

Deskriptivní statistika (kategorizované proměnné)

Nejprve malé opakování:

- **Deskriptivní statistika** se zabývá popisem dat, jejich sumarizaci a prezentací.
- **Kategorizované proměnné** jsou všechny proměnné, jejichž hodnoty se nacházejí v určitých kategoriích. Jedná se tedy o nominální, ordinální i kardinální proměnné (pouze ale kardinální poměrové).

Různé druhy proměnných umožňují různé druhy popisu.

Popis nominálních proměnných

U nominálních proměnných zjišťujeme:

- **rozložení četností** variant znaku (pomocí tabulek četností),
- nejčastěji zastoupenou kategorii – **modus** (modálních kategorií někdy může být více než 1),
- **variační poměr**, který se vypočítá tak, že od jedné odečteme podíl četnosti modální kategorie a velikosti souboru.

Popis ordinálních proměnných

U ordinálních proměnných zjišťujeme:

- rozložení četností variant znaku (pomocí tabulek četností),
- nejčastěji zastoupenou kategorii – **modus** (modálních kategorií někdy může být více než 1),
- **medián** (mediánovou kategorii),
- variační poměr,
- další vlastnosti, kterými se ale nebudeme dopodrobna zabývat.

Popis a kontrola dat

Prvním úkolem výzkumníka je popis výběrového souboru. Charakteristikou vzorku by měla začít každá analýza i analytická kapitola v bakalářské či diplomové práci. Zajímá nás například:

- Kolik je ve výběrovém souboru jednotek?

- Kolik je v souboru mužů a žen?
- Kolik je v souboru lidí se ZŠ/SŠ/VŠ vzděláním?
- Jak je v souboru distribuován věk?

Toto rozložení může být vyjádřeno v **absolutních, relativních, či kumulativních relativních četnostech.**

- **Absolutní četnost** udává absolutní číslo – hodnotu četnosti varianty proměnné v souboru.
Například: V souboru je 1456 mužů a 1201 žen.
- **Relativní četnost** udává **podíl** četnosti varianty proměnné v souboru.
Například: V souboru je 24 % osob se základním vzděláním.
- **Kumulativní relativní četnost** udává kumulativní podíly variant proměnné v souboru (nejsou použitelné pro nominální proměnné).
Například: V souboru je 36 % respondentů, kteří mají alespoň maturitu (tedy nejen úspěšní středoškoláci s maturitou, ale také vysokoškoláci se všemi variantami diplomů).

Popis a kontrola kategorizovaných dat

Tabulky četností

Pro zobrazení základních hodnot popisu rozložení hodnot kategorizovaných proměnných (tedy proměnných nominálních a ordinálních s menším počtem variant odpovědí) se používá tzv. **tabulka četností**. Ta obsahuje jak absolutní, tak relativní četnosti hodnot proměnných. Takto vypadá správná a kompletní tabulka četností:

Jaké je Vaše vzdělání?		Četnost odpovědí	Relativní četnost	Validní relativní četnost
Validní hodnoty	Základní	46	7,5 %	7,6 %
	Základní vyučen /střední bez maturity	62	10,1 %	10,2 %
	Střední s maturitou	307	50,1 %	50,5 %
	Pomaturitní nástavba, VOŠ	40	6,5 %	6,6 %
	Vysokoškolské	153	25,0 %	25,2 %
	Celkem validní hodnoty	608	99,2 %	100,0 %
Chybějící hodnoty (neví, neodpověděl/a)	Chybějící hodnoty	5	0,8 %	
Celkem		613	100,0 %	

V praxi se často používá jen zkrácená verze tabulky obsahující pouze validní četnosti:

Jaké je Vaše vzdělání?	Četnost odpovědí	Validní relativní četnost
Základní	46	7,6 %
Základní vyučen /střední bez maturity	62	10,2 %
Střední s maturitou	307	50,5 %
Pomaturitní nástavba, VOŠ	40	6,6 %
Vysokoškolské	153	25,2 %

Před počítáním četností je ale potřeba zkontrolovat data. Kontrolujeme, zda se nachází v platném intervalu (například proměnná pohlaví nabývá v našem souboru pouze hodnot 1 a 2, všechny ostatní varianty by měly být omyly).

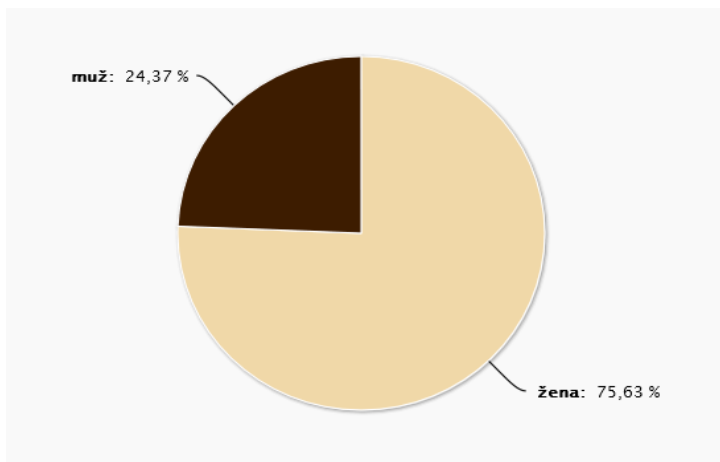
Grafy četností

Pro znázornění rozložení četností se využívají i grafy znázorňující četnosti hodnot proměnných. Nejznámějšími variantami jsou koláčový a sloupcový graf.

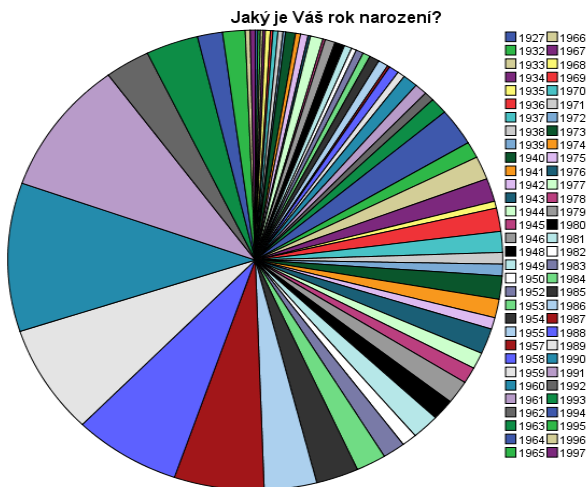
Koláčový graf je vhodný:

- pro třídění prvního stupně (jedna datová řada),
- pro porovnání četností u nominálních proměnných, které nemají příliš mnoho hodnot (méně než 7),
- pokud hodnoty, které chcete vykreslit, nejsou nulové,
- pokud hodnoty představují část celku.

Příklad proměnné, kde je vhodné využít koláčový graf:



Příklad proměnné, kde NENÍ vhodné využít koláčový graf:

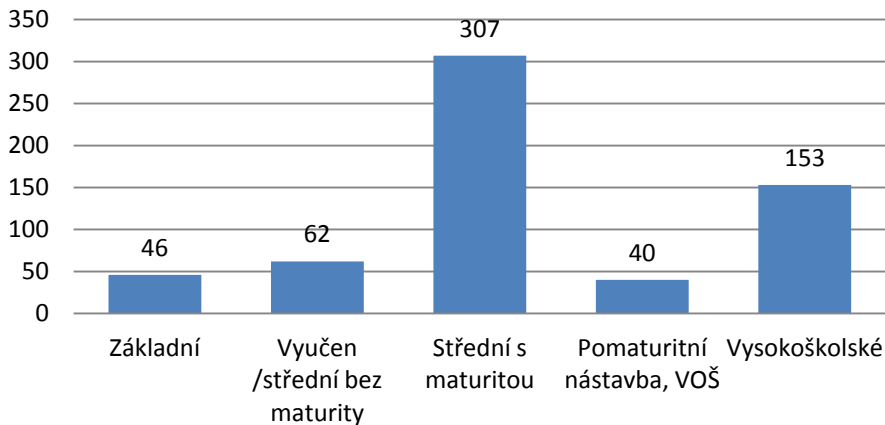


Sloupcový graf je vhodný pro:

- porovnání položek,
- ordinální proměnné a kardinální proměnné s menším počtem kategorií,
- znázornění změn za časové období (třídění druhého stupně).

Příklad sloupcového grafu:

Jaké je Vaše vzdělání?



Grafy se v Excelu vkládají pomocí funkce „Grafy“ na listu „Vložení“.

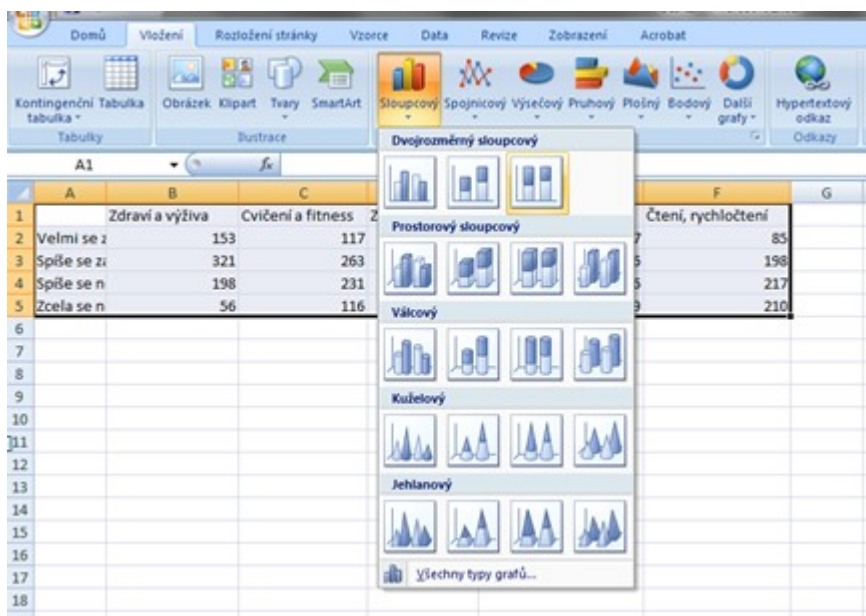
Porovnání rozložení četností

Pro zobrazení porovnání rozložení četností u baterií otázek se používají **skládané sloupcové grafy**.

Skádaný sloupcový graf můžete vytvořit tak, že si připravíte tabulku s absolutními validními četnostmi u jednotlivých kategorií:

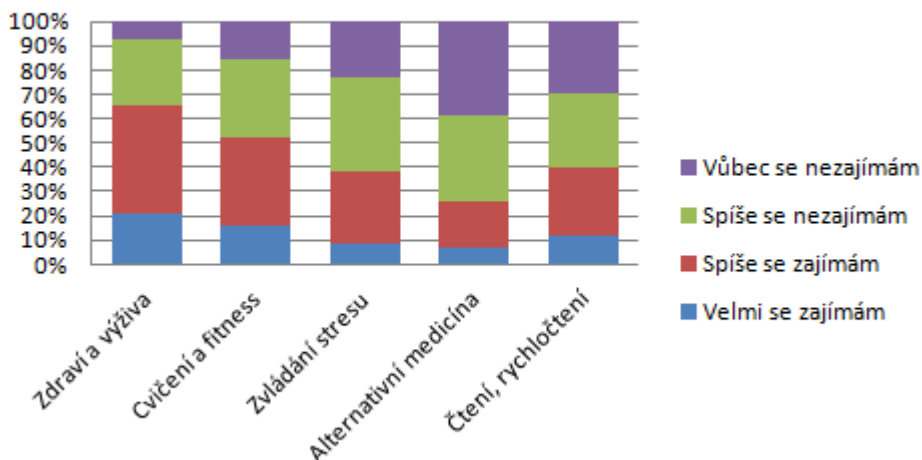
	A	B	C	D	E	F	G
1		Zdraví a výživa	Cvičení a fitness	Zvládání stresu	Alternativní medicína	Čtení, rychločtení	
2	Velmi se z	153	117	64	47	85	
3	Spíše se z	321	263	210	136	198	
4	Spíše se n	198	231	280	256	217	
5	Zcela se n	56	116	169	279	210	
6							
7							
8							

Tabulku si označíte a zvolíte možnost „Vložení“ – „Grafy“ – „Sloupcový“.



Výsledkem je skládaný sloupcový graf, který přehledně ukazuje rozdíly v rozložení jednotlivých proměnných.

Zájem o jednotlivé oblasti



Modus a medián

Pro připomenutí z minulého semestru si uvedme, v čem se liší MODUS a MEDIÁN (obě udávají tzv. míry centrální tendence a často se pletou):

MODUS je hodnota, která se v datech vyskytuje nejčastěji.

MODÁLNÍ KATEGORIE je tedy nejpočetněji zastoupená kategorie.

MEDIÁN dělí řadu výsledků seřazených podle velikosti na dvě stejně početné poloviny.

MEDIÁNOVÁ KATEGORIE je ta, ve které je dosaženo 50% všech údajů, postupujeme-li od první kategorie výše.

Jestliže je počet položek ve výzkumném souboru lichý, pak platí:

$$\text{Medián} = x_{(n+1)/2}$$

Jestliže je počet položek ve výzkumném souboru sudý, pak platí:

$$\text{Medián} = 0,5(x_{n/2} + x_{n/2+1})$$

Představte si otázku na počet dětí. Odpovědi respondentů jsou $\{0, 1, 1, 2, 2, 3, 5\}$.

- V souboru jsou dvě modální kategorie (tedy kategorie s nejvyšším počtem výskytů) – jsou to hodnoty 1 a 2.
- Mediánová kategorie je 2. Medián je na rozdíl od aritmetického průměru málo citlivý k odlehlým (extrémním) hodnotám. Pokud by byly odpovědi respondentů $\{0, 1, 1, 2, 2, 3, 5, 10\}$, medián stále zůstává roven 2.

Modus a medián v Excelu

V Excelu existují na výpočet mediánu a modu jednoduché příkazy MEDIAN a MODE. Syntaxe zápisu je snadná:

- =MEDIAN(datová oblast) – např. =MEDIAN(A1:A730)
- =MODE(datová oblast) – např. =MODE(A1:A730)

(Příkazy vypočítají medián a modus ze sloupce A, řádků 1-730.)

Tipy pro vytváření grafů

Levine a Stephan (2010) shrnují několik tipů pro prezentaci dat prostřednictvím grafů v akademickém prostředí:

- vždy si vyberte ten nejjednodušší graf,
- vždy používejte popisek grafu,
- popište obě osy,
- vyvarujte se ilustrací a zbytečného používání grafiky na pozadí nebo okrajích grafu,
- vyvarujte se používání módních piktoqramů, které by mohly ztížit čitelnost dat,
- vertikální osa by měla začínat nulou (pokud nezačíná negativními hodnotami).

V neakademickém prostředí (např. pro účely marketingu) je využití grafiky vhodné, v prostředí akademickém je na prvním místě čitelnost dat. 3D efekty a vkládání obrázků mohou znemožnit čtení hodnot dat. Další tipy pro vytváření grafů najdete třeba [zde](#).

Spojité proměnné

Spojité (nekategorizované) proměnné jsou ty proměnné, které mohou nabývat všech hodnot z daného intervalu. Může jednat o plat, věk, počet obyvatel města, délku pracovní zkušenosti v měsících...

Aritmetický průměr

Aritmetický průměr je třetí mírou centrální tendence. U kardinálních dat lze jako míry centrální tendence využívat všechny tři:

- modus,
- medián,
- aritmetický průměr.

Aritmetický průměr je ukazatelem „průměrné“ hodnoty, nemusí být ale vždy ukazatelem nejvhodnějším – vhodné je jej kombinovat s mediánem. Aritmetický průměr je totiž velmi citlivý na extrémní hodnoty. I jedna extrémní hodnota může výrazně posunout aritmetický průměr.

Příklad: V roce 2010 byl podle serveru Platy.cz průměrný měsíční plat 23 300 Kč. Medián byl však na hodnotě 21 000 Kč. Znamená to, že průměr vychýlil menší počet jedinců s výrazně vyšším platem.

Průměrný měsíční plat (v Kč)	Medián (Kč)	Rozdíl (v %)
23 300	21 000	11%

Zdroj: Platy.cz

Pro připomenutí:

Modus se používá, pokud:

- rozdělení má více vrcholů,
- chceme zjistit nejčastější hodnoty.

Medián používáme, pokud:

- jsou data ordinální nebo kardinální,
- chceme znát střed rozložení dat,
- (v kombinaci s průměrem) pokud soubor obsahuje extrémní hodnoty,
- jestliže je rozložení dat zešikmené.

Aritmetický průměr je vhodné používat, pokud

- jsou data kardinální,
- rozložení je symetrické,
- chceme použít statistické testy. (Hendl 2009)

Minimum, maximum a rozpětí

První charakteristiky nekategorizovaných dat, na které se díváme už při fázi čištění dat, jsou **minimální** a **maximální hodnoty**. Z nich také snadno spočítáme **rozpětí**.

Rozpětí je nejjednodušší míra variability a snadno se vypočítá jako rozdíl mezi nejvyšší a nejnižší hodnotou.

Např. Je-li minimální hodnota 18 a maximální 1024, rozpětí hodnot proměnné v souboru je 1006.

Rozptyl a směrodatná odchylna

Rozptyl je definován jako střední hodnota kvadrátů odchylek od střední hodnoty (průměru). Vyjadřuje variabilitu rozdělení souboru náhodných hodnot kolem její střední hodnoty. Při průměrování odchylek dělíme číslem $n-1$.

S rozptylem úzce souvisí **směrodatná odchylna**. Ta se vypočítá jako odmocnina z rozptylu. Vrací tedy míru rozptýlenosti do měřítka původních dat. V podstatě nám říká, uvnitř jakého intervalu okolo průměru leží zvolené procento případů – tedy čím je směrodatná odchylna menší, tím lépe pro aritmetický průměr.

Hendl (2009) srozumitelně vysvětluje, jak dochází k výpočtu směrodatné odchylny:

1. Nejprve si vypočítáme všechny odchylky od průměru (např. při hodu kostkou vždy spočítáme odchylku konkrétní hozené hodnoty od celkového průměru).
2. Umocnění na druhou převede záporné odchylky na kladná čísla. Zároveň zvýrazní váhu extrémnějších odchylek.
3. Sečteme kvadratických odchylek.
4. Dělením číslem $n-1$ získáme průměrnou kvadratickou odchylku.
5. Odmocnina (v případě směrodatné odchylny) převede výsledek do původního měřítka dat.

Pro názornost si pojdme ukázat příklad, který dobře znáte – hodnocení vyučujících na KISKu a směrodatnou odchylnu tohoto hodnocení.

Zajímavost předmětu	není vůbec zajímavý	***X(*)**...	je velmi zajímavý
Přínosnost předmětu	není vůbec přínosné	***X(*)*...	je velmi přínosné
Obtížnost obsahu	velmi snadný(*)**X**	velmi obtížný
Náročnost na přípravu	velmi snadný(*)X*...	velmi obtížný
Dostupnost studijních zdrojů	velmi špatně dostupné(*)**X*	velmi dobře dostupné
Jak učitel učí	velmi špatný	***X(*)**...	vynikající
Učitel jako odborník	není odborníkem(*)**X*	je odborníkem

Zajímavost předmětu	není vůbec zajímavý(*)...*X	je velmi zajímavý
Přínosnost předmětu	není vůbec přínosné(*)...*X*	je velmi přínosné
Obtížnost obsahu	velmi snadný	**X**(*)....	velmi obtížný
Náročnost na přípravu	velmi snadný	*X**(*)....	velmi obtížný
Dostupnost studijních zdrojů	velmi špatně dostupné(*)**X**	velmi dobře dostupné
Jak učitel učí	velmi špatný(*)...*X	vynikající
Učitel jako odborník	není odborníkem(*)...X	je odborníkem

Průměrné hodnocení proměnné „Učitel jako odborník“ je u obou vyučujících podobné – jeden vyučující má průměrné hodnocení 9, druhý má průměrné hodnocení 10.

Směrodatná odchylna (zvýrazněná hvězdičkami) nám ale poskytne rychlou další informaci – říká nám, jak moc se hodnocení všech respondentů pohybovalo kolem průměru. Vidíme, že zatímco v druhém případě se hodnocení výjimečně shodovalo a studující se shodli na tom, že učitel je skutečný odborník, v prvním případě nebyla shoda zdaleka tak velká.

Rozptyl a směrodatná odchylka v Excelu

- rozptyl – příkaz **VAR**
- směrodatná odchylka – příkaz **SMODCH.VÝBĚR**