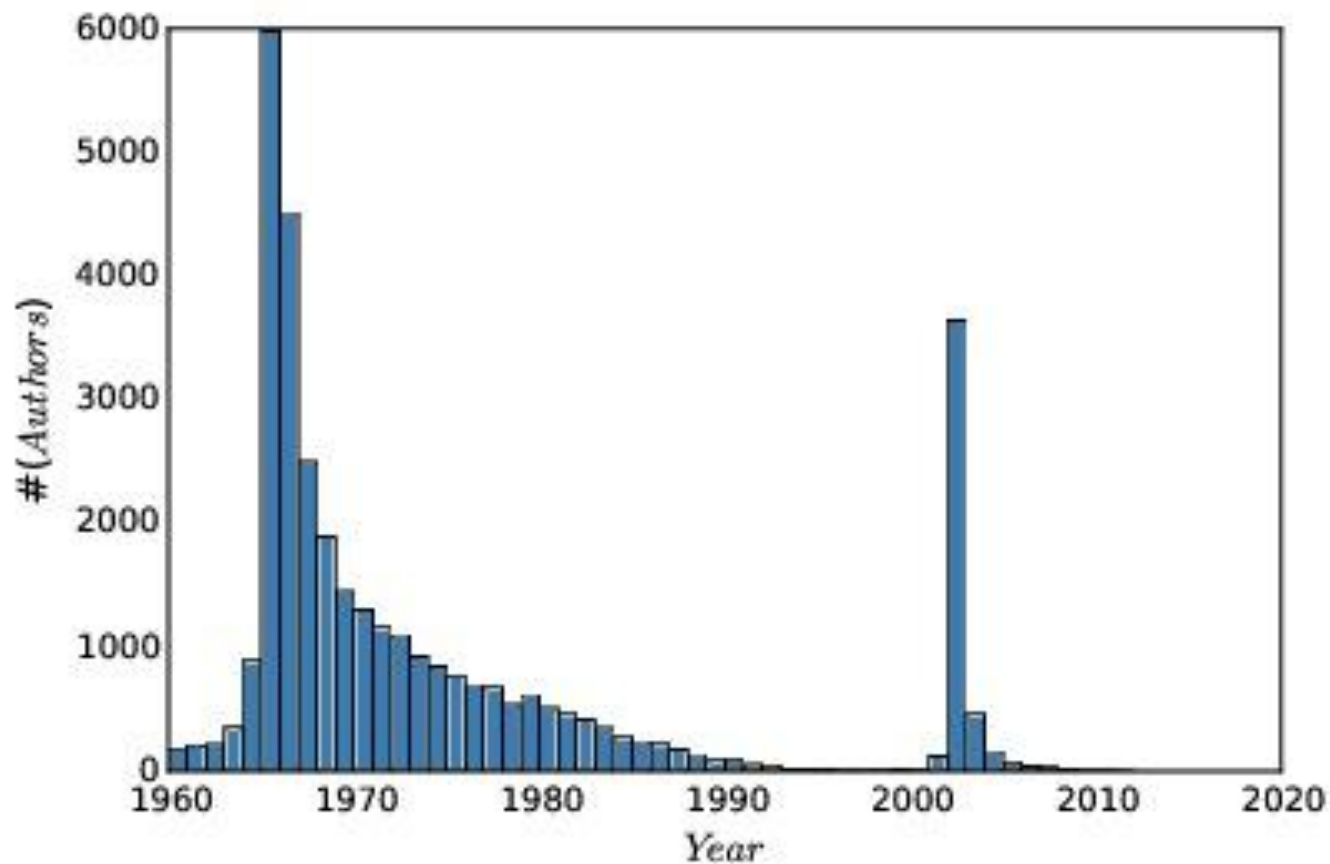
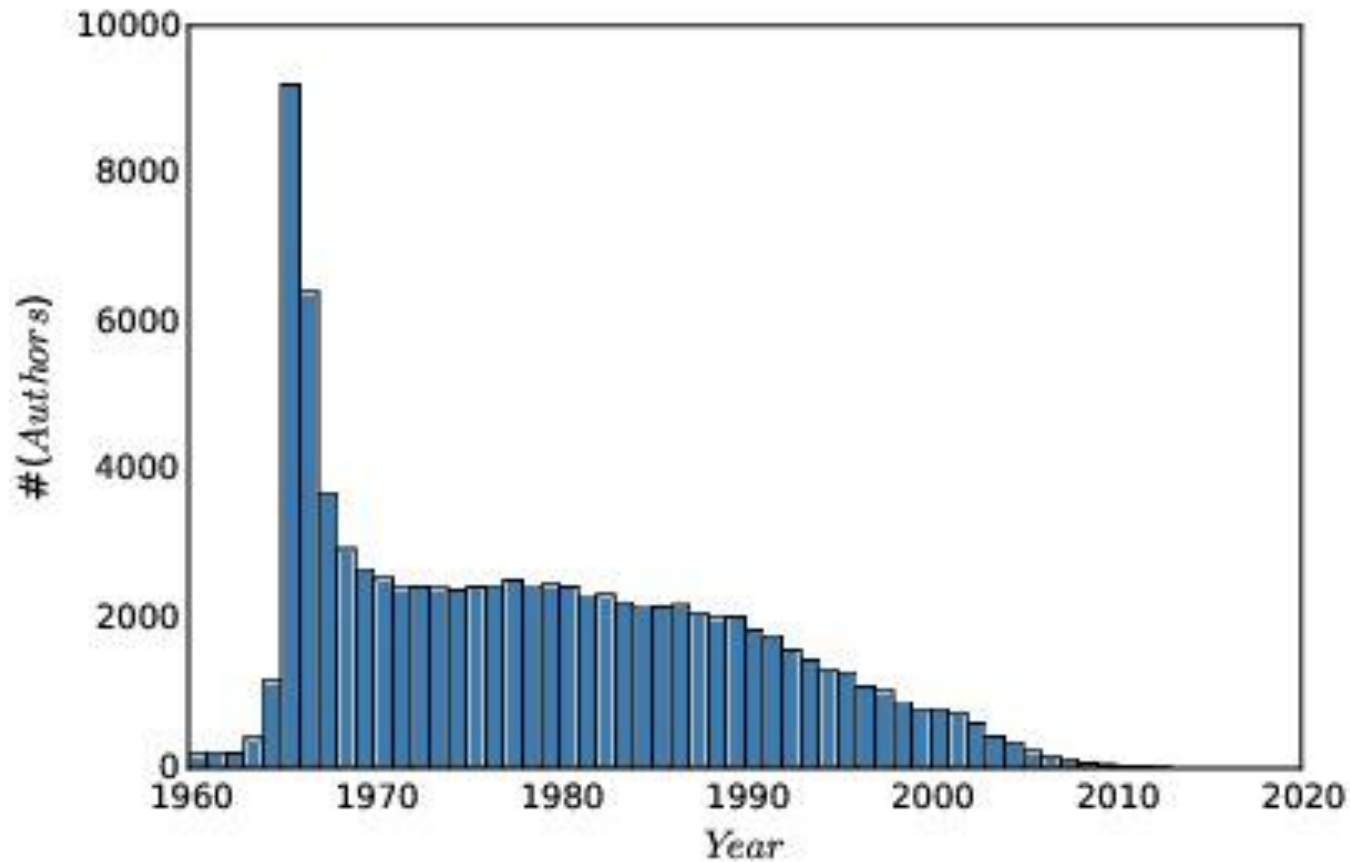


Kedy začalo publikovať 100000 najcitovanejších vedcov?

Kedy začalo publikovať 100000 najcitovanejších vedcov?



Kedy začalo publikovať 100000 najcitovanejších vedcov?

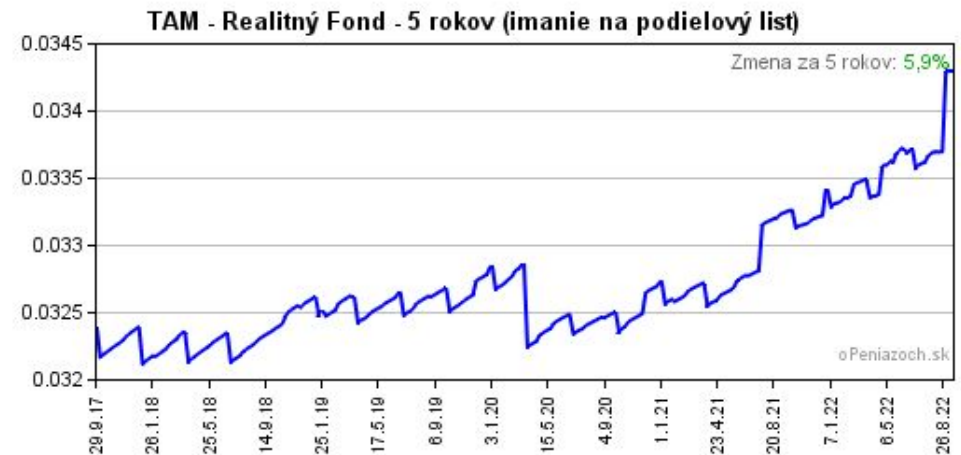
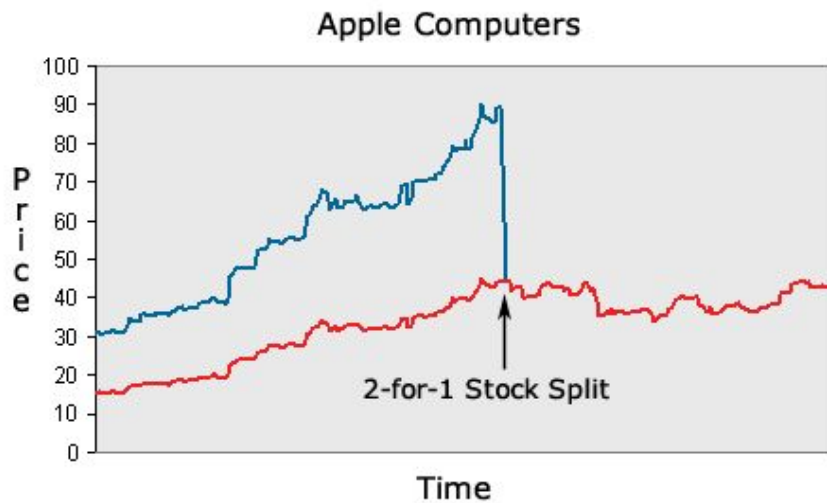


# Problémy s konverziami

- Pri meraných hodnotách, čiarka alebo bodka nemusí byť vždy desatinné číslo (napr. feet,in lb,oz a pod.)
- locale pri veľkých hodnotách (čiarky, bodky, medzery...)
- reprezentácia reťazcov (napríklad - má niekoľko rôznych UTF-8 hodnôt)
  - skús soundex
- časové zóny
- dátumy a staré dátumy
- voľné dni (stock prices vs. počasie)
- kurzy a hodnota peňazí

September 1752						
Su	M	Tu	W	Th	F	Sa
-	-	1	2	14	15	16
17	18	19	20	21	22	23
24	25	26	27	28	29	30

- Česká republika dáva na školstvo oveľa viac peňazí ako Slovensko
- Ceny akcií dlhodobo korelujú s cenami ropy
- Ktorý film zarobil najviac peňazí v histórii? Gone with the Wind \$390 mil.  
Harry Potter (8 filmov) \$7.7 mld.

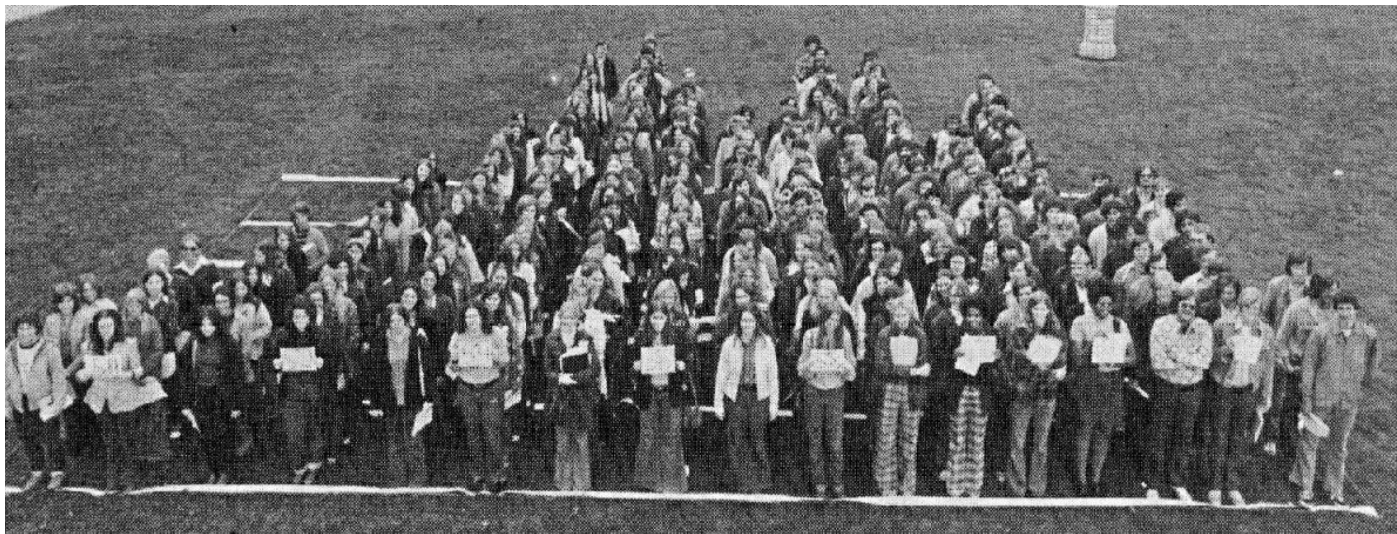
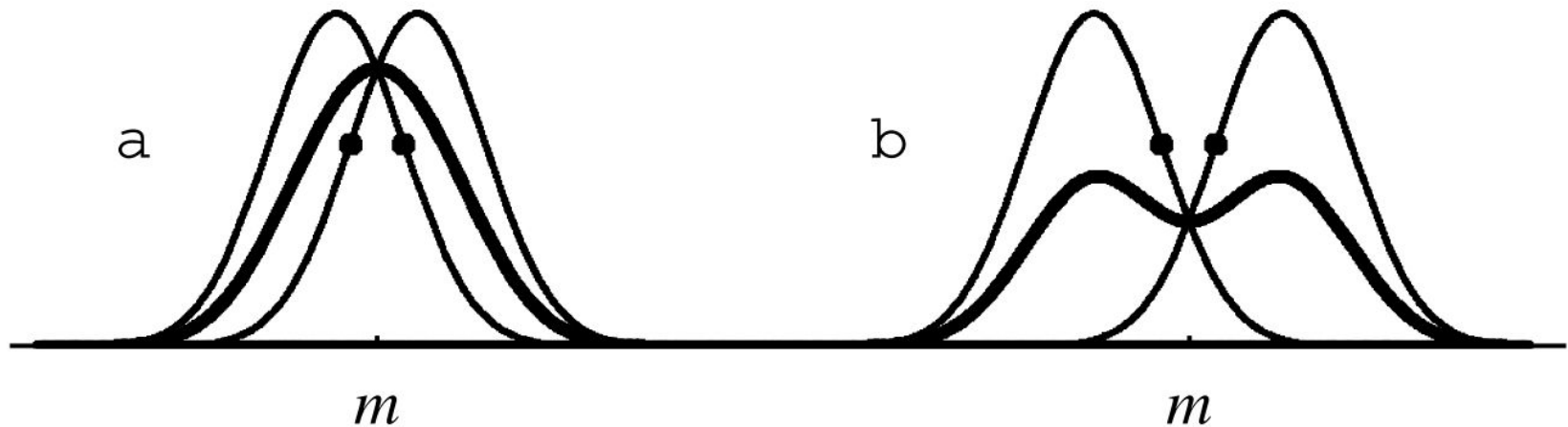


# Bimodálne distribúcie

- Bimodálna distribúcia čohokoľvek je podozrivá
- Môže vzniknúť napr. chybným (alebo chýbajúcim) prepočtom jednotiek
- Z-skóre je celkom dobrá “bezrozmerná” veličina

$$z = \frac{x - \mu}{\sigma}$$

# Výška lidí MÔŽE BYŤ bimodálna



# Chýbajúce pozorovania (ascertainment bias)

- Vedľajšie účinky očkovania sú oveľa horšie ako priebeh COVIDu



Cite as: M. Pavelka *et al.*, *Science*  
10.1126/science.abf9648 (2021).

# The impact of population-wide rapid antigen testing on SARS-CoV-2 prevalence in Slovakia

**Martin Pavelka<sup>1,2,3\*</sup>, Kevin Van-Zandvoort<sup>4,5</sup>, Sam Abbott<sup>4,5</sup>, Katharine Sherratt<sup>4,5</sup>, Marek Majdan<sup>6</sup>, CMMID COVID-19 working group<sup>5</sup>, Inštitút Zdravotných Analýz<sup>2</sup>, Pavol Jarčuška<sup>7</sup>, Marek Krajčí<sup>1</sup>, Stefan Flasche<sup>4,5†</sup>, Sebastian Funk<sup>4,5†</sup>**

<sup>1</sup>Slovak Ministry of Health, Bratislava, Slovakia. <sup>2</sup>Inštitút Zdravotných Analýz (Institute of Health Analyses), Bratislava, Slovakia. <sup>3</sup>Department of Public Health and Policy, London School of Hygiene and Tropical Medicine, London, UK. <sup>4</sup>Department for Infectious Disease Epidemiology, London School of Hygiene and Tropical Medicine, London, UK. <sup>5</sup>Centre for Mathematical Modelling of Infectious Diseases, London School of Hygiene and Tropical Medicine, London, UK. <sup>6</sup>Institute for Global Health and Epidemiology, Faculty of Health Sciences and Social Work, Trnava University, Trnava, Slovakia. <sup>7</sup>Faculty of Medicine, Pavol Jozef Šafárik University, Košice, Slovakia.

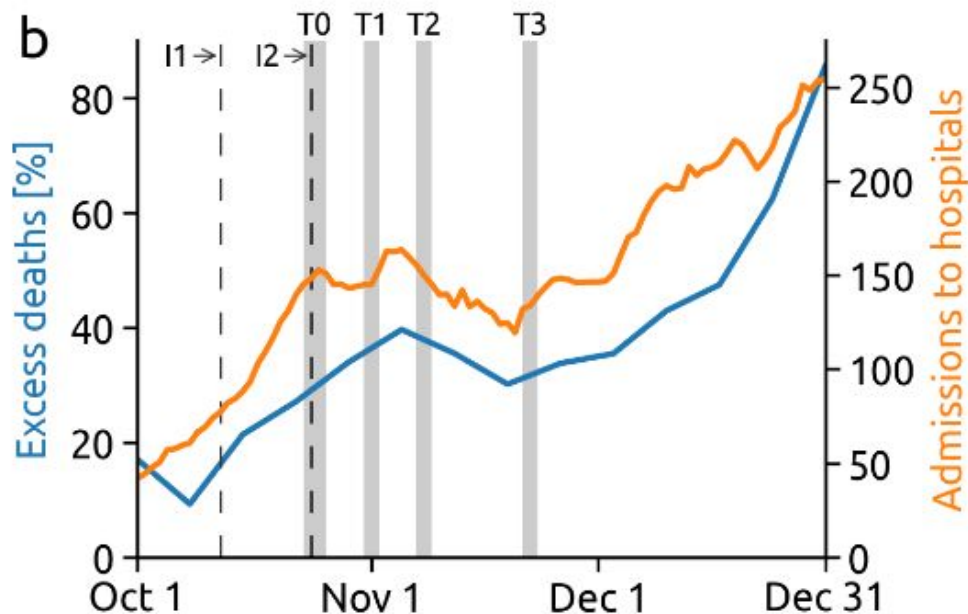
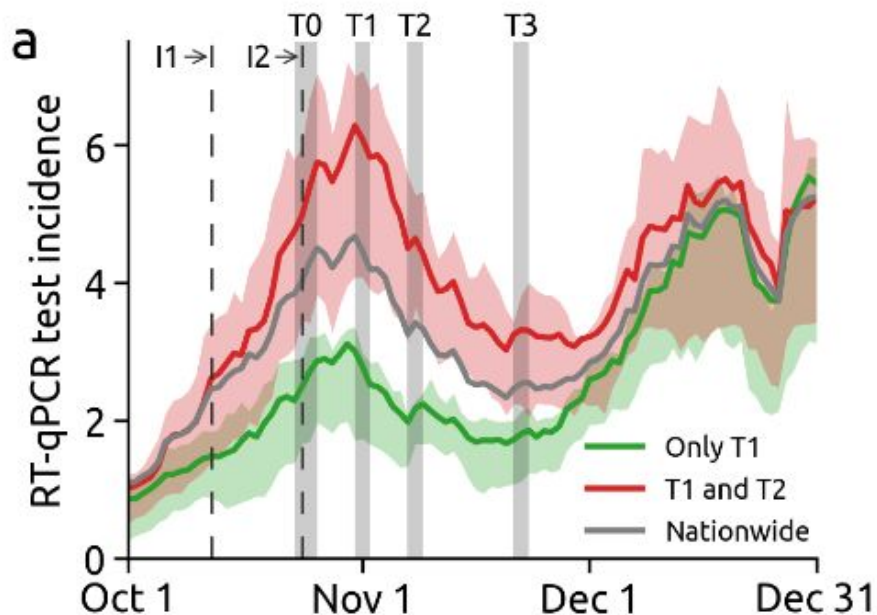
\*Corresponding author. Email: martin.pavelka@health.gov.sk

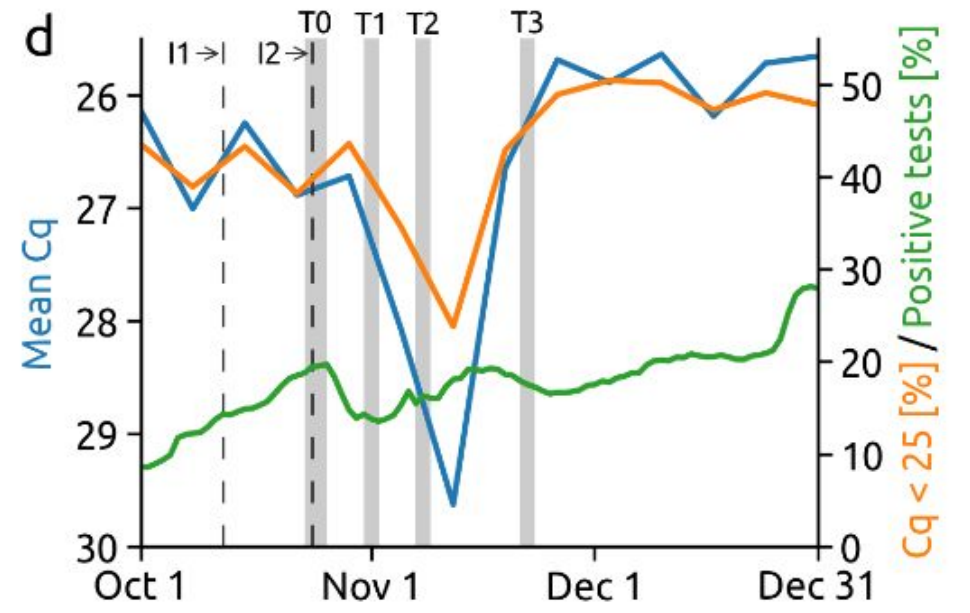
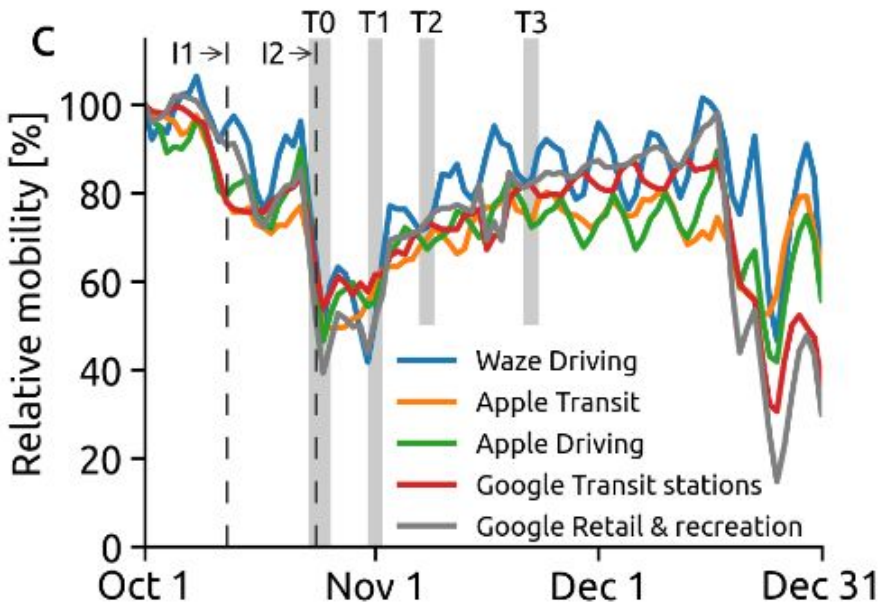
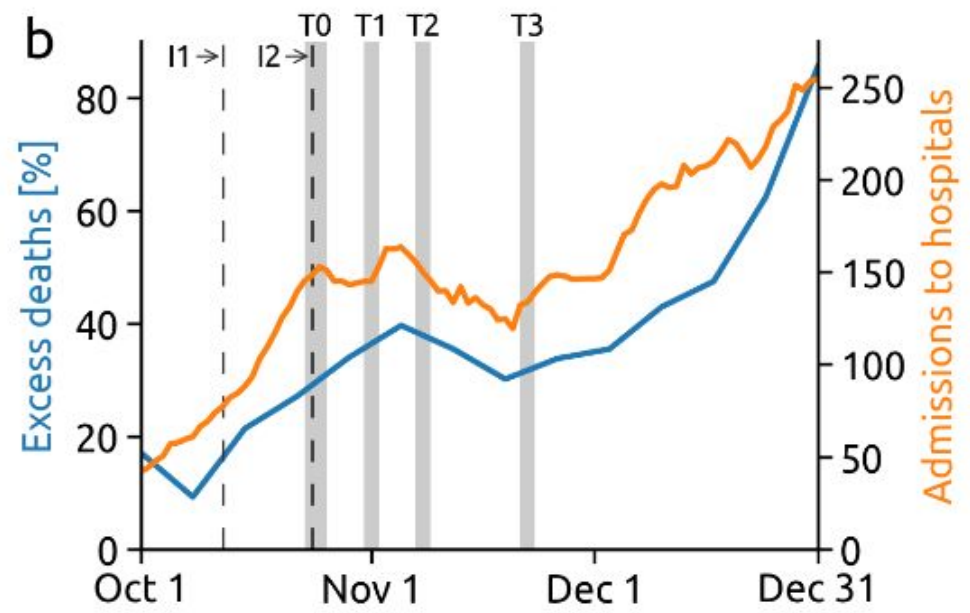
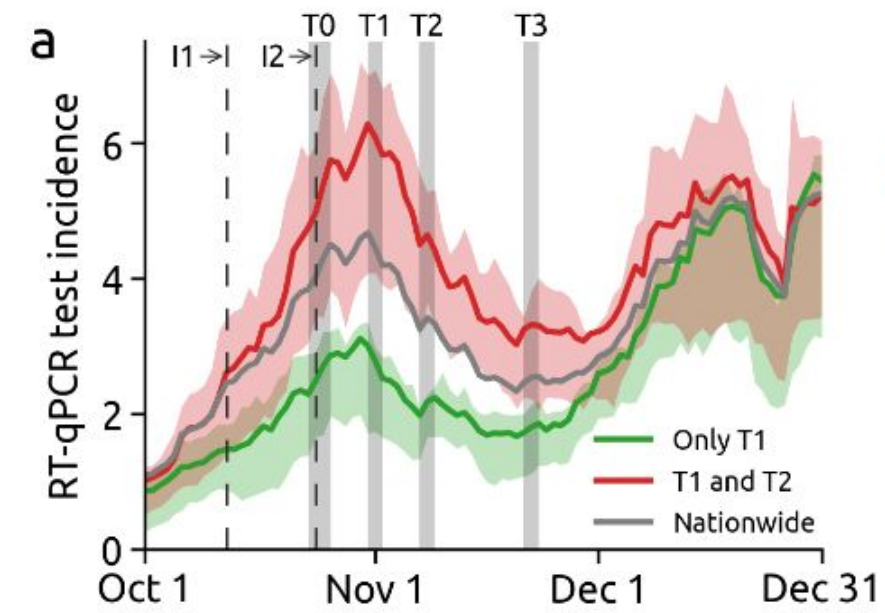
†These authors contributed equally to this work.

**Slovakia conducted multiple rounds of population-wide rapid antigen testing for SARS-CoV-2 in late 2020, combined with a period of additional contact restrictions. Observed prevalence decreased by 58% (95% CI: 57-58%) within one week in the 45 counties that were subject to two rounds of mass testing, an estimate that remained robust when adjusting for multiple potential confounders. Adjusting for epidemic growth of 4.4% (1.1-6.9%) per day preceding the mass testing campaign, the estimated decrease in prevalence compared to a scenario of unmitigated growth was 70% (67-73%). Modelling indicated that this decrease could not be explained solely by infection control measures, but required the additional impact of isolation and quarantine of household members of those testing positive.**

# Výskyt ochorenia sa znížil o 58% medzi prvým (T1) a druhým kolom (T2) plošného testovania

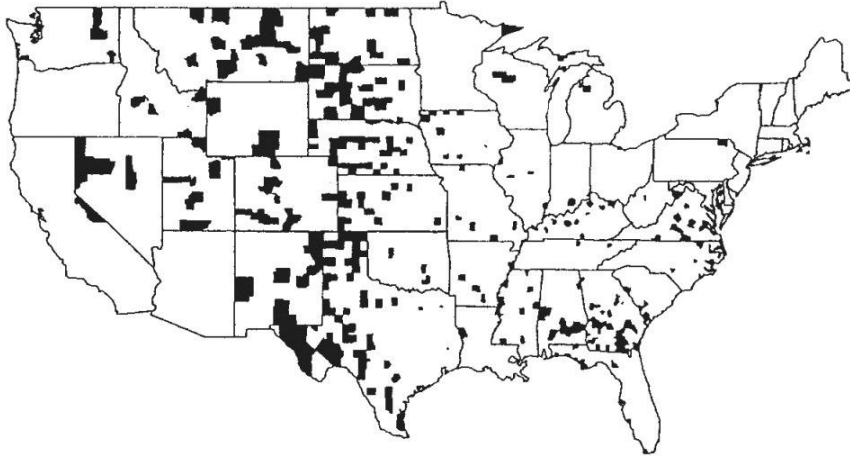
- V prvom kole boli testovaní všetci, ktorí prišli
- V druhom kole (o týždeň neskôr) boli testovaní všetci, ktorí prišli, **okrem**:
  - tých, čo boli pozitívni v prvom kole
  - ich rodinných príslušníkov a úzkych kontaktov (už boli v karanténe)
  - tých, čo nechceli riskovať dvojtýždňové “domáce väzenie” pre celú rodinu



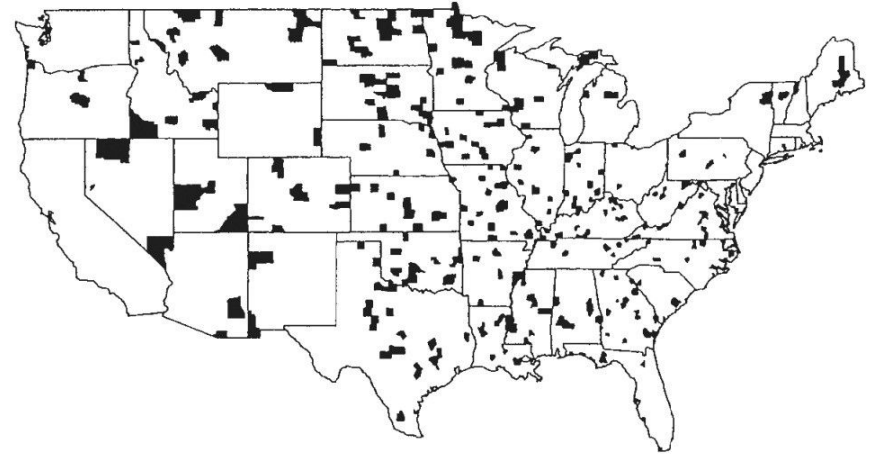


# Odláhlé hodnoty (outliers)

- Najväčší evidovaný stavec dinosaura je o 50% väčší ako všetky ostatné nájdené stavce dinosaurov
  - skutočnosť alebo artefakt?
  - BTW, príslušné múzeum už nevie lokalizovať príslušný exponát
- Na minimá, maximá, hodnoty s príliš veľa štandardnými odchýlkami od priemeru sa treba pozerať kriticky
- Zistiť **prečo** máme odláhlé hodnoty, nie ich jednoducho mazať
  - Zmazanie odláhlých hodnôt môže viesť k lepšiemu modelu (napr. neobvyklé chyby merania)
  - Zmazanie odláhlých hodnôt môže viesť k horšiemu modelu (napr. príliš jednoduchý model dát)



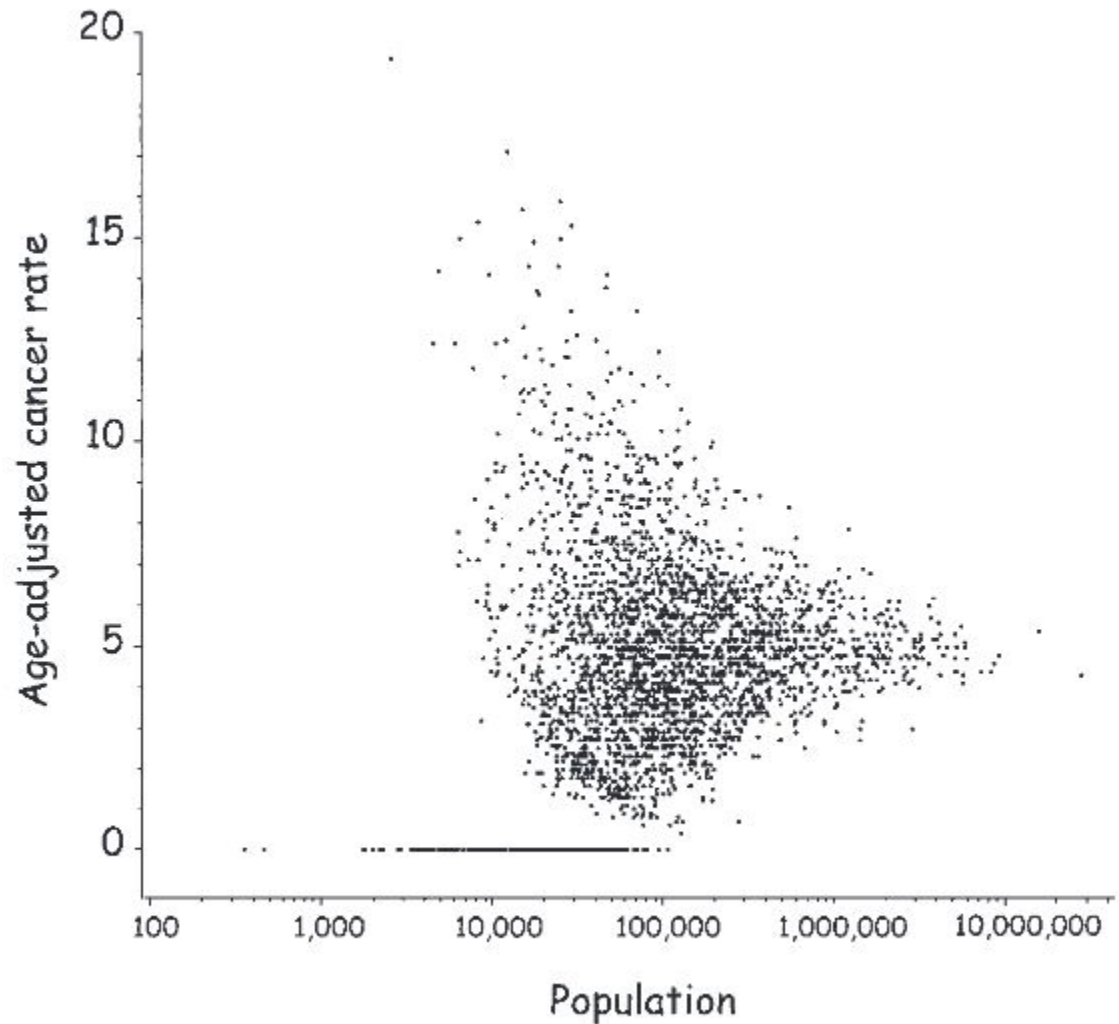
Counties with the **LOWEST**  
kidney cancer death rates  
(1980-1989)



Counties with the **HIGHEST**  
kidney cancer death rates  
(1980-1989)

Gelman, Nolan: Teaching Statistics

**Figure 1.3.**  
Age-adjusted kidney cancer rates for all U.S. counties in 1980–1989 shown as a function of the log of the county population.



Howard Wainer: Picturing the  
Uncertain World, 2009

# Chýbajúce dáta

- Väčšina metód sa nedokáže vyrovnat' s chýbajúcimi dátami (napr. rok úmrtia osoby, ktorá ešte žije, chýbajúce meranie, zjavne nezmyselná hodnota, ...)
- **Riešenie 1:** zahod' záznamy, ktoré majú chýbajúce dáta  
ak máme veľa atribútov, môže sa stať, že takmer v každom zázname bude niečo chýbať
- **Riešenie 2: imputácia chýbajúcich hodnôt**
  - priemer hodnôt daného atribútu
  - náhodná existujúca hodnota daného atribútu
  - hodnota získaná predikciou z iných atribútov (napr. lineárnou regresiou)

Ovplyňuje zvolený spôsob imputácie výsledky analýzy?

## Zhrnutie: Čistenie dát je veľmi dôležitým (a zložitým) krokom pri ich analýze

- Chyby v dátach často príčinou chybných analýz
- Náhodné chyby (napr. chyby merania) sú obvykle OK
  - pozor na početnosť hodnôt / štandardné odchýlky!
- Artefakty (systematické chyby) robia problémy:
  - chyby unifikácie dát
  - nepozorovateľné hodnoty (ascertainment bias)
  - odľahlé hodnoty (outliers)
- Chýbajúce hodnoty