

Experimental Design and Statistics - AGA46E

M. Maciak (Czech University of Life Sciences, Prague)

Lab Session 3 - Summer Term 2015

1 Prehled nekterych prikazu

- `read.table()` a `read.csv()` pro nacitani vlastniho datoveho suboru do Rka; Prvni prikaz je vseobecny pro libovolny datovy soubor, druhy prikaz (`csv`) slouzi specialne pro nacitani datovych souboru s koncovkou `csv` (doporuceny typ datoveho souboru);
- `data` je take mozne nacist primo z "clipboardu" pocitace a to pomocou kombinace tlacitek CTRL + C a prikazu:

```
read.table(file = "clipboard", sep="\t", header=TRUE)
```

- dodatecne parametre, ktere podrobneji specifikuju datovu soubor, ktery chceme nacist, jsou napr. `sep` pro definovani oddelovace sloupce, `dec` pro definovani desetinnych cisel (symbol pro desetinou carku), `header` pro upresneni jestli data obsahuju hlavicku ci nikoli a pripadne parametr `skip`, ak je za potreby nacitavat data pozdeji nez hned v prvni radku;
- funkce pro vykreslovani obrazku: `plot()`, `boxplot()`, `barplot()`, `pie()`, `dotchart()`, `hist()` - tyhle funkce zaroven oteviraji v Rku graficke okno, do ktereho sa postupne zakresluje obrazek;
- dodatecne graficke funkce v Rku: `lines()`, `abline()`, `segments()`, `points()`, `text()`, nebo `legend()` - tyhle funkce lze pouzit pouze v pripade, ze graficke okno je uz otevrene; Funkce slouzi k dokresleni dodatecnych objektu jiz do existujiciho obrazku.

2 Nahodne Generatory v Rku

- pouzite google a vyhledejte prikazy, ktere se v Rku pouzivaji ke generovani nahodnych cisel z Binomickeho, Poissonoveho, ci Geometrickeho rozdeleni;
- dodatecne parametre jsou nutne k spravne fungovani. Jake parametre to jsou a jak funguji? (vyzkousejte ruzne nastaveni pro dodatecne parametre a pomoci graficky nastroju - napr. funkce `hist()` zjistite, jaky vlyv maji na generovane hodnoty);
- vygenerujte dostatecne dlouhu posloupnost pro nejake rozdeleni a pomoci obrazku overte, jestli tvar distribucni funkce zodpoveda s distribucni funkci daneho rozdeleni. (pouzite napr. prikazy `hist()`, `ecdf()`, nebo `stepfun()`);
- zopakujte to same, ale zakazdym pro delsi posloupnost generovanu z daneho rozdeleni; (napr. pro rozsah nahodniho vyberu $n = 10, 20, 50, 100, 1000, 10000, \dots$)

3 Diskretní nahodní veličiny a faktory v Rku

- nektěre diskretní nahodní veličiny muzu byt take reprezentovány jako faktory - faktorové proměny; Každý faktor má určité počty úrovní (hladin) a každá úroveň má své pojmenování - label; (např. namísto hodnot 0 a 1 pro muže a ženy, můžeme stejně použít kódování 'male' a 'female', nebo 'M' a 'F')
- pojmenování faktorů a jejich příslušných úrovní by mělo byt intuitivní a přímé;
- jaký je hlavní rozdíl mezi diskretní interpretací hodnot a faktorovou interpretací hodnot? (např. mezi hodnotami 0 a 1 a označením 'M' a 'F', v obou případech pro rozlišení mužů a žen)

```
> genderVector <- rbinom(20, 1, 0.5) # random vector of zeros and ones
> genderFactor <- as.factor(genderVector)
```

- porovnejte diskretní interpretaci `genderVector` a faktorovou interpretaci `genderFactor` a také `class(genderVector)` s `class(genderFactor)`;
- můžeme změnit pojmenování faktorových úrovní (labels) z numerických znaků 0 a 1 na cokoli jině, např. "male" a "female", nebo "M" a "F", atd. (povšimnete si, že pak už se nejedná o numerické hodnoty);

```
> levels(genderFactor) <- c("male", "female")
> genderFactor
```

- využijte náhodné generátory v Rku a nageenerujte nějaké sekvence náhodných diskretních hodnot; Následně z hodnot udelejte faktory s určitým počtem úrovní a nakonec tyto úrovně přejmenujte (použijte nějaké intuitivní názvy příslušných úrovní);

Například:

```
> accidents <- rpois(31, 2)
> accidents[accidents >= 3] <- 3
> accidents[accidents >= 1 & accidents <= 2 ] <- 1
> accidents[accidents == 0] <- 0
> accidentsFactor <- factor(accidents, labels = c("None", "Minor", "> 3"))
```

- použijte datový soubor `Orange` (k dispozici přímo v instalaci Rka) a ověřte, které proměnné jsou faktory a jaké jsou jejich příslušné úrovně;

```
> Orange
> levels(Orange[,1])
```

- pomocí příkazu `levels()` změňte jména příslušných úrovní;
- vytvořte vlastní datový soubor (tabulku) např. pomocí příkazu `matrix`, nebo `data.frame()`, který bude obsahovat alespoň tři faktorové proměnné, každou s alespoň dvěma úrovněmi;
- proč je důležité (až kritické) správně rozpoznávat diskretní proměnné a faktorové proměnné?

4 Realne data s diskretnymi a faktorovymi promennymi

- stahnite a nactete datovy soubor `passengerData.csv`;
- vyuzijte zakladni popisne charakteristiky, abyste ste dozvedeli neco o datech; (jaka je jejich struktura, jake hodnoty nabyvaji a pod.)
- ktere promenne jsou faktory a ktere ne?
- vyuzijte graficke nastroje v Rku a udelejte nekolik obrazku;
- pomoci dodatecnym funkci udelejte obrazky hezci; (napr. dodatecne parametre `pch`, `xlab` a `ylab`, `main`, `col`, prikaz `legend()`, a dalsi)