

## שאלות לדוגמה למבחן:

1. נתונות ארבע תצפיות חד מימדיות שערכיהן 1,2,2.5,3. מריצים עליהן אלגוריתם  $k$ -means עם  $k=2$ .

a. מהן החלוקות הסופיות האפשריות?

b. מהי החלוקה האופטימאלית?

2. נתונות תצפיות חד מימדיות  $X_1, \dots, X_n$  שהן עירוב של שני משתנים נורמליים עם תוחלת 0 ושונות לא ידועה.

a. מהי פונקציית  $\log likelihood$  של הנתונים. (לכתוב תוך שימוש בפרמטרים  $(\pi, \sigma_1, \sigma_2)$ )

b. תארו אלגוריתם מסוג EM שאומד את הפרמטרים  $\pi, \sigma_1, \sigma_2$ . הסבירו את תשובתכם.

3. נתון מודל רגרסיה רב ממדית  $Y_i = \beta_0 + \sum_{j=1}^p X_{i,j} \beta_j + \varepsilon_i \quad i = 1, \dots, n$

יהיו  $\hat{\beta}_0^L, \hat{\beta}_1^L, \dots, \hat{\beta}_p^L$  אומדי LASSO עבור פרמטר  $\lambda$  מסויים.

הראו שמתקיים:

$$\hat{\beta}_0^L = \bar{Y} - \sum_{j=1}^p \hat{\beta}_j^L \bar{X}_j$$

$$\text{כאשר } \bar{X}_j = \frac{1}{n} \sum_{i=1}^n X_{i,j} \text{ ו- } \bar{Y} = \frac{1}{n} \sum_{i=1}^n Y_i$$

4. נתון מודל רגרסיה פשוטה:  $Y_i = \beta_0 + \beta_1 X_i + \varepsilon_i \quad i = 1, \dots, n$ . שוקלים שני מודלים: מודל A עם חותך

בלבד (כלומר  $Y_i = \beta_0 + \varepsilon_i$ ) ומודל B המודל המלא (עם  $X$  ועם חותך). הראו שלפי הקריטריון של  $C_p$

$$\text{נבחר במודל A אם ורק אם } \hat{\beta}_1^2 < \frac{2\hat{\sigma}^2}{\sum_{i=1}^n (X_i - \bar{X})^2}$$

הדרכה:

a. הראו שנבחר במודל A אם ורק אם  $TSS + 2\hat{\sigma}^2 < RSS + 4\hat{\sigma}^2$  כאשר

$$TSS = \sum_{i=1}^n (Y_i - \bar{Y})^2, \quad RSS = \sum_{i=1}^n (Y_i - \hat{\beta}_0 - \hat{\beta}_1 X_i)^2$$

b. העזרו בתוצאה מהכיתה כדי להראות ש  $TSS - RSS = TSS r^2$  כאשר

$$r^2 = \frac{[\sum_{i=1}^n (X_i - \bar{X}) Y_i]^2}{\sum_{i=1}^n (X_i - \bar{X})^2 \sum_{i=1}^n (Y_i - \bar{Y})^2}$$

c. הראו ש  $TSS r^2 = \sum_{i=1}^n (X_i - \bar{X})^2 \hat{\beta}_1^2$