

Εξόρυξη Δεδομένων

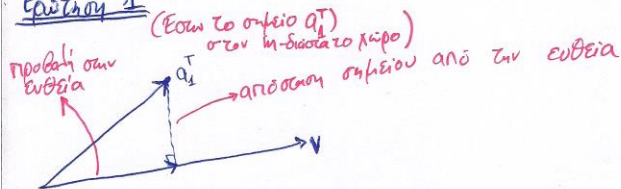
Δεύτερη Σειρά Ασκήσεων

Ευάγγελος Τζώρτζης

AM: 3088

Εξόρυξη Δεδομένων / 2^η σειρά Ασκήσεων / Ευάγγελος Τζώρτζης AM: 3088

Ερώτηση 1



Λόγω του Πυθαγόρειου θεωρήματος, για το σημείο που περιγράφεται από το a_i^T και κάποιας γραμμής πάνω στην οποία βρίσκεται το διάνυσμα v , ισχύει ότι:

$$(a_i^T)^2 = (\text{προβολή})^2 + (\text{απόσταση από ευθεία})^2$$

$$(\text{απόσταση})^2 = a_{i1}^2 + a_{i2}^2 + a_{i3}^2 + \dots + a_{im}^2 - (\text{προβολή})^2 \quad ①$$

Για να ελαχιστοποιηθεί το άθροισμα των τετραγώνων των αποστάσεων των σημείων από την ευθεία, μπορούμε να ελαχιστοποιήσουμε το $\sum_{i=1}^n (a_{i1}^2 + a_{i2}^2 + \dots + a_{im}^2) \quad ②$ μείον το άθροισμα των τετραγώνων των προβολών των σημείων στη γραμμή. Όμως το ② άθροισμα είναι σταθερό, άρα μπορούμε να μεγιστοποιήσουμε το άθροισμα των τετραγώνων των προβολών των σημείων στην ευθεία.

Το μήκος της προβολής της γραμμής a_i του πίνακα A πάνω στο μοναδιαίο διάνυσμα v είναι $|a_i \cdot v|$. Άρα το άθροισμα των τετραγώνων των προβολών είναι $|A \cdot v|^2$. Η βέλτιστη ευθεία είναι αυτή που μεγιστοποιεί το $|A \cdot v|^2$ και επομένως ελαχιστοποιεί το άθροισμα των τετραγώνων των αποστάσεων των σημείων από τη γραμμή.

Έχοντας τα παραπάνω στο νου, ορίζουμε το πρώτο singular vector v_1 του A , το οποίο είναι διάνυσμα στήλης, ως τη καλύτερη ευθεία από την αρχή των αξόνων για τα m στοιχεία n διάστασης τα οποία είναι οι γραμμές του πίνακα A . Έτσι έχουμε $v_1 = \arg \max_{|v|=1} |A \cdot v|$.

Ερώτηση 2

$$A. \quad P(X=x) = (a-1) \cdot x^{-a}$$

$$P(a|x) = \frac{P(x|a) \cdot P(a)}{P(x)}$$

$$P(x_i|a) = (a-1) \cdot x_i^{-a}$$

Πιθανότητα όλων των σημείων (υποθέτουμε ανεξάρτητα των x)

$$P(x|a) = \prod_{i=1}^n P(x_i|a) = \prod_{i=1}^n (a-1) \cdot x_i^{-a} = L(a)$$

Log Likelihood:

$$LL(a) = \log \left[\prod_{i=1}^n (a-1) x_i^{-a} \right] = \sum_{i=1}^n \log(a x_i^{-a} - x_i^{-a})$$

$$\frac{\partial LL(a)}{\partial(a)} = \sum_{i=1}^n \frac{1}{\frac{a}{x_i^a} - \frac{1}{x_i^a}} \left[\left(\frac{a}{x_i^a} \right)' - \left(\frac{1}{x_i^a} \right)' \right] =$$

$$= \sum_{i=1}^n \frac{1}{\frac{a}{x_i^a} - \frac{1}{x_i^a}} \left[\frac{x_i^a - a \cdot \ln(x_i) \cdot x_i^a}{x_i^{2a}} - \ln(x_i) \cdot x_i^{-a} \cdot (-1) \right] =$$

$$= \sum_{i=1}^n \frac{x_i^a}{a-1} \cdot \left[\frac{1-a \cdot \ln(x_i)}{x_i^a} + \frac{\ln(x_i)}{x_i^a} \right] =$$

$$= \sum_{i=1}^n \frac{1-a \ln(x_i) + \ln(x_i)}{a-1} = 0 \Rightarrow \frac{1}{a-1} \sum_{i=1}^n [1-a \ln(x_i) + \ln(x_i)] = 0 \Leftrightarrow$$

$$\frac{1}{a-1} \left[n - a \sum_{i=1}^n \ln(x_i) + \sum_{i=1}^n \ln(x_i) \right] = 0 \Leftrightarrow \frac{1}{a-1} \left[n + (1-a) \sum_{i=1}^n \ln(x_i) \right] = 0 \Leftrightarrow$$

$$\frac{n}{a-1} - \sum_{i=1}^n \ln(x_i) = 0 \Leftrightarrow \frac{n}{a-1} = \sum_{i=1}^n \ln(x_i) \Leftrightarrow a = 1 + n \cdot \left[\sum_{i=1}^n \ln(x_i) \right]^{-1}$$

Εξορκισμός Δεδομένων / 2^η Σειρά Ασκήσεων / Εισαγωγή Τυπικής AM: 3088

Ερώτηση 2

$$B. \quad P(X_i|L1) = (a_1 - 1) X_i^{-a_1}, \quad P(X_i|L2) = (a_2 - 1) X_i^{-a_2}$$

$$P(X_i) = \pi_1 \cdot P(X_i|L1) + \pi_2 \cdot P(X_i|L2) \quad \boxed{\frac{\partial P(X_i|L1)}{\partial a_1} = \frac{\ln(X_i) - \ln(X_i) \cdot a_1 + 1}{X_i^{a_1}}}$$

$$P(L1|X_i) = \frac{\pi_1 \cdot P(X_i|L1)}{P(X_i)} = \delta_{iL1}, \quad P(L2|X_i) = \delta_{iL2}$$

$$LL(\theta) = \sum_{i=1}^n \log(\pi_1 \cdot P(X_i|L1) + \pi_2 \cdot P(X_i|L2)), \quad \theta = (a_1, a_2, \pi_1, \pi_2)$$

Για το a_1 θα υπολογίσω το $\frac{\partial LL(\theta)}{\partial a_1} = 0$, και για το a_2 $\frac{\partial LL(\theta)}{\partial a_2}$

$$\frac{\partial LL(\theta)}{\partial a_1} = \sum_{i=1}^n \frac{1}{\pi_1 \cdot P(X_i|L1) + \pi_2 \cdot P(X_i|L2)} \cdot \pi_1 \cdot \frac{\partial P(X_i|L1)}{\partial a_1} =$$

$$= \sum_{i=1}^n \frac{\pi_1}{P(X_i)} \cdot \left(\frac{\ln(X_i) - \ln(X_i) \cdot a_1 + 1}{X_i^{a_1}} \right) = \sum_{i=1}^n \frac{\pi_1}{P(X_i)} \cdot (\ln(X_i) \cdot (a_1 - 1) + 1) \cdot X_i^{-a_1} =$$

$$= \sum_{i=1}^n \frac{\pi_1}{P(X_i)} \cdot \left[\underbrace{(-\ln(X_i)) \cdot (a_1 - 1) \cdot X_i^{-a_1}}_{P(X_i|L1)} + X_i^{-a_1} \right] = \sum_{i=1}^n \left[\frac{\pi_1 \cdot P(X_i|L1)}{P(X_i)} \cdot (-\ln(X_i)) \right] + \sum_{i=1}^n \frac{X_i^{-a_1} \cdot \pi_1}{P(X_i)} =$$

$$= \sum_{i=1}^n -\ln(X_i) \cdot \delta_{iL1} + \sum_{i=1}^n \frac{\underbrace{P(X_i|L1)}_{(a_1-1) \cdot X_i^{-a_1}} \cdot \pi_1}{P(X_i)} = 0 \Leftrightarrow$$

$$\Rightarrow \sum_{i=1}^n \frac{\delta_{iL1}}{(a_1 - 1)} = \sum_{i=1}^n \ln(X_i) \cdot \delta_{iL1} \Rightarrow \frac{1}{(a_1 - 1)} = \frac{\sum_{i=1}^n \ln(X_i) \delta_{iL1}}{\sum_{i=1}^n \delta_{iL1}} \Rightarrow$$

$$\Rightarrow a_1 - 1 = \frac{\sum_{i=1}^n \delta_{iL1}}{\sum_{i=1}^n \ln(X_i) \delta_{iL1}} \Rightarrow$$

$$a_1 = 1 + \frac{\sum_{i=1}^n \delta_{iL1}}{\sum_{i=1}^n \ln(X_i) \delta_{iL1}}$$

παράγωγο

$$a_2 = 1 + \frac{\sum_{i=1}^n \delta_{iL2}}{\sum_{i=1}^n \ln(X_i) \delta_{iL2}}$$

Εξορμή Αεδομένων / 2^η Σειρά Ασκήσεων / Ελάγγελος Τζατζής ΑΜ: 3088

Ερώτηση 2

B. (συνέχεια)

Υποδοχιστος των π_1, π_2 με χρήση πολλαπλασιαστών Lagrange

• λοχσα $\pi_1 + \pi_2 = 1$

$$f(\theta, \lambda) = LL(\theta) - \lambda(\pi_1 + \pi_2 - 1) \quad \frac{\partial f}{\partial \lambda} = 0 \Rightarrow \pi_1 + \pi_2 = 1$$

$$\frac{\partial f}{\partial \pi_1} = \sum_{i=1}^n \frac{1}{P(X_i)} \cdot P(X_i | L_1) - \lambda = 0 \Rightarrow$$

$$\lambda = \sum_{i=1}^n \frac{P(X_i | L_1)}{P(X_i)} \Rightarrow \pi_1 \cdot \lambda = \sum_{i=1}^n \frac{\pi_1 \cdot P(X_i | L_1)}{P(X_i)} = \sum_{i=1}^n P(L_1 | X_i) \quad (1)$$

$$\text{αντιστοιχα: } \pi_2 \cdot \lambda = \sum_{i=1}^n P(L_2 | X_i) \quad (2)$$

Ανο ① και ②:

$$(\pi_1 + \pi_2) \lambda = \sum_{i=1}^n [P(L_2 | X_i) + P(L_1 | X_i)] = n \Rightarrow \boxed{\lambda = n}$$

$$\text{αρα: } \begin{cases} \pi_1 = \frac{1}{n} \sum_{i=1}^n P(L_1 | X_i) \\ \pi_2 = \frac{1}{n} \sum_{i=1}^n P(L_2 | X_i) \end{cases}$$