

# Statistická analýza dat

Jméno: \_\_\_\_\_

Podpis: \_\_\_\_\_

Cvičení	
Zkouška (písemná + ústní)	$\geq 25$
Celkem	$\geq 50$
Známka	

**Pokyny k vypracování:** doba řešení je 120min, jasně zodpovězte pokud možno všechny otázky ze zadání, pracujte s pojmy používanými v předmětu, můžete používat kalkulátory.

**Statistiké minimum.** (10 b) Zodpovězte následující otázky:

- (a) (5 b) Uvažujte náhodný vektor  $\mathbf{X}$ . Definujte kovarianční a korelační matici. Jaké mají tyto matice vlastnosti? K čemu se dají použít?
- (b) (5 b) Vysvětlete význam pojmu matoucí proměnná (confounding variable). Uveďte příklad a naznačte vliv na model.

**Analýza rozptylu.** (10 b) Odpovězte na otázky níže.

(a) (2 b) K čemu se používá parametrická jednostupňová analýza rozptylu (parametric one-way ANOVA)? Formulujte její nulovou a alternativní hypotézu.

(b) (3 b) Jaké má tato metoda předpoklady? Jak je budete testovat? Co se stane, pokud splněny nejsou?

(c) (3 b) Podrobně popište výstupní tabulku ANOVA testu na konci posloupnosti příkazů níže.

```
F<-unlist(mapply(rep,times=c(8,9,10),x=c(1,2,3)))
O<-F+rnorm(n=27,mean=0,sd=2)
summary(aov(O ~ as.factor(F)))
   Df Sum Sq Mean Sq F value Pr(>F)
F      2    11.41    5.707   1.953  0.164
Residuals 24   70.14    2.923
```

(d) (2 b) K čemu slouží následný post-hoc test? Na jakém principu je založen?

**Diskriminační analýza.** (10 b) Níže diskutujte vlastnosti lineární a kvadratické diskriminační analýzy (LDA a QDA).

(a) (2 b) Z jaké myšlenky obě metody vycházejí? Napište definiční vztah.

(b) (2 b) Jaký je základní rozdíl mezi LDA a QDA? Z čeho plyne?

(c) (1 b) Předpokládejte, že řešíte problém s lineární bayesovskou rozhodovací hranicí. Která z metod dosáhne vyšší přesnosti nad trénovacími daty? Která nad testovacími? Proč?

(d) (1 b) Předpokládejte, že řešíte problém s nelineární bayesovskou rozhodovací hranicí. Která z metod dosáhne vyšší přesnosti nad trénovacími daty? Která nad testovacími? Proč?

(e) (1 b) Uvažujte obecnou klasifikační úlohu. S rostoucím počtem trénovacích příkladů relativní testovací klasifikační přesnost QDA vzhledem k LDA poroste, bude klesat nebo se nebude měnit? Proč?

(f) (3 b) Máte určit, zda na akci firmu s loňským ročním výnosem 4% bude vyplacena dividenda. Z burzovní analýzy velkého počtu firem víte, že firem, které vyplácí dividendu, je 80% a jejich průměrný roční výnos je 10%. Firmy bez dividendy mají průměrný výnos 0%. Rozdělení výnosů v obou skupinách je normální s rozptylem  $\hat{\sigma}^2 = 0.36$ . Budete aplikovat LDA, nebo QDA? Nemusíte důsledně počítat pravděpodobnost, stačí přesně zapsat.

**Multivariátní regrese.** (10 b) Sestavujete multivariátní lineární model. Závisle proměnných je velký počet, hledáte model, který minimalizuje kritérium ( $y_i$  je hodnota závisle proměnné v i-tém vzorku,  $x_{ij}$  je hodnota j-té nezávisle proměnné v i-tém vzorku):

$$\sum_{i=1}^m (y_i - \hat{\beta}_0 - \sum_{j=1}^p \hat{\beta}_j x_{ij})^2 + \lambda \sum_{j=1}^p \hat{\beta}_j^2$$

Parametr  $\lambda$  nejprve nastavíte na 0, poté jej postupně zvyšujete. S nárůstem  $\lambda$

(a) (2 b) trénovací reziduální součet čtverců (residual sum of squares, RSS)

- i) zpočátku poroste, od jisté doby ale začne klesat a vytvoří invertovanou U křivku,
- ii) zpočátku bude klesat, od jisté doby ale začne růst a vytvoří U křivku,
- iii) bude stále růst,
- iv) bude stále klesat,
- v) zůstane konstantní.

(b) (2 b) testovací RSS

- i) zpočátku poroste, od jisté doby ale začne klesat a vytvoří invertovanou U křivku,
- ii) zpočátku bude klesat, od jisté doby ale začne růst a vytvoří U křivku,
- iii) bude stále růst,
- iv) bude stále klesat,
- v) zůstane konstantní.

(c) (2 b) variance

- i) zpočátku poroste, od jisté doby ale začne klesat a vytvoří invertovanou U křivku,
- ii) zpočátku bude klesat, od jisté doby ale začne růst a vytvoří U křivku,
- iii) bude stále růst,
- iv) bude stále klesat,
- v) zůstane konstantní.

(d) (2 b) zaujetí (bias)

- i) zpočátku poroste, od jisté doby ale začne klesat a vytvoří invertovanou U křivku,
- ii) zpočátku bude klesat, od jisté doby ale začne růst a vytvoří U křivku,
- iii) bude stále růst,
- iv) bude stále klesat,
- v) zůstane konstantní.

(e) (2 b) neredukovatelná chyba (irreducible error  $\epsilon$ )

- i) zpočátku poroste, od jisté doby ale začne klesat a vytvoří invertovanou U křivku,
- ii) zpočátku bude klesat, od jisté doby ale začne růst a vytvoří U křivku,
- iii) bude stále růst,
- iv) bude stále klesat,
- v) zůstane konstantní.

**Robustní statistika.** (10 b) Odhadněte dvěma různými metodami robustně rozptýlenost (scale) ze vzorku  $\{-1.84, 1.18, 0.0499, -0.751, -0.00707, -2.05, -1.47, -0.0520, -0.991, -0.945\}$ .

(a) (2 b) Metoda 1 (popis a aplikace na vzorek):

(b) (2 b) Metoda 2 (popis a aplikace na vzorek):

(c) (2 b) Dejte tyto odhady do vztahu s obvyklým odhadem standardní odchylky.

(d) (2 b) Popište kritéria, jež jsou určující pro kvalitu robustního odhadu rozptýlenosti.

(e) (2 b) Diskutujte výhody a nevýhody vámi zvolených metod podle kritérií popsaných v předchozím bodě.