva z těchto vzorků jsou infikované. Navrhněte postup, jak pomocí co nejmenšího počtu testů zjistit, o který vzorek se jedná, za předpokladu, že v průběhu jednoho testu můžete smíchat více vzorků dohromady a test vyjde pozitivní právě tehdy, když jeden ze vzorků byl infikován.

Řešení. Každému ze vzorků přiřadíme index $X = X_{q-1} \dots X_0$, který odpovídá zápisu jeho pořadového čísla v soustavě základu 3, tedy $X_p \in \{0, 1, 2\}$. Současně provedeme $3q + \binom{q}{2}$ testů. Jednu skupinu testů budou obsahovat testy, ve kterých budeme vždy pro každou pozici testovat všechny vzorky,

kde $X_i = j$, $j = 1, 2, 3$. Těchto testů je $3q$. Druhou skupinu testů budou obsahovat testy, ve kterých pro každou dvojici indexů (i, j) , $0 \leq i < j < q$, testujeme všechny vzorky, pro něž $X_i = X_j$.

Pokud žádný z první skupiny testů nevyšel pozitivní, testovaná skupina neobsahuje žádný infikovaný vzorek. Pokud pro každou pozici vyšel právě jeden pozitivní test, obsahuje testovaná skupina právě jeden infikovaný test a jeho index na každé pozici odpovídá testu, který byl pozitivní.

Pokud u právě jedné pozice vyšly dva pozitivní testy a pro ostatní pozice jen jeden pozitivní test, obsahuje testovaná skupina dva infikované vzorky a jejich indexy určíme snadno.

Zbývá tedy diskutovat případ, že pro dva a více pozic vyšly dva pozitivní testy (tři nemohou vyjít, protože soubor obsahuje nejvýše dva infikované vzorky).

Předpokládejme, že pro pozice i, j vyšly první skupině testu vždy dva pozitivní testy. Mohou nastat pouze dva případy:

1. Pro oba indexy vyšly pozitivní testy pro stejné dvě číslíce (označme je a, b), mohou tedy nastat dvě situace:
 - a. Oba infikované vzorky mají na pozicích i, j stejné číslíce a test s indexem (i, j) z druhé skupiny testů je pozitivní.
 - b. Oba infikované vzorky mají na pozicích i, j různé číslíce a test s indexem (i, j) z druhé skupiny testů je negativní.
2. Pro oba indexy vyšly pozitivní testy pro různé dvojice číslíc. Protože ale máme k dispozici jen tři různé číslíce, musí být testy pro jednu z číslíc pozitivní pro oba indexy. Předpokládejme tedy, že pro index i vyšly pozitivní testy pro číslíce a, b a pro index j vyšly pozitivní testy pro číslíce b, c . Opět mohou nastat jen dvě situace:
 - a. Jeden infikovaný vzorek má v indexu na pozicích i, j číslíce (a, c) a druhý vzorek má na těchto pozicích číslíce (b, b) a test s indexem (i, j) z druhé skupiny testů je pozitivní.
 - b. Jeden infikovaný vzorek má v indexu na pozicích i, j číslíce (a, b) a druhý vzorek má na těchto pozicích číslíce (b, c) a test s indexem (i, j) z druhé skupiny testů je negativní.

Diskuze. Přejít od binární soustavy k soustavě o základu tři přinesl významné vylepšení testu, např. pro výše zmiňovaných 1 000 vzorků potřebujeme 42, tedy o 23 testů méně, než u předchozího algoritmu. Navíc v každém testu postačí testovat třetinu všech vzorků, oproti polovině u předchozího testu. Efektivita celého testu vychází z číselného základu,

který byl zvolen tři právě proto, aby bylo možné rychle testovat případy dvojic.

4. Nalezení tří vadných vzorků

Eppstein, Goodrich & Hirschberg (2007) představili ve svém článku algoritmus, který je schopen nalézt až tři vadné vzorky. My zde v tomto článku představíme mírně vylepšenou verzi tohoto algoritmu, která spočívá v tom, že algoritmus modifikujeme pro soustavy o různém základu.

Mějme $n = k^q$ vzorků, které budou očíslovány pomocí čísel v soustavě o základu k , tedy pomocí čísel ve tvaru $X = X_{q-1} \dots X_0$ s číslicemi $X_p \in \{0, 1, \dots, k-1\}$. Nyní X představuje konkrétní číslo z množiny $\{0, 1, \dots, n-1\}$, p pozici číslice v zápise čísla (tedy $p \in \{0, 1, \dots, q-1\}$), a v je číslicí v dané soustavě (tedy $v \in \{0, 1, \dots, k-1\}$).

Nyní sestavíme matici M tak, že bude mít n sloupců (tedy každý sloupec odpovídá jednomu vzorku) a každý její řádek odpovídá jednomu testu. Řádky budou obsahovat pouze nuly a jedničky. Pokud je ve sloupci nula, není příslušný vzorek zahrnut do testu, pokud je ve sloupci jednička, naopak příslušný vzorek je součástí testu.

Matice M má $k^2 \binom{q}{2}$ řádků. Řádek (p, p', v, v') matice M je spojen s pozicemi číslic p a p' (kde $p < p'$) a číslicemi v a v' . $M[(p, p', v, v'), X] = 1$, právě když $X_p = v$ a $X_{p'} = v'$.

Nyní označme $\text{test}_M(p, p', v, v')$ výsledek (1 pro pozitivní a 0 pro negativní) testu, který odpovídá řádku (p, p', v, v') v matici M . Pro úplnost dodefinujeme pro $p' > p$, že $\text{test}_M(p', p, v, v') = \text{test}_M(p, p', v, v')$.

Pomocí funkce test_M můžeme definovat následující tři funkce:

$\text{test}_B(p, v)$ má hodnotu 1 pokud nějaký test ukazuje na pozitivní vzorek, který má číslici v na pozici p a hodnotu 0, pokud takový test neexistuje. Jinými slovy $\text{test}_B(0, v) = 0$, právě když $\sum_{i=0}^{k-1} \text{test}_M(0, 1, v, i) = 0$, jinak $\text{test}_B(0, v) = 1$.

Pro $p > 0$, $\text{test}_B(p, v) = 0$, právě když $\sum_{i=0}^{k-1} \text{test}_M(0, p, i, v) = 0$, jinak $\text{test}_B(p, v) = 1$.

$\text{test1}(p)$ představuje počet číslic, pro které byl detekován pozitivní výskyt na pozici p . Tedy, $\text{test1}(p) = \sum_{i=0}^{k-1} \text{test}_B(p, i)$.

$\text{test2}(p, p')$ označíme počet uspořádaných dvojic číslic, pro které je zaznamenán pozitivní výskyt na uvedených pozicích. Jinými slovy

$$\text{test2}(p, p') = \sum_{i=0}^{k-1} \sum_{j=0}^{k-1} \text{test}_M(p, p', i, j).$$

Nyní určíme počet vadných vzorků a jejich číselné označení.

Nechť $T = \max(\text{test1}(p))$.

Pokud $T = 0$, nejsou v souboru žádné vadné vzorky.

Když $T = 1$, pak $\text{test1}(p) = 1$ pro všechna p . Označme X_p takový prvek, pro nějž $\text{test}_B(p, X_p) = 1$. Pak existuje právě jeden vadný vzorek $X = X_{q-1} \dots X_0$.

Když $T = 3$, pak existují právě 3 vadné vzorky. Nechť p je takové, že $\text{test1}(p) = 3$ a $X_{p_1}, X_{p_2}, X_{p_3}$ takové, že $\text{test}_B(p, X_{p_i}) = 1$. Pak pro každé p' různé od p a X_{p_i} existuje právě jedno $X_{p'_i}$ tak, že $\text{test}_M(p, p', X_{p_i}, X_{p'_i}) = 1$. Hledané vadné vzorky pak odpovídají číslům $X_i = X_{(q-1)_i} \dots X_{o_i}$.

Zbývá tak vyřešit případ, kdy $T = 2$. Mohou nastat dva případy, $\max(\text{test2}(p, p')) = 2$ a $\max(\text{test2}(p, p')) = 3$.

Nechť $\max(\text{test2}(p, p')) = 3$. Potom existují právě tři vadné vzorky.

Nechť p, p' je takové, že $\text{test2}(p, p') = 3$ a $X_{p_1}, X_{p_2}, X_{p'_1}, X_{p'_2}$ takové, že $\text{test}_M(p, p', X_{p_1}, X_{p'_1}) = 1$ a $\text{test}_B(p, X_{p_2}) = 1$. Potom pro každé p' různé od p a X_{p_2} existuje právě jedno $X_{p'}$ takové, že $\text{test}_M(p, p', X_{p_2}, X_{p'}) = 1$. Pak jeden hledaný vzorek odpovídá číslu $X = X_{(q-1)_i} \dots X_{(o_i)}$. U dalších dvou vadných vzorků známe číselnici na pozici p (která je u obou vzorků stejná) a na pozici p' , kde se číselnice vadných vzorků liší. Vhodnou kombinací obou vlastností můžeme opět s použitím funkce test_M určit u vadných vzorků číselnici na všech ostatních pozicích.

A nakonec zbývá vyřešit případ, kdy $T = 2$ a $\max(\text{test2}(p, p')) = 2$. Ukážeme, že v tomto případě existují právě dva vadné vzorky, jejich určení pak není náročné.

Předpokládejme, že existují 3 vadné vzorky a p je takové, že $\text{test1}(p) = 2$. Potom jeden vzorek (nazvěme ho A) má na pozici p číselnici v a další dva vzorky (nazvěme je B a C) zde mají číselnici u (různou od v). Protože vzorky B a C jsou různé, musí existovat pozice p' , na které mají rozdílné číselnice. Potom ale $\text{test2}(p, p') = 3$, což je ve sporu s předpokladem $\max(\text{test2}(p, p')) = 2$.

5. Diskuze

Algoritmus jsme představili v podobě, kdy si můžeme volit základ soustavy, ve které číslujeme jednotlivé vzorky. V situaci, kdy počet vzorků není přesně k^q , můžeme doplnit sestavu neutrálními vzorky do nejbližší mocniny čísla k a následně postupovat přesně dle předloženého algoritmu. Algoritmus je samozřejmě neefektivnější v okamžiku, když počet vzorků je mocninou zvoleného základu. V následující tabulce pro mocniny malých

čísel ukazujeme, kolik testů je při různém základu potřeba pro nalezení tří vadných vzorků.

	16	25	27	32	64	81	125	128	243	256	512	625	729	1024
$k = 2$	24			40	60			84		112	144			180
$k = 3$			27			54			90				135	
$k = 4$	16				48					96				160
$k = 5$		25										150		

Z uvedené tabulky vidíme, že v některých případech může být počet testů větší, než počet testovaných vzorků. Potom je samozřejmě efektivnější testovat každý vzorek samostatně. S narůstajícím počtem vzorků se použití navrženého algoritmu, oproti samotnému testování každého vzorku samostatně, stává stále výhodnějším. Vidíme také, že počet potřebných kroků je velmi závislý na bázi, kterou volíme, například pro 625 vzorků je výhodnější použít bázi 3 (a testovat vlastně 729 vzorků), než použít bázi o základu 5, která by se v tomto případě nabízel, ale vyžadovala o 15 testů více.

6. Závěr

Práce s daty a algoritmizace patří mezi základní témata, kterým se má věnovat nově koncipovaný předmět Informatika. V článku jsme představili několik algoritmů, které vychází z problémů praxe a které se uvedeným tématům věnují. Jejich výhodou je, že jsou relativně jednoduché, využívají pouze základní matematické znalosti (především zápis čísla v soustavě o jiném základu), zároveň velmi srozumitelně otevírají otázky spojené se správností a efektivitou algoritmu, které jsou pro rozvoj algoritmického myšlení zásadní.

Literatura

- [1] *Dorfman, R.*: The detection of defective members of large populations. The Annals of Mathematical Statistics, roč. 14 (1943), č. 4, s. 436–440. Dostupné z: <https://www.jstor.org/stable/2235930>
- [2] *Eppstein, D., Goodrich, M. T., Hirschberg, D. S.*: Improved combinatorial group testing algorithms for real-world problem sizes. SIAM Journal on Computing, roč. 36 (2007), č. 5, s. 1360–1375.
- [3] *Shasha, E. D.*: Kybernetické hlavolamy Dr. Ecce. Mladá fronta, Praha, 2005.