

9. REGRESE

9.1 Základní analýza podle typu dat

Kardinální:

- korelační analýza

Kategorizovaná (ordinální, nominální):

- log-lineární analýza (nejlépe v programu LEM)

Faktorová a clusterová analýza hledají faktory v proměnných nebo clustery v hodnotách. Obě jsou ale primárně určeny pro kardinální data, pro ordinální a nominální proměnné se používá obdobná analýza latentních tříd (latent class analysis).

9.2 Volba správného typu regrese

Volí se podle typu závisle proměnné:

Kardinální:

- lineární regrese, ve Statě příkaz REGRESS, pokud jsou data intervalová (nemají neutrální hodnotu, tj. chybí nula)
- Poissonova regrese, ve Statě příkaz POISSON pokud jsou data poměrová (mají rozsah od nuly do nekonečna, např. počet pokusů o sebevraždu)

Ordinální a Nominální:

- binární logistická regrese (hodnoty 0,1), ve Statě příkaz LOGIT
- ordinální logistická regrese (uspořádané hodnoty, vzdálenosti mezi jednotlivými hodnotami jsou stejné), ve Statě příkaz OLOGIT
- multinomická logistická regrese (hodnoty nemusí být uspořádány a mezi nimi nemusí být stejná vzdálenost), ve Statě příkaz MLOGIT.

9.3 Lineární regrese

Odhadem koeficientů regresního modelu hledáme rovnici přímky, která ideálním způsobem proloží body jednotlivých pozorování. Základní rovnice přímky má podobu:

$$y = a + bx + e$$

kde y je závisle proměnná, a je konstanta vypočtená z regresního modelu (průsečík přímky s osou y pro $x=0$), b je koeficient regresního modelu vypočtený z regresního modelu (sklon přímky proti ose x), x je hodnota nezávisle proměnné, e jsou nevysvětlené rezidua.

Přidáváním dalších nezávisle proměnných x_1, x_2, x_3, \dots se rovnice komplikuje následovně:

$$y = a + b_1x_1 + b_2x_2 + b_3x_3 + \dots + e$$

REGRESS – odhadne koeficienty regresního modelu zadaného za příkazem REGRESS. Na první místo se píše závisle proměnná (vysvětlovaná), za ni se v libovolném pořadí zapisují nezávisle proměnné (vysvětlované, prediktory, determinanty)

```
regress god
regress god age
```

- R-squared: z kolika procent model vysvětluje sociální realitu. Po vynásobení stovkou dostaneme míru vysvětlování v procentech.
- Prob > F: statistická významnost celého modelu (platí obvyklá hranice 0.05)
- P>|t|: statistická významnost jednotlivých vysvětlujících proměnných (platí obvyklá hranice 0.05)

Kategorizované proměnné je potřeba zadat jako tzv. dummy proměnné, což v praxi znamená, že pro každou hodnotu proměnné je vytvořena nová proměnná nabývající hodnot 0 a 1. Např. pro proměnnou `dny_v_tydnu` by bylo vytvořeno sedm dummy proměnných `dny_v_tydnu_pondeli` (nabývající hodnoty 1 v pondělí a hodnot 0 v jiných dnech), `dny_v_tydnu_uty` (nabývající hodnoty 1 v úterý a hodnot 0 v jiných dnech) apod. Při odhadu regresních modelů stačí jednoduše před každou kategorizovanou (tj. ordinální nebo nominální) zapsat písmeno `i` s tečkou. Stata pak zvolí první hodnotu (tedy např. pondělí) jako tzv. referenční a vysvětluje efekt následujících hodnot ve srovnání s touto referenční hodnotou.

```
regress god age i.v291
```

Pokud nás zajímá souvislost dvou či více proměnných, necháme odhadnout tzv. interakci. Stata pak vypočte, jak podobu výsledné přímky ovlivňuje kombinace dvou nezávisle proměnných (např. kombinace věku a pohlaví). Interakce se zadává jednoduše pomocí znaku # mezi dvěma proměnnými.

```
regress god age i.v291 age#i.v291
```

Jelikož věk není vhodné pojímat jako kategorizovanou proměnnou (Stata vypočte interakci pro každou hodnotu věku a každé pohlaví, např. tedy 17 let pro muže, 17 let pro ženy, 18 let pro muže, 18 let pro ženy ...). Pomocí znaku c a tečka můžeme Statě nařídit, aby danou proměnnou považovala za spojitou.

```
regress god age i.v291 c.age#i.v291
```

PREDICT – vypočte hodnoty proměnné podle posledního vypočteného modelu. Důležitými parametry jsou XB, který odhadne hodnoty lineárního modelu, a RES, který vypočte hodnoty reziduálů (rozdíl mezi naměřeno a vypočtenou hodnotou)

```
predict novapromenna, xb res
```