

Εξόρυξη Δεδομένων

Πρώτη Σειρά Ασκήσεων

Ταφλαμπάς Νικόλαος 4500

Άσκηση 1

Ερώτημα Α)

1)

Ο αλγόριθμος, αρχικοποιεί μιας λίστα K θέσεων για την συλλογή τυχαίων δειγμάτων (δεξαμενή ή *reservoir*). Έπειτα, αρχίζει να διαβάζει την σειρά δεδομένων μας, τελικού και αγνώστου μεγέθους N .

Τα πρώτα K δεδομένα της σειράς καταγράφονται στην δεξαμενή. Τα υπόλοιπα δεδομένα, έχουν πιθανότητα K/i , να αντικαταστήσουν κάποιο στοιχείο της δεξαμενής, με i = τον αριθμό των στοιχείων που έχουμε προσπέραση.

2)

Θεωρώ:

N = "αριθμος στοιχειων για δειγματοληψια"

R = "Δεξαμενη/Reservoir"

K = "μεγεθος δεξαμενής"

I = "Μετρητης περασμένων στοιχειων"

X = "Δειγμα"

Αρχικά, θα υπολογίσουμε μερικές πιθανότητες:

$$P(\text{"να επιλεγεί ένα στοιχείο"}) = K/I = P(X \in R)$$

$$P(\text{"να αντικαταστήσει ένα συγκεκριμένο στοιχείο της δεξαμενής"}) = 1/K = P(X \rightarrow R[x])$$

Αρχικά, ας δούμε για $I > K$. Αρκεί να δείξουμε ότι η πιθανότητα για ένα στοιχείο να είναι στο τέλος του αλγορίθμου μέσα στην δεξαμενή, πρέπει να είναι K/N

$P(\text{"να ανήκει στο τέλος στην δεξαμενή"}) =$
 $P(\text{"να επιλεγθεί"}) * P(\text{"να μην αλλάξει σε κάποιο μελλοντικό γύρο"}) =$
 $P(X \in T) * [P(\text{"να μην επιλεγθεί το στοιχείο i"}) \vee P(\text{"να επιλεγθεί αλλά να μην αλλάξει το στοιχείο μας"})]$

$$P(X \in R') = P(X \in R) * \prod_{i=x+1}^N [(1 - P(X \in R) + (P(X \in R) * P(X \rightarrow R[x])))]$$

$$\begin{aligned}
 P(X \in R') &= \frac{k}{i} * \left[\left(1 - \frac{k}{i+1} \right) + \left(\frac{k}{i+1} * \frac{k-1}{k} \right) \right] * \dots \\
 &= \frac{k}{i} * \left[\frac{i}{i+1} * \frac{i+1}{i+2} * \dots * \frac{N-1}{N} \right] \\
 &= \frac{k}{N}
 \end{aligned}$$

Αρά όντως ισχύει. Τώρα για τα πρώτα K στοιχεία, θεωρούμε:

$$\frac{K}{i} = 1 \Leftrightarrow K = i$$

Έπειτα, κάνοντας την από πάνω διαδικασία για $K/i = 1$, καταλήγουμε με:

$$P(X \in R') = \frac{i}{N} \Leftrightarrow \frac{K}{N}$$

Πράγμα που ισχύει για κάθε ένα από τα πρώτα K στοιχεία. Έτσι καταλήγουμε ότι ΟΛΑ τα στοιχεία έχουν πιθανότητα K/N να επιλεγθούν.

Ερώτημα Β)

Σαφώς, θα βασιστούμε στον κλασικό Reservoir Sampling για $K=1$.

Ο αλγόριθμος μας, αρχικά, ορίζει έναν μετρητή περασμένου βάρους, που θα κρατάει το άθροισμα όλων των βαρών των στοιχείων που διαβάζουμε.

Με το πέρασμα ενός στοιχείου, αυξάνουμε τον μετρητή κατά w_i και ελέγχουμε την πιθανότητα w_i/W όπου w_i το βάρος του στοιχείου και W η τιμή του μετρητή μας. Άμα πέτυχει, τότε αλλάζουμε το στοιχείο του Reservoir στο στοιχείο που ελέγχουμε.

Το ερώτημα, φυσικά, είναι αν δουλεύει. Αν όντως λειτουργεί ορθά, τότε η πιθανότητα ΚΑΘΕ στοιχείου, πρέπει να είναι w_i/W' , όπου W' είναι το σύνολο όλων των βαρών μας

Αρχικά, ας υποθέσουμε τυχαίο στοιχείο i , με ακόλουθες πιθανότητες:

$$P(\text{"να επιλεγθαι το στοιχειο } i\text{"}) = \frac{w_i}{w_i + \sum_{k < i} w_k} = \frac{w_i}{\sum_{k \leq i} w_k}$$

Τώρα έστω στοιχείο j που έρχεται μετά το i . Ξέροντας βάση το από πάνω την πιθανότητα να επιλεγθεί, θα βρούμε την πιθανότητα να ΜΗΝ επιλεγθεί:

$$P(\text{να ΜΗΝ επιλεγθαι το στοιχειο } j) = 1 - \frac{w_j}{\sum_{k \leq j} w_k}$$

$$= \frac{\sum_{k \leq j} w_k - w_j}{\sum_{k \leq j} w_k} = \frac{\sum_{k \leq i} w_k}{\sum_{k \leq j} w_k}$$

Αυτό μπορεί να γραφτεί για κάθε στοιχείο μετά το j. Άρα η πιθανότητα του i να μείνει στην δεξαμενή είναι:

$$P(i \in R') = P(\text{"να επιλεγθαι το στοιχειο } i\text{"}) * \prod_{k=i+1}^N P(\text{να ΜΗΝ επιλεγθει το στοιχειο } k)$$

$$= \frac{w_i}{\sum_{k \leq i} w_k} * \frac{\sum_{k \leq i} w_k}{\sum_{k \leq j} w_k} * \dots * \frac{\sum_{k \leq N-1} w_k}{W'} = \frac{w_i}{W'}$$

Αυτό μπορεί να αποδειχθεί για κάθε στοιχείο i, συνεπώς ο αλγόριθμος μας δουλεύει σωστα!