

Design experimentu a statistika - AGA46E

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

Final Project - zadani - Duben 2015

Zakladni instrukce

- Finalni projekt je individualni prace ve ktere kazdy student samostatne analyzuje vlastni datovy soubor. Kazdy student ma k dispozici individualni data vygenerovane na zaklade datumu narozeni. Finalni projekt musi byt spracovan za pomoci vhodneho statistickeho softwaru (preferovany je statisticky software R, ale mozne jsou i jine varianty). Program **Microsoft Excel neni vhodny statisticky software!**

- Data potrebne k spracovani finalniho projektu jsou simulovana data s fiktivnimi vysledkami hodnoceni studentov Ceske Zemedelske Univerzity v Praze. Datovy soubor obsahuje nekolik tisíc fiktivnich studentu a nekolik promennych (diskretnych a spojitych), ktere byli pro kazdeho studenta zaznamenany.

Datovy soubor finalProjectData2015.RData (k stazeni primo po kliknuti na predchozi odkaz se jmenem souboru) je k dispozici take v univerzitnem systemu Moodle a na webové stránce predmetu.

- **Dulezite:** Ke spracovani finalniho projektu nepouzivejte cely datovy soubor, pouze mensi (individualni) datovy soubor vygenerovany pomoci automatickeho R scriptu *generator.r*, ktery je stejne soucasti tohoto finalniho projektu (k stazeni v systemu Moodle a take na webové stránce predmetu: *generator.r*). Podrobny postup, jak si kazdy student/studentka samostatne vygeneruje prislusny datovy soubor, je detailne popsany dalsi sekci.
- **Finalni projekt je nutne odovzdat nejpozdeji dva dni pred skouskou! Student muze reseni odovzdat jako "paper copy" nebo ho zaslat emailem na adresu *maciak@af.czu.cz*. V pripade, ze student posle reseni emailem, je nutne, aby byl finalni report ulozeny ve formatu PDF (pouze jeden soubor PDF, ktery obsahuje vsechnu potrebne nalezitosti - popis, obrazky, tabulky, vysledky a finalni interpretaci). Nazev PDF souboru musi byt intuitivni (napr. *'jmenoprijmeni.pdf'*).**
- Kazdy projekt bude vyhodnocen individualne. Maximum ohodnoceni je 100 %, co predstavuje 10 bodu do zapoctoveho hodnoceni. Podrobne informace o hodnoceni predmetu a zpusobu ziskani zapoctu jsou v systemu Moodle, nebo na stránce predmetu);

Popis promennych

- data obsahuji simulovane vysledky studentu Ceske Zemedelske Univerzity v Praze.
- v celkovem datovem souboru je zaznamenano 6500 studentu a sledovanych je celkove 16 ruznych promennych (spojitych, diskretnych i faktorovych).
- kazdy student dostane ke spracovani pouze mensi cast datovehou souboru, ktery bude automaticky vygenerovany pomoci souboru *generator.r*. Datovy soubor pro kazdeho studenta by mel obsahovat cca 700 pozorovani (radku) a cca 9-10 promennych (sloupcu). Presne hodnoty poctu pozorovani a poctu promennych jsou ruzna pro ruznych studentov.

- popis jednotlivých proměnných je uvedený níže. Ne všechny proměnné jsou ale nutně obsazeny v souboru, který bude mít k dispozici každý student - student pouze dostane výběr některých proměnných.
 1. *Gender* - pohlaví studenta: 'male' nebo 'female';
 2. *Age* - věk studenta udávány v letech;
 3. *Nationality* - národnost studenta - kódována pouze jako 'czech' pro studenta/studentku české národnosti a 'other' jinak;
 4. *Expenses* - celkové měsíční náklady na živobytí odhadnuté studentem v tisících korunách;
 5. *Housing* - informace o tom, jestli student bydlí na koleji - 'dormitory', nebo s rodiči - 'parents', nebo ve vlastním - 'own';
 6. *Class* - informace o tom, ve kterém ročníku je daný student aktuálně zapsán;
 7. *Faculty* - informace o tom, na které ČZU fakultě student aktuálně studuje;
 8. *WeekHours* - přibližný průměrný počet hodin, které je daný student týdně na univerzitě/fakultě;
 9. *WorkHours* - přibližný průměrný počet hodin, které student potřebuje týdně na samostudium a přípravu do školy;
 10. *WeekClasses* - přibližný průměrný počet vyučovacími hodinami za jeden týden;
 11. *Score1*, *Score2* a *FinalScore* - celkový průměrný prospěch v procentech za dva různé semestry (1 a 2) a pak celkové hodnocení za všechny semestry;
 12. *Grade1*, *Grade2* a *FinalGrade* - celkový průměrný prospěch propočítaný na výslednou známku - uváděné za dva různé semestry (1 a 2) a pak celkové hodnocení za všechny semestry dohromady;

Generování vlastního datového souboru

- k vygenerování vlastního datového souboru je za potřeby samostatný R-kový skript *generator.r*, který je k stažení přímo po kliknutí na předchozí odkaz, nebo je k dispozici v systému Moodle nebo na webové stránce předmětu.
- soubor *generator.r* si stáhnete a uložíte do pracovního adresáře. Do stejného pracovního adresáře uložíte i datový soubor *finalProjectData.RData*.
- otevřete si program R a nastavíte správný pracovní adresář (pomocí příkazů `setwd()` a `getwd()` nebo vyklikáním v hlavním menu programu).
- otevřete, jestli jste ve správném adresáři - zadejte v R-ku příkaz


```
> list.files()
```

a měli byste vidět seznam všech souborů, které se nacházejí v daném pracovním adresáři. Jestliže je pracovní adresář správný, musí mezi nimi také být soubory *generator.r* a *finalProjectData.RData*.
- v R-ku zadejte následující příkaz a potvrďte tlačítkem 'enter':


```
> source("generator.r")
```
- následně postupujte podle instrukcí na obrazovce:
 1. zadejte prvních 6 čísel svého rodného čísla (bez lomítka, pouze prvních 6 čísel) a následně potvrďte tlačítkem 'enter';
 2. vyberte si jednu z ponukovaných možností: zadejte hodnotu 1 když preferujete hokej, hodnotu 2, když preferujete fotbal a hodnotu 3, když preferujete tenis a opět potvrďte tlačítkem 'enter';
 3. TO JE VŠE!

4. príslušný datový súbor s jmenom *mojeData* bol automaticky vytvorený v R-ku (napíšte napr. príkaz `ls()` a uvidíte objekt s názvom *mojeData*).
5. príslušný datový súbor bol rovnako uložený do Vášho pracovného adresára - meno súboru je *mojeDataVASECISLO.txt*. Tento súbor môžete kedykoľvek opakovane načítať v programe R pomocou príkazu `read.table()`.
6. **ke spracovaní konečného projektu použijete úroveň významnosti, ktorú Vám vygeneroval softvér pri generovaní dátového súboru (príslušná úroveň spoľahlivosti sa automaticky zobrazí na obrazovke)!**
7. generovanie dátového súboru môžete kedykoľvek zopakovať (dostanete rovnaké dáta), stačí, keď zadáte prvých 6 čísel rodového čísla a správnu voľbu športu, ktorú preferujete;

Pozadavky ke spracovaniu

Spracovanie konečného projektu je celkom individuálne. Je však potrebné, aby bol konečný projekt uceleným reportom - úvod, popis dát, hypotézy, metodológia, spracovanie, výsledky, a konečná interpretácia. Celkový objem reportu nie je stanovený, bez problému je možné odovzdať report na cca 2-3 stranách, a report bude kompletný.

Je však potrebné splniť nasledujúce:

1. uviesť niekoľko popisných štatistík, popísať dáta, urobiť niekoľko obrázkov, vysvetliť, predpokladať možné závery...
2. aspoň dva intervaly spoľahlivosti (podľa Vašej voľby);
3. aspoň dva štatistické testy (opäť podľa Vašej individuálnej voľby);
4. aspoň jeden párový a jeden nepárový *t*-test;
5. aspoň dvakrát využiť metódu analýzy rozptylu interpretovať výsledky;
6. aspoň jednou použiť lineárny regresný model;
7. **V úvode reportu, ktorý odovzdáte, uvádzte 6 čísel rodového čísla, ktoré ste použili k simulácii a rovnako tak Vašu voľbu ohľadom preferovaného športu.**