

NORMÁLNÍ ROZLOŽENÍ, ZÁKLADY TESTOVÁNÍ HYPOTÉZ A STATISTICKÁ INFERENCE.

ZUR357 Statistická analýza dat --

9. listopadu 2017

AUTOMATIC RECODE

Transformace string values do numerických hodnot

CATEGORIZE VARIABLES

Kategorizuje kardinální znaky podle percentilů tak, aby každá kategorie obsahovala přibližně stejný počet případů

COUNT VALUES

Sčítání výskytu určité hodnoty

Např. Které tituly pravidelně čtete (0=nečte, 1 čte)

- Respekt, Reflex, Euro, Týden.....

Kolik strategií používají děti, aby se vypořádali s obtěžc

140. Again, thinking about this time, did you do any of these things?

PLEASE TICK AS MANY BOXES AS NEEDED

- | | | |
|---|--|--------------------------|
| A | I stopped using the internet for a while | <input type="checkbox"/> |
| B | I deleted any messages from the person who sent it to me | <input type="checkbox"/> |
| C | I changed my filter/contact settings | <input type="checkbox"/> |
| D | I blocked the person who had sent it to me | <input type="checkbox"/> |
| E | I reported the problem (eg clicked on a 'report abuse' button, contact an internet advisor or "Internet service provider (ISP)") | <input type="checkbox"/> |
| F | None of these things | <input type="checkbox"/> |
| G | Don't know | <input type="checkbox"/> |



RANK CASES

Vytvoří novou proměnnou, která určuje pořadí případu podle původní proměnné

COMPUTE

Libovolné operace zadané strukturovaným příkazem

Podmínka IF...

Změna rozložení

Centrování

Sumační indexy

Komplexní operace

144. In the PAST 12 MONTHS, how often, have these things happened to you?

PLEASE TICK ONE BOX ON EVERY LINE

		Very often	Fairly often	Not very often	Never/ almost never
A	I have gone without eating or sleeping because of the internet	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
B	I have felt bothered when I cannot be on the internet	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
C	I have caught myself surfing when I'm not really interested.....	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
D	I have spent less time than I should with either family, friends or doing schoolwork because of the time I spent on the internet.....	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
E	I have tried unsuccessfully to spend less time on the internet	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>

PRÁCE S PODSOUBORY

Procedura SELECT CASES

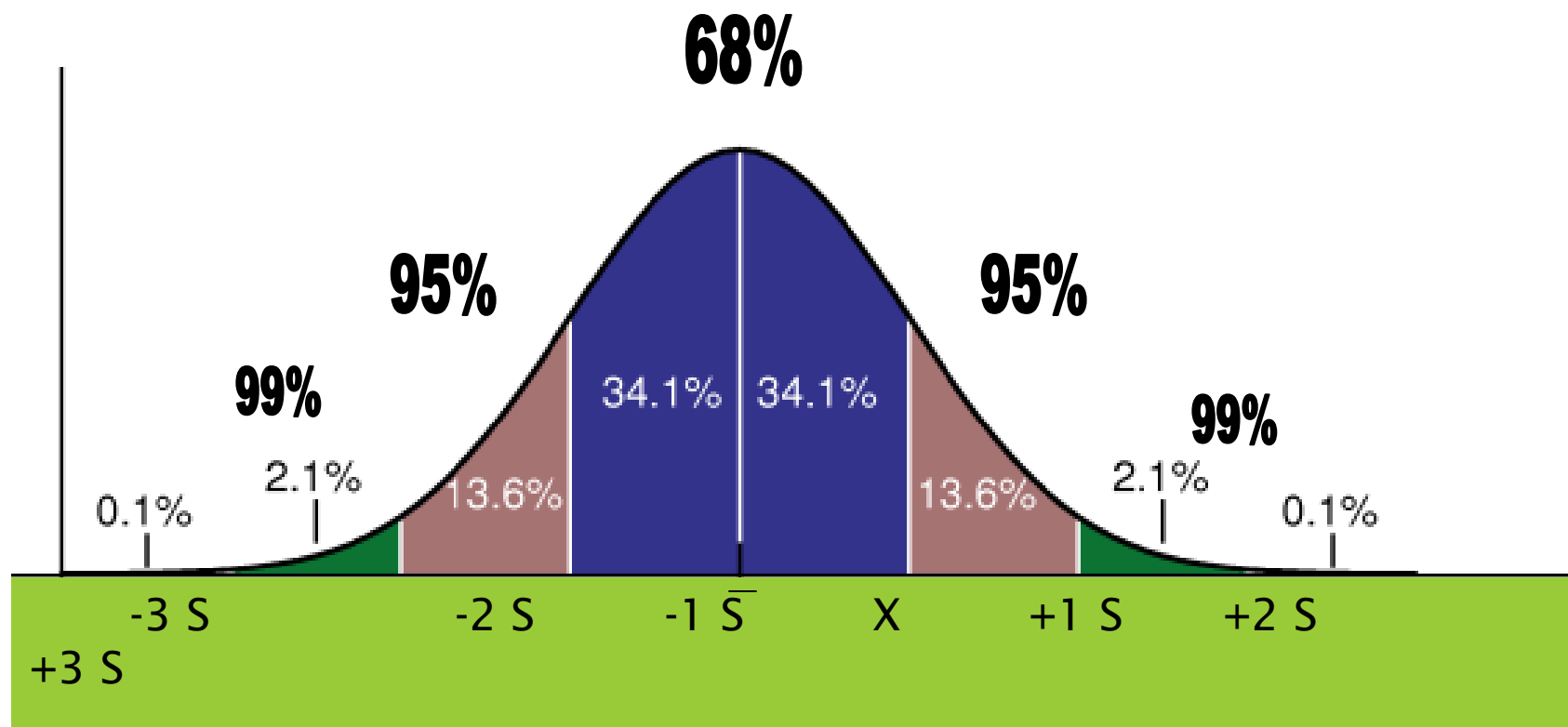
- Náhodně vybrané případy - redukce souboru (Random Sample of Cases)
- Výběr vedený výzkumnou otázkou (Select IF)

ÚKOL 1. Ověřte H1: V Lidových novinách vystupuje v roli hlavního aktéra častěji politik než odborník na dané téma.

Procedura SPLIT FILE

ÚKOL 2. Porovnejte výskyt odborníků a politiků jako aktérů v jednotlivých denících (H2: MFD bude využívat odborníků častěji než LN.)

NORMÁLNÍ ROZLOŽENÍ (NORMAL DISTRIBUTION; GAUSSOVA KŘIVKA)

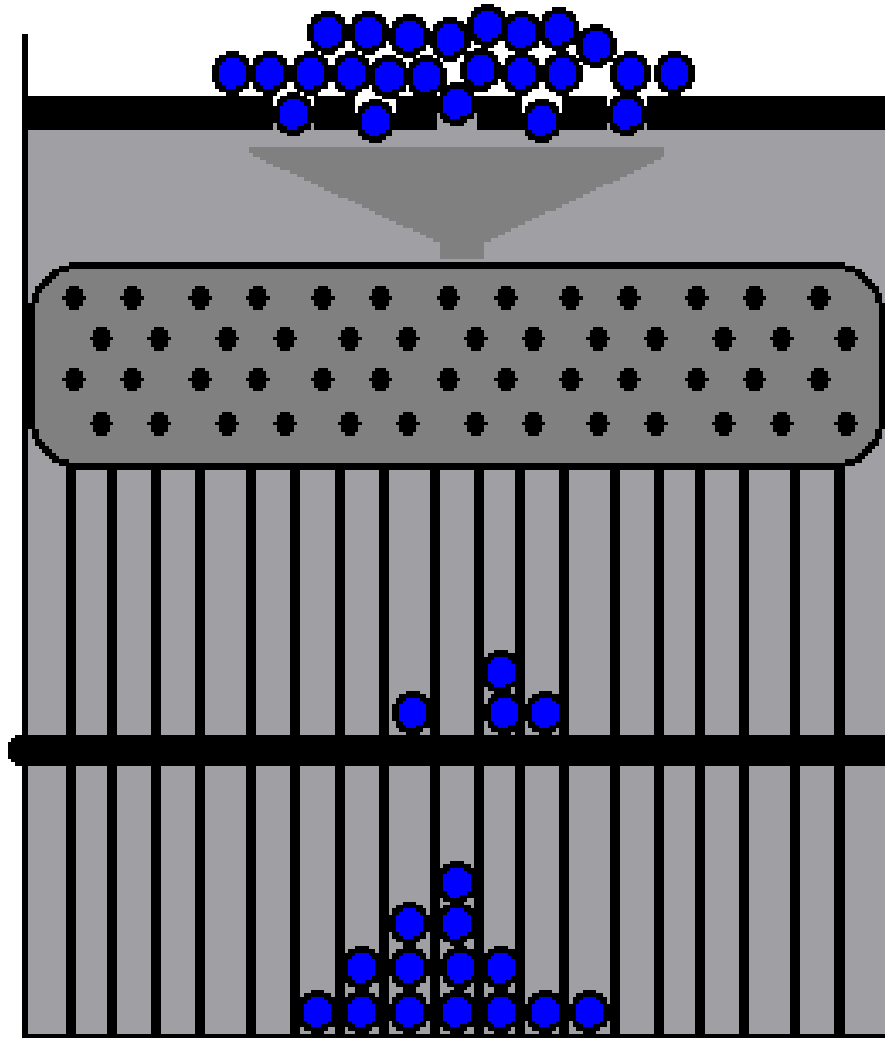


směrodatná (standardní) odchylka

průměr

SIR FRANCIS GALTON (1822-1911)

The Quincunx



Hopper

1. Beads are put into the hopper

Funnel

2. The beads drop through the funnel one by one

Pins

3. The beads bounce randomly through a series of pins

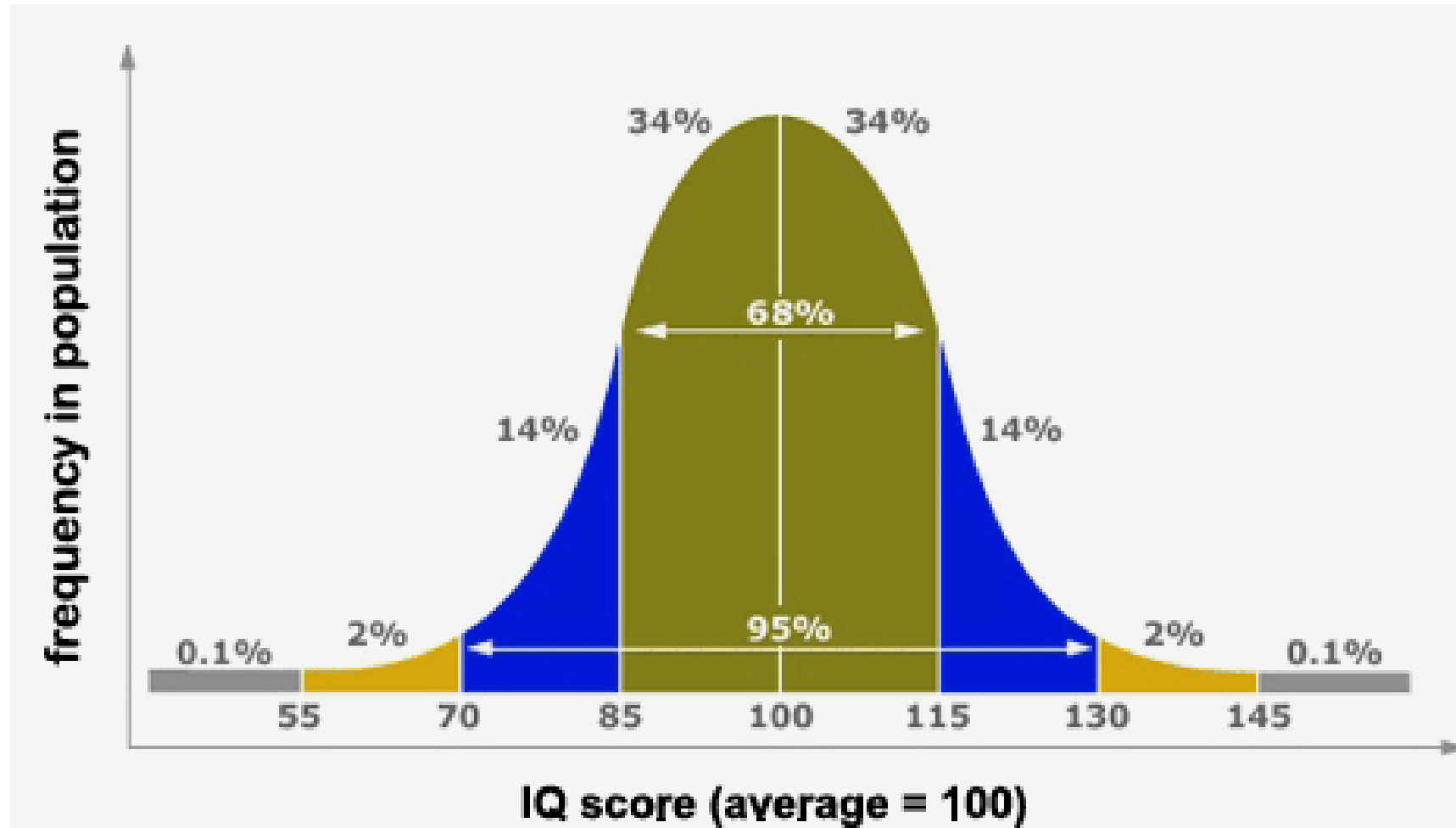
2nd Run

4. Each bead drops into a numbered slot

1st Run

5. The results of one run can be saved and compared with the results of a later trial

IQ TEST



TESTY NORMÁLNÍHO ROZDĚLENÍ

Histogram s křivkou normálního rozdělení

Posouzení šikmosti a špičatosti $<0; 1>$

- 0 = normální rozdělení
- 1 = rozdělení není normální
- _____ | | (neplatí pro $N > 200$)

Kolmogorovův-Smironovův test (K-S test): Procedura Explore

- H_0 : Rozložení proměnné se neliší od normálního.
- Kritická hranice sig. $> 0,05$

Důsledky:

- použití neparametrického testu
- transformace proměnných
- Použití parametrického testu

STANDARDIZOVANÉ (NORMOVANÉ) NORMÁLNÍ ROZDĚLENÍ

Základem pro inferenční statistiky

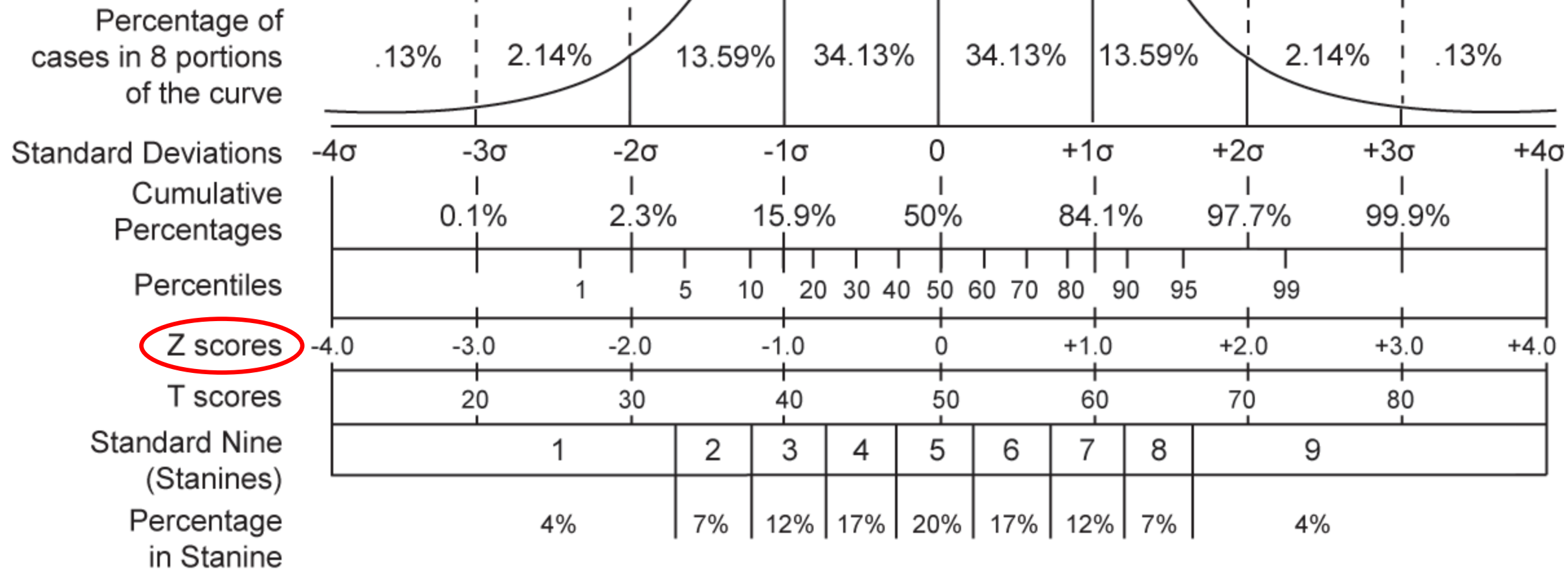
Převádí hodnoty na tzv. z-skóry, tzn. vyjadřuje hodnotu jako vzdálenost od průměru ve směrodatných odchylkách

— (Procedura DESCRIPTIVES)

Užití:

- Umožňuje porovnávat proměnné měřené na různých škálách a v různých jednotkách
- Standardizace dat

*Normal,
Bell-shaped Curve*



INFERENČNÍ (VÝBĚROVÁ) STATISTIKA

Inference = statistické usuzování (z výběrového souboru na základní soubor, z parametru na statistiku)

Odpovídají data v našem vzorku populaci? S jakou pravděpodobností?

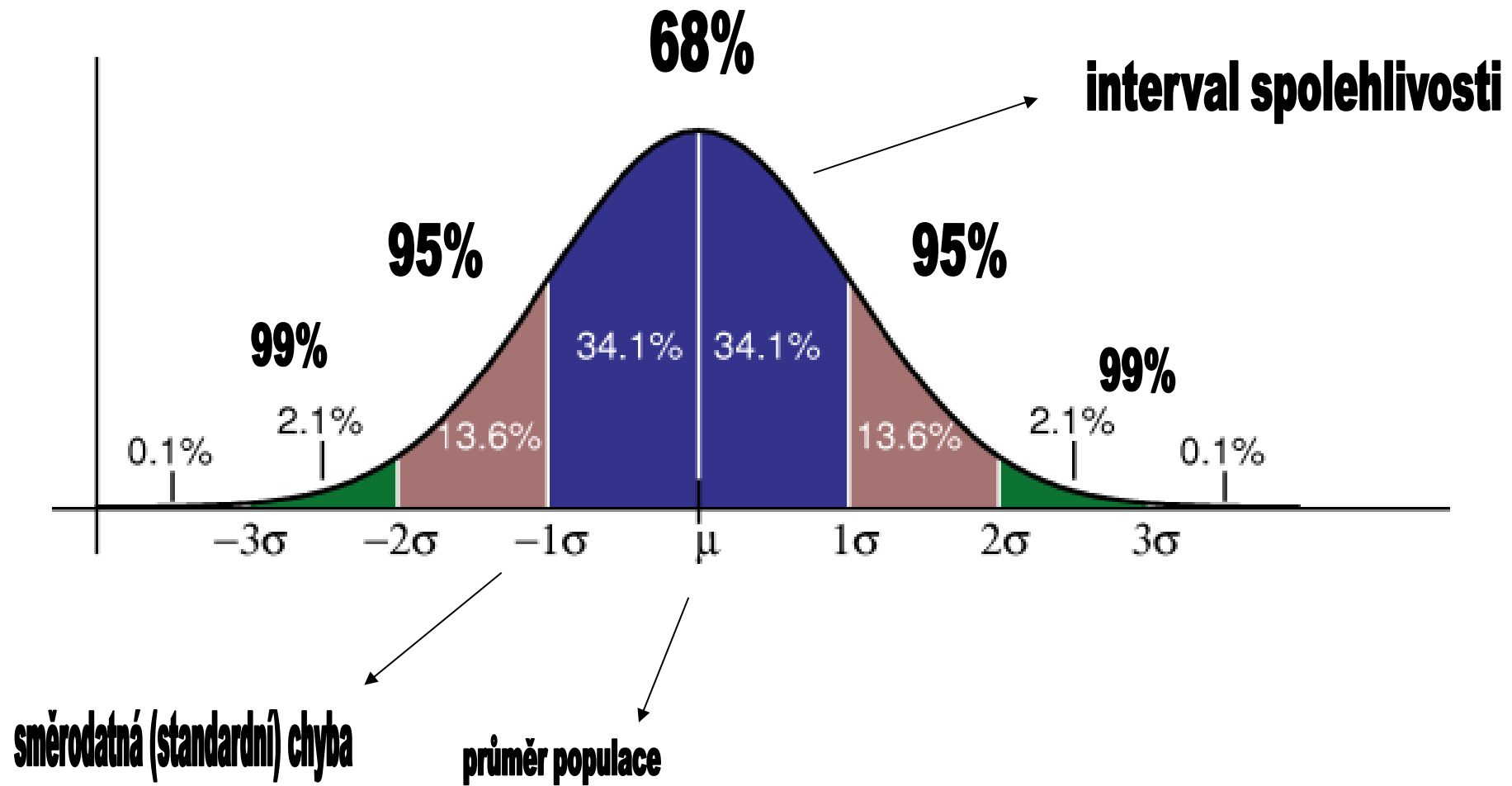
Zahrnuje

- (1) odhad parametrů – metody odhadu pro bodové a intervalové odhady
- (2) testování hypotéz – testy statistické významnosti

Předpoklady:

- Pravděpodobnostní výběr
- Vysoká návratnost (response rate)

TEORETICKÉ NORMALNÍ ROZLOŽENÍ (POPULACE)



CENTRÁLNÍ LIMITNÍ VĚTA

Centrální limitní věta (central limit theorem) říká, že když provedeme mnoho výběrů o určité velikosti založených na pravděpodobnostním principu, pak se rozložení (distribuce) výběrových průměrů přiblíží normálnímu rozdělení a celkový průměr těchto průměrů se bude podobat průměru v populaci. A to nezávisle na tom, jak jsou hodnoty proměnné rozloženy v populaci.

Standardní chyba průměru je pak směrodatnou odchylkou rozdělení tohoto rozdělení průměrů, vyjadřuje výběrovou chybu

STANDARDNÍ CHYBA PRŮMĚRU (STANDARD ERROR OF THE MEAN)

$$SE_{\bar{x}} = \frac{s}{\sqrt{n}}$$

← standardní odchylka

← velikost vzorku

Pro nominální data se používá **standardní chyba proporce** (standard error of proportion),

$$SE(p) = \sqrt{\frac{p \cdot q}{n}}$$

Výběrová chyba (sample error) na hladině $p=0.05$

1.96 * SE

Výpočet horního limitu konfidenčního intervalu na 95% hladině významnosti: **CI = m + (SE * 1.96)**

Výpočet spodního limitu konfidenčního intervalu na 95% hladině významnosti: **CI = m - (SE**

INTERVAL SPOLEHLIVOSTI (CONFIDENCE INTERVAL)

Parametry odhadujeme s určitou mírou pravděpodobnosti/rizika, tzv. hladina významnosti

- typicky 95% HV \rightarrow 5% riziko, příp. 99%HV
- Odvozeno od normálního rozdělení ($\pm 4 \rightarrow$ 99% případů, $\pm 3 \rightarrow$ 95% případů)
- Vyšší míra jistoty vede k širšímu CI
- - výběrový průměr
- - z-skór požadované úrovně pravděpodobnosti, pro je to 1.96
- - standardní směrodatná chyba, SD výběrových průměrů

Výpočet pro průměr:

(PROCEDURA DESCRIPTIVES)

$\bar{x} \pm z \cdot \frac{s}{\sqrt{n}}$, kde $\frac{s}{\sqrt{n}}$ je standardní chyba průměru

Výpočet pro relativní četnost (%):

relativní četnost

$\sqrt{\frac{p(1-p)}{n}}$, kde p je pozorovaná

TESTOVÁNÍ HYPOTÉZ

Testujeme, zda:

- Vzorek pochází z populace s určitým rozdělením (reprezentativita)
- Zda dva výběry pocházejí z téže populace (např. rozdíly mezi muži a ženami)
- Zda ne/existuje vztah mezi proměnnými

Nulová hypotéza H_0 – předpokládá neexistenci rozdílu, buď ji zamítnout lze nebo nelze

- Teoretická nulová hypotéza (Např. Neexistuje rozdíl mezi platy žen a mužů.)
- Statistická nulová hypotéza (Např. Rozdíl mezi průměrným platem mužů a průměrným platem žen je roven nule.)

Alternativní hypotéza H_1 – předpokládá rozdíl

- Oboustranná (two-tailed) - Mezi proměnnými pohlaví a příjem bude vztah, příjem se bude lišit.
- Jednostranná (one-tailed) – Ženy budou v průměru vydělávat méně.

URČETE H_0 NEBO H_1

H_0 : Respekt a Reflex se nebudou odlišovat v hodnocení (pozitivní/negativní) prezidenta Zemana.

H_0 : Mezi mírou prokrastinace a vlastnictvím účtu na Facebooku nebude žádný vztah.

H_1 : Děti, jimž rodiče kontrolují telefon budou méně často obětmi kyberšikany.

H_1 : Mezi mírou narcismu a počtem zveřejněných fotografií na sociálních sítích bude vztah. .

POSTUP TESTOVÁNÍ

Volba testu

Volba testovacího kritéria se známým rozložením

Výpočet hodnoty testového kritéria (data)

Interpretace výsledku

$HV < 0,05 \rightarrow$ zamítáme nulovou hypotézu, neboť pravděpodobnost, že bychom získali taká data, kdyby platila H_0 , je malá

- Sig. v SPSS
- p - v odborné literatuře
- Arbitrární charakter HV
- Statistická významnost \neq věcná významnost

Např. Kolmogorovův-Smirnovův test normality rozložení, srovnání CI pro dva populační průměry