

# Parametrické testy pro dva výběry

B02907 Informační a komunikační technologie



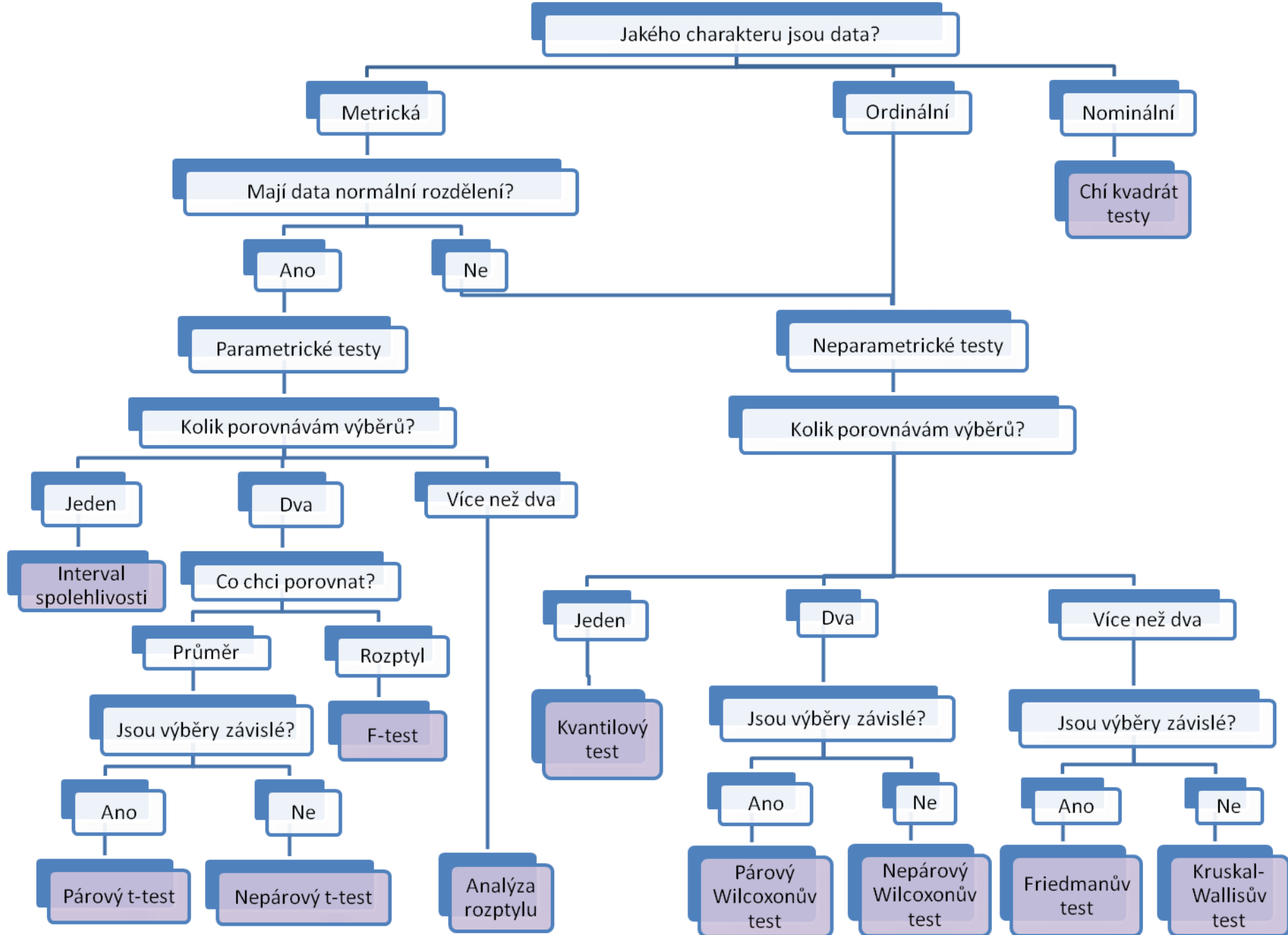
Lubomír Štěpánek,  
Ústav biofyziky a informatiky  
1. LF UK



# Upozornění!

- dole v poznámkách jsou u většiny snímků rozšiřující a vysvětlující komentáře
- u některých statistických metod budete odkazováni na statistické tabulky, které jsou volně přístupné online na adrese <http://new.euromise.org/czech/tajne/ucebnice/html/html/node15.html>
- (obvykle bude ještě na příslušném snímku odkaz zopakován; autor vynaložil značné úsilí, aby se symbolika v prezentacích shodovala se symbolikou v tabulkách, proto by neměla být orientace v tabulkách problémem)
- z předložených prezentací se můžete učit, můžete je kopírovat či jinak měnit, ale bez dovození autora/autorů je nesmíte použít do svých publikací ☺
- předložené prezentace nejsou bezchybnou statistickou kuchařkou, proto ne zcela doporučuji se na ně ve svých pracích odkazovat, nebo je dokonce citovat ☺
- pokud se budu sám odkazovat na vhodnou literaturu, myslím tím nejspíše následující dvě knihy:
  - Zvára: Biostatistika. Karolinum, Praha 1988
  - Zvárová et al.: Biomedicínská statistika I. Základy statistiky pro biomedicínské obory
- dotazy a konzultace možné a vlastně i doporučeny

(Lubomír Štěpánek, stepanek.lub@seznam.cz)

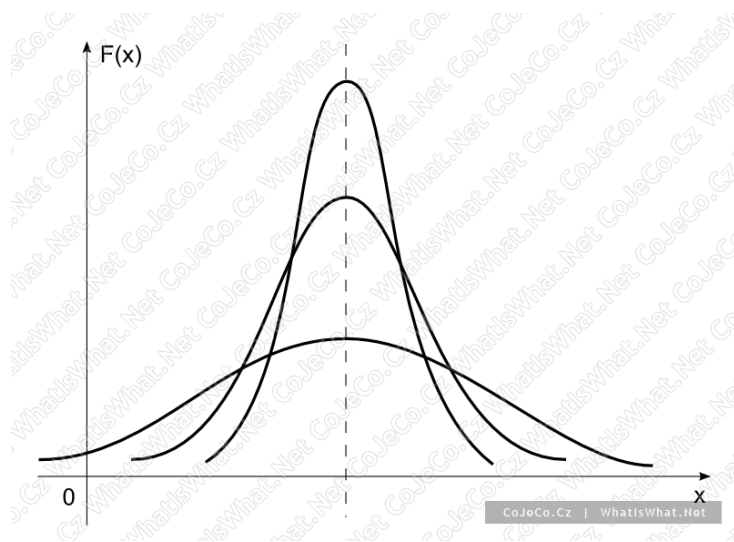


# Jakého charakteru jsou data?

- metrická data
  - vyjádřena reálnými čísly a jednotkami
  - souvisle nebo nesouvisle pokrývají číselnou osu, hodnoty lze mezi sebou porovnat ve smyslu kolikrát větší/menší
  - např. *hodinová mzda (Kč), tělesná výška (cm), počet dětí v rodině (jednotkou „1“)*
- ordinální data
  - vyjádřena celými čísly bez jednotky nebo slovně
  - jednotlivé hodnoty lze mezi sebou porovnat ve smyslu menší X větší, lepší X horší atd., ale ne jinak
  - např. bolest *malá, střední, intenzivní*; grade nádoru *I, II, III*; pořadí závodníka *1., 2., 3., ...* (nebo *první, druhý, třetí*)
- nominální (alternativní) data
  - vyjádřena slovně, někdy čísly, ale ty mají pouze symbolický charakter
  - hodnoty jsou vůči sobě neporovnatelné (rovnocenné)
  - např. kraje ČR *{Jihočeský, ..., Zlínský}*, pohlaví *{muž, žena}*, interleukiny *{1, 2, 3, ...}*

# Mají data normální rozdělení?

- metrická data, jejichž histogram četností pro jednotlivé hodnoty znaku je Gaussova křivka
- symetrická, popisujeme průměr a odchylku
- splňují podmínky pro veličinu *šikmost* a špičatost



# Kolik porovnáváme výběrů?

- (jeden výběr)
  - výběrové charakteristiky (míra polohy a variability)
    - popisná statistika, populační odhady (intervaly spolehlivosti)
  - *např. průměrný věk úmrtí obyvatel žijících za polárním kruhem je 75 let*
- dva výběry (ev. více)
  - porovnávání výběrových statistik mezi soubory
  - *např. je průměrný věk úmrtí obyvatel žijících za polárním kruhem vyšší než na rovníku?*

# Dva a více výběrů

- pro porovnání výběrových charakteristik (průměru, odchylky) nutná *hypotéza*
- hypotéza říká, co porovnáváme (1) a jaké očekáváme výsledky (2)
  - např. „**průměrný věk úmrtí** (1) *obyvatel žijících za polárním kruhem* je **vyšší než** (2) *na rovníku*“
- hypotéza buďto platí (pravděpodobnost 1), nebo neplatí (0), my můžeme skutečnost s určitou pravděpodobností odhalit

# Hypotéza

- nulová hypotéza  $H_0$  – porovnávané charakteristiky výběrů nejsou signifikantně rozdílné
  - „oba výběry jsou vlastně shodné, jsou jakoby vybrány z jedné a téže populace“
- alternativní hypotéza  $H_1$  – porovnávané charakteristiky výběrů jsou signifikantně rozdílné
  - „výběry nemohou být shodné, je-li v rozdíl v některé jejich charakteristice“



# Hypotéza

- *oboustranná*

- „výběry mají rozdílné... (zkoumaná vlastnost)“
- těžší prokázat
- ale nevyžaduje předchozí argumentaci, nic dopředu nepředpokládáme

- *jednostranná*

- „jeden výběr má menší/větší ... než druhý...“
- lehčí prokázat
- ale vyžaduje dopředu argumentaci, proč můžeme vyloučit opačnou nerovnost

# Chyba prvního a druhého typu

- chyba prvního typu (pravděpodobnost alfa)
  - pravděpodobnost, že za platnou přijmeme na základě statistické metody alternativní hypotézu, která ve skutečnosti neplatí
- chyba druhého typu (pravděpodobnost beta)
  - pravděpodobnost, že zamítneme statistickou metodou alternativní hypotézu, která ve skutečnosti platí

# Co chci porovnat?

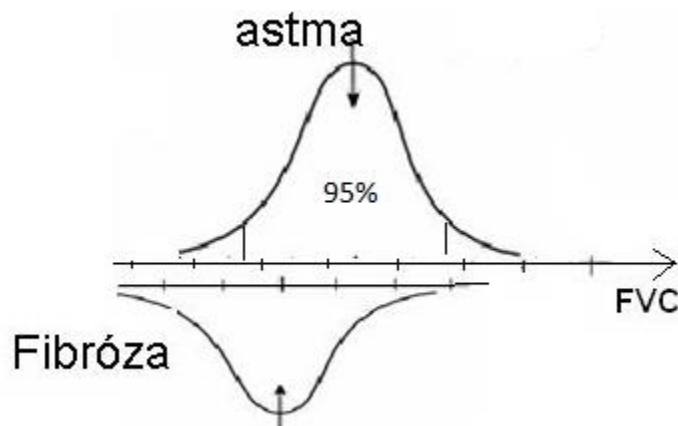
- míra polohy
  - průměr, medián
  - např. „je průměrná výška 18letých chlapců větší než průměrná výška 18letých dívek?“ (porovnání průměrů), „je střední doba přežití po infarktu myokardu delší než po cévní mozkové příhodě?“ (porovnání mediánů), atd.
- míra variability
  - rozptyl (směrodatná odchylka), kvantily
  - např. „je přesnost měření jednoho přístroje lepší než přesnost druhého?“, „je hmotnost 95% percentilu mužů vyšší než 95% percentilu žen?“ atd.

# Jsou výběry spárované?

- dva či více výběrů mohou být na sobě nezávislé (nespárované), nebo závislé (spárované)
- nespárované – obvykle se jedná o výběry zcela různých prvků, často mají výběry různou velikost
  - např. pacienti s diabetem X pacienti zdraví atd.
- spárované – prvky v různých situacích, proto jsou na sobě závislé, všechny výběry mají stejnou velikost
  - např. pacienti před léčbou X ti samí po léčbě atd.

# Jak myslí parametrický test

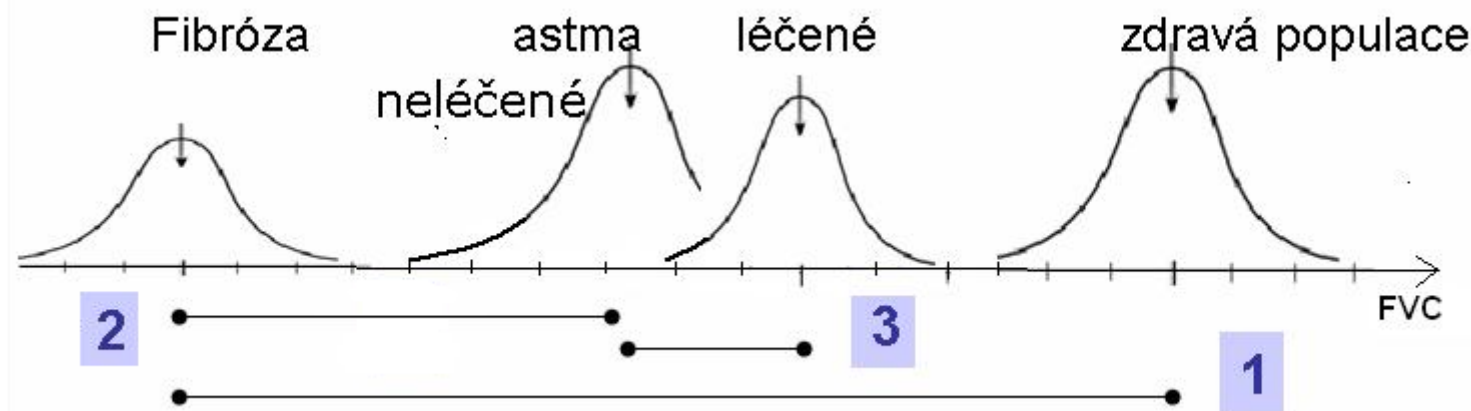
- máme dva výběry, oba s normálním rozdělením dat
- dle  $H_0$  se předpokládá, že oba vybrány ze stejné populace
- jeden výběr lze porovnávat vůči druhému, který jakoby tvoří prvnímu jeho „populaci“, z níž byl vybrán
- pak musí mít s 5% chybou stejný průměr (t-test) a odchylku (F-test)



# Parametrický test hypotézy

- máme-li alespoň dva výběry, které vykazují normální rozdělení
- máme-li hypotézu
- máme-li zvolenou hladinu významnosti
- → pak můžeme provést tzv. parametrický test hypotézy

# Porovnání průměrů



## Situace:

(0. Srovnání výběru s populací (normou) – u-test, interval spolehlivosti

*Příklad:* hodnoty FVC u zdravých (normální populace) a fibrózy)

**1. Srovnání dvou výběrů (nezávislé, různé osoby) – t-test**

*Příklad:* hodnoty FVC astmatiků ve srovnání s fibrózou

hmotnost novorozenců matky N - D

**2. Srovnání dvou výběrů (závislé, tytéž v různých situacích) – párový t-test**

*Příklad:* FVC u skupiny astmatiků před léčbou a po ní

# Nepárový (dvoubýběrový) t-test

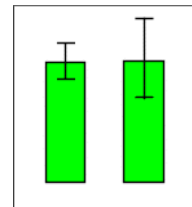
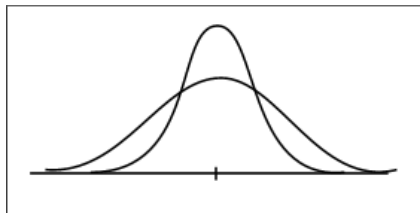
- požadujeme metrická data, normální rozdělení dat, dva nezávislé výběry se stejným nebo podobným rozptylem (shodnost rozptylů lze obejít)
- porovnává průměr znaku zkoumaného u dvou nezávislých výběrů
- $H_0$ : průměry v obou výběrech jsou shodné,  $H_1$ : průměry v obou výběrech jsou odlišné
- $H_0$  např. „průměrná tělesná výška maturantů a maturantek je shodná“ „průměrné FVC astmatiků a zdravých jedinců je shodné“ atd. proti odpovídajícím  $H_1$
- tabulkový procesor Excel: vložit funkci -> TTEST
  - parametry: Pole1 – množina dat prvního výběru, Pole2 – množina dat druhého výběru, strany – 1 pro jednostrannou, 2 pro oboustrannou alternativní hypotézu (tu doporučuji), Typ – 2 pro nepárové výběry se stejným rozptylem, 3 pro nepárové výběry s různým rozptylem
  - výsledek: hladina významnosti  $p$ ; je-li  $p \leq 0,05$ , přijímáme alternativní hypotézu  $H_1$ , jinak přijímáme  $H_0$
- shodnost rozptylů ověříme F-testem (viz dále)
- R-Project: `t.test(x, y, alternative = "two.sided", mu = 0, var.equal = FALSE, conf.level = 0.95)`



# Párový t-test

- požadujeme metrická data, normální rozdělení dat, dva závislé výběry
- porovnává průměr znaku zkoumaného u dvou závislých výběrů
- $H_0$ : průměry v obou výběrech jsou shodné,  $H_1$ : průměry v obou výběrech jsou odlišné
- $H_0$  např. „průměrná tělesná hmotnost rodiček před porodem a po porodu je shodná“, „měření stejného fenoménu oběma metodami dává průměrně stejné výsledky“, „průměrný tlak hypertoniků před léčbou a po léčbě je srovnatelný“ atd. proti odpovídajícím  $H_1$
- matematický algoritmus viz vhodná literatura
- tabulkový procesor Excel: vložit funkci -> TTEST
  - parametry: Pole1 – množina dat prvního výběru, Pole2 – množina dat druhého výběru, strany – 1 pro jednostrannou, 2 pro oboustrannou alternativní hypotézu (tu doporučuji), Typ – 1 pro spárované výběry
  - výsledek: hladina významnosti  $p$ ; je-li  $p \leq 0,05$ , přijímáme alternativní hypotézu  $H_1$ , jinak přijímáme nulovou hypotézu  $H_0$
- R-Project: `t.test(x, y, mu=0, alternative = "two.sided", paired= TRUE)`

# Porovnání rozptylů



*Grafická reprezentace*

## Srovnání rozdílu rozptylů – F-test

### Situace:

- Srovnání přesnosti měření

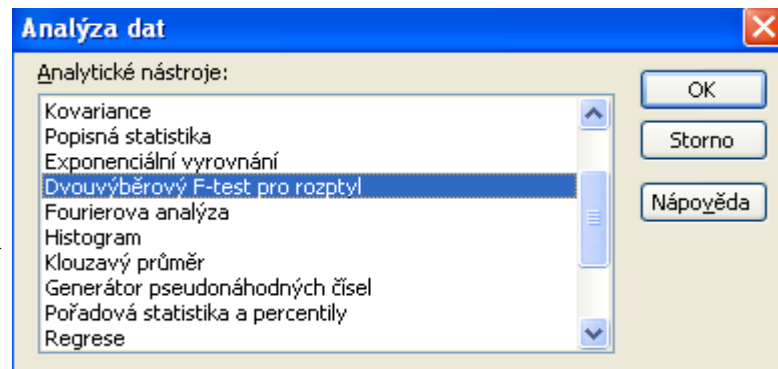
*Př: Tatáž sada krevních vzorků vyšetřena v různých laboratořích*

- Vyrovnanost účinku léků

*Př: V průměru shodný efekt, ale různá variabilita*

*? Proč u někoho výborný efekt, u někoho slabý*

*volba F-testu v Excelu,* →  
*je-li přítomna Analýza dat, jinak přes FTEST*

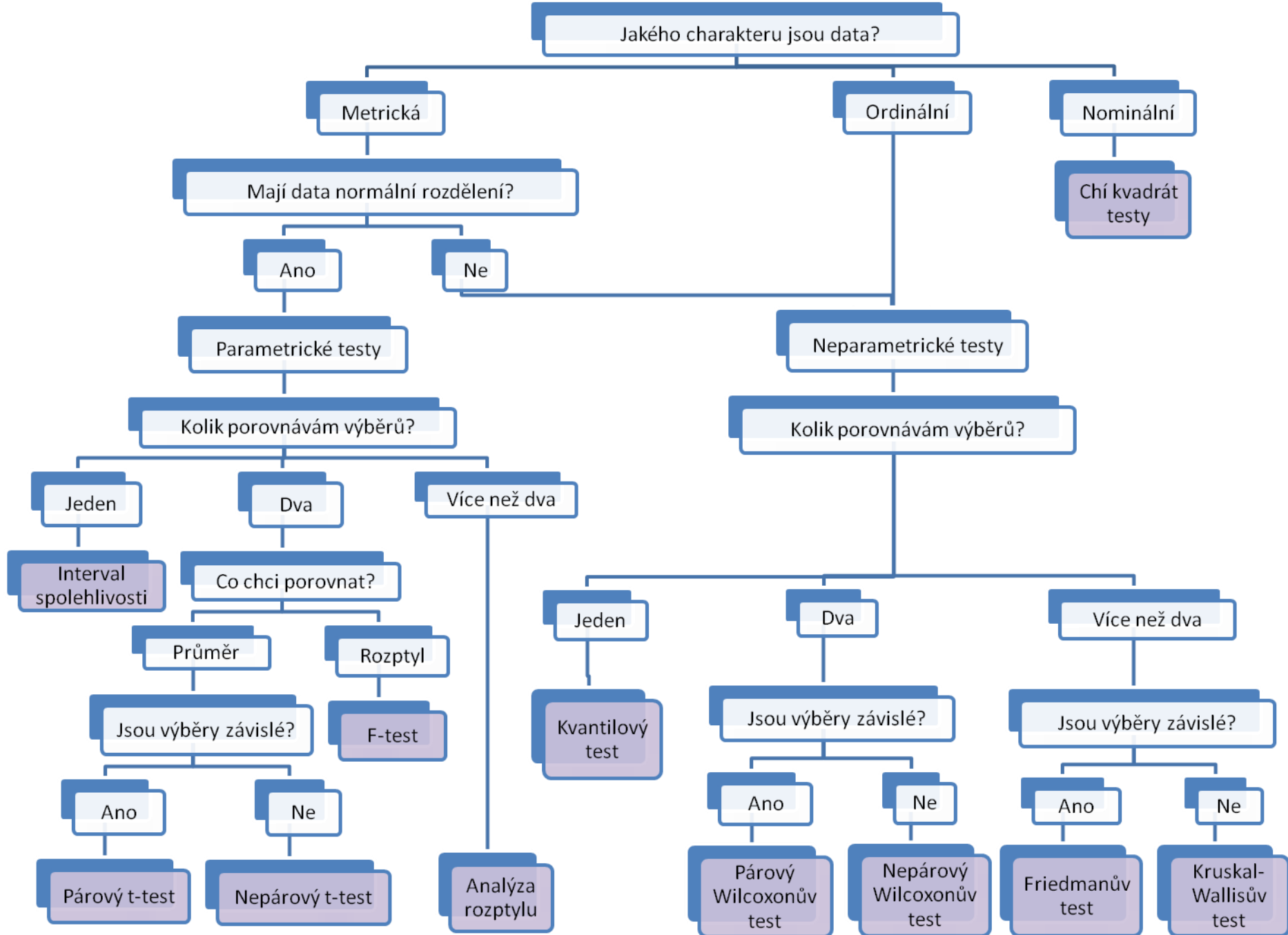


# F-test

- požadujeme metrická data, normální rozdělení dat, dva nezávislé výběry
- porovnává rozptyl znaku zkoumaného na obou výběrech
- $H_0$ : rozptyly v obou výběrech jsou shodné,  $H_1$ : rozptyly v obou výběrech jsou odlišné
- $H_0$  např. „přesnost obou střelců při střelbě na terč je shodná“, „obě metody měří ve stejné šíři“ atd. proti odpovídajícím  $H_1$
- matematický algoritmus viz vhodná literatura
- tabulkový procesor Excel: vložit funkci -> FTEST
  - parametry: Pole1 – množina dat prvního výběru, Pole2 – množina dat druhého výběru
  - výsledek: hladina významnosti  $p$ ; je-li  $p \leq 0,05$ , přijímáme alternativní hypotézu  $H_1$ , jinak přijímáme nulovou hypotézu  $H_0$
- F-test se neužívá tak často jako t-test, ale je vhodný pro určení shodnosti/rozdílnosti rozptylů dvou výběrů, které chceme porovnat t-testem (abychom věděli, který typ t-testu použít)
- R-Project: `var.test(x, y, ratio = 1, alternative = "two.sided", conf.level = 0.95)`

# Poznávací znamení podle výchozího formátu dat

- nepárový t-test
  - data (čísla) jsou v nesymetrické tabulce  $2 \times \max(n_1, n_2)$
- párový t-test
  - data jsou v symetrické tabulce  $2 \times n$
- F-test
  - stejné jako u t-testu



lubomir.stepanek@lf1.cuni.cz  
lubomir.stepanek@fbmi.cvut.cz