

# Υπογραμμικοί Αλγόριθμοι

## Μάθημα 4 - 22/10/2019

Επιμέλεια Διαφανειών: Αλέξανδρος Ζέρντεβ

### Γραμμική σκιαγράφηση

initialize() :  $x \leftarrow \vec{0}$

update(i,Δ) :  $x_i \leftarrow x_i + \Delta$

Σχεδιάζω  $\Phi$  και κρατάω το  $y = \Phi x$  και το  $\Phi$  στη μνήμη.

- Χώρος:  $\Phi \in \mathbb{R}^{m \times n}$  + χώρος για το  $y$
- Χρόνος Ανανέωσης την  $t$ -οστή στιγμή:

$$- y(t) = \Phi x^{(t)} = \Phi(x^{(t-1)} + \Delta e_i) = \underbrace{y(t-1)}_{\Phi x^{(t-1)}} + \underbrace{\Delta \cdot \Phi e_i}_{i\text{-th column}}$$

-  $e_i$  είναι το διάνυσμα που έχει 1 στην  $i$ -οστή θέση, και 0 οπουδήποτε αλλού.

Για παράδειγμα,  $e_3 = \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}$

- Άρα Ο χειρότερος χρόνος για να υπολογίσεις μια στήλη του  $\Phi$

- Δεν μπορούμε να σώσουμε όλο το  $\Phi$ . Είναι πολύ μεγάλο, μπορεί μέχρι  $mn$  λέξεις!
- Για το  $\Phi$  χρειαζόμαστε συνήθως μία ή δύο μεθόδους: μία μέθοδο που υπολογίζει στο στοιχείο  $\Phi_{i,j}$ , ή/και μία μέθοδο που υπολογίζει την στήλη  $\Phi_i = \Phi e_i$ .

## AMS σκιαγράφημα:

$$\underbrace{y}_{\in \mathbb{R}^{\Theta(\varepsilon^{-2})}}} = \underbrace{\begin{bmatrix} \pm 1 & \pm 1 & \pm 1 & \pm 1 & \pm 1 & \pm 1 \\ \pm 1 & \pm 1 & \pm 1 & \pm 1 & \pm 1 & \pm 1 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ \pm 1 & \pm 1 & \pm 1 & \pm 1 & \pm 1 & \pm 1 \end{bmatrix}}_{\Phi} \begin{bmatrix} t_1 \\ t_2 \\ \cdot \\ t_n \end{bmatrix}$$

$\Theta(\varepsilon^{-2})$  ανά 4 ανεξάρτητες συναρτήσεις  $\sigma_r[n] \rightarrow \{-1, 1\}$  για τον  $r$ -στό μετρητή.

Για να πάρω την  $i$ -στήλη κοιτάω τα  $\sigma_1(i), \sigma_2(i), \sigma_3(i) \dots \rightarrow \begin{bmatrix} \sigma_1(i) \\ \sigma_2(i) \\ \sigma_3(i) \\ \vdots \\ \sigma_{\Theta(\varepsilon^{-2})}(i) \end{bmatrix}$

- Υπολογίζω κάθε  $\sigma_r(i)$  σε  $O(1)$ , όλη τη στήλη σε  $O(\varepsilon^{-2})$
- Χώρος για το  $\Phi$  είναι  $O(\varepsilon^{-2} \lg n)$  δυφία (ή  $O(\varepsilon^{-2})$  λέξεις)

## CountMin σκιαγράφημα:

- Στην προηγούμενη διάλεξη είδαμε το CountMin σκιαγράφηματος σαν  $R = \Theta(\log n)$  γραμμές από μετρητές, και κάθε  $i \in [n]$  ήταν "συνδεδεμένο" με ένα μετρητή σε κάθε γραμμή.
- # γραμμών (χώρος ανά συνάρτησή)
  - Σε μορφή πινάκων το  $\Phi$  του Countmin σκιαγράφημα είναι  $R$  πίνακες κάθετα τοποθετημένοι, έτσι ώστε κάθε πίνακας να υποδεικνύει ποιά στοιχεία συνεισφέρουν σε ποιον μετρητή.

$$\begin{array}{c}
O(\frac{1}{\varepsilon}) \\
- \\
O(\frac{1}{\varepsilon}) \\
- \\
O(\frac{1}{\varepsilon})
\end{array}
\begin{array}{c}
\left[ \begin{array}{cccccccc}
. & . & . & i & . & j & . & . & . \\
. & . & . & 0 & 1 & . & 1 & . & . & . \\
. & . & . & 1 & 0 & . & 0 & . & . & . \\
. & . & . & 0 & 0 & . & 0 & . & . & . \\
- & - & - & - & - & - & - & - & - & - \\
. & . & . & 0 & 1 & . & 0 & . & . & . \\
. & . & . & 0 & 0 & . & 1 & . & . & . \\
. & . & . & 1 & 0 & . & 0 & . & . & . \\
- & - & - & - & - & - & - & - & - & - \\
. & . & . & 0 & 1 & . & 0 & . & . & . \\
. & . & . & 1 & 0 & . & 1 & . & . & . \\
. & . & . & 0 & 0 & . & 0 & . & . & . \\
\end{array} \right]
\end{array}
\begin{array}{c}
\left[ \begin{array}{c}
x \\
x_1 \\
x_2 \\
. \\
. \\
. \\
. \\
. \\
. \\
. \\
. \\
x_n
\end{array} \right]
\end{array}$$

$$\Phi = \begin{bmatrix} \Phi^{(1)} \\ - \\ \Phi^{(2)} \\ \dots \\ \Phi^{(R)} \end{bmatrix}$$

Το που είναι οι άσοι στη γραμμή  $l$ -οστή γραμμή του  $\Phi^{(r)}$ ,  $r \in [R]$  δείχνει ποιά στοιχεία  $i \in [n]$  συνεισφέρουν στον  $l$ -οστό μετρητή στην  $r$ -οστή επανάληψη. επαναλήψεις.

### ε-βαρέα στοιχεία

Υπενθύμιση:Θέλουμε να κρατήσουμε δομή δεδομένων για  $x \in R^n$  η οποία βρίσκει μια λίστα  $L$  η οποία περιέχει όλα τα  $i$  ώστε  $|x_i| \geq \varepsilon \|x\|_1$  και κανένα  $i$  ώστε  $|x_i| < \frac{\varepsilon}{2} \|x\|_1$ .

- cash-register (2η Σειρά Ασκήσεων)

$$- \Delta = 1$$

- strict turnstile

$$- x_i \geq 0$$

- Turnstile

–  $x_i \in R$

### CountMin για strict turnstile (προηγούμενη διάλεξη)

- Χώρος  $O(\varepsilon^{-1} \lg n)$  λέξεις
- Πιθανότητα λαθους  $\frac{1}{n^{100}}$

Σήμερα: Υπάρχει ντετερμινιστικός αλγόριθμος για το πρόβλημα  $\varepsilon$ -βαρέων στοιχείων με χώρο

$$\varepsilon^{-2} \cdot \left( \frac{\lg n}{\lg \frac{1}{\varepsilon} + \lg \lg n} \right)^2$$

λέξεις.

Σημειώνουμε ότι είναι απαραίτητη η τετραγωνική εξάρτηση στο  $\varepsilon$  για ντετερμινιστικούς αλγόριθμους.

#### Ορισμός:

Ένας πίνακας  $A \in R^{m \times n}$  λέγεται  $\varepsilon$ -(μη συνοχικός) αν  $\|A_i\|_2 = 1$  και  $|\langle A_i, A_j \rangle| \leq \varepsilon$  ή

$$A^T A = I + \varepsilon J \in R^{n \times n} \quad J = \begin{bmatrix} \in [-1, 1] & \in [-1, 1] & \in [-1, 1] & \in [-1, 1] & \in [-1, 1] & \in [-1, 1] \\ \in [-1, 1] & \in [-1, 1] & \in [-1, 1] & \in [-1, 1] & \in [-1, 1] & \in [-1, 1] \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ \in [-1, 1] & \in [-1, 1] & \in [-1, 1] & \in [-1, 1] & \in [-1, 1] & \in [-1, 1] \end{bmatrix}$$

**Λήμμα:** Έστω ένας  $A \in R^{m \times n}$  ( $\varepsilon/2$ )-(μη συνοχικός) πίνακας τον οποίο μπορώ να αποθηκεύσω σε χώρο  $s_1$  μπορώ να βρω μια στήλη του σε χρόνο  $t_u$ . Τότε υπάρχει ντετερμινιστικός αλγόριθμος ροής για τα  $\varepsilon$ -βαρέα στοιχεία με χώρο  $S + m$ , χρόνο ανανέωσης  $t_u + m$  και πολυπλοκότητα χρόνου ερωτήματος  $O(nt_u)$ , που έρχεται.

Απόδειξη: Θα χρησιμοποιήσω  $\Phi = A$ . Θα υπολογίσω κάθε  $i \in [n]$  μια εκτίμηση του  $x_i$ . Αν η εκτίμηση είναι  $\geq \varepsilon \|x\|_1$  θα βάλω το  $i$  στη λίστα  $L$ . Στο strict turnstile μοντέλο μπορώ να υπολογίσω το  $\|x\|_1$  με ένα μετρητή, διότι  $\sum_{i=1}^n x_i = \sum_{i=1}^n |x_i| = \|x\|_1$ .

Ο εκτιμητής για κάθε  $i \in [n]$  είναι

$$x'_i = \langle y, A_i \rangle = \langle Ay, Ae_i \rangle = (Ax)^T Ae_i = x^T (A^T A) e_i = x^T (I + \varepsilon J) e_i =$$



το μοναδικό μη μηδενικό στοιχείο, παίρνουμε έναν  $\varepsilon$ - (μη συνοχικό) με  $\Theta(\varepsilon^{-2} \lg n)$  γραμμές. Αλλά είναι δύσκολο να τον βρούμε! Θα κάνουμε κάτι πιο εξεζητημένο χρησιμοποιώντας τους περίφημους κώδικες Reed-Solomon, ή αλλιώς πολυώνυμα χαμηλού βαθμού.

### Κώδικες Read-Solomon

Πάρε  $q = t$ , έστω  $q$  πρώτος αριθμός.

Έστω μια αρίθμηση όλων των πολυωνύμων βαθμού  $d$  από το  $\mathbb{Z}_q$  στο  $\mathbb{Z}_q$ .

$$g(x) = \left( \sum_{i=0}^d a_i x^i \right) \pmod q$$

$$g : \{0, 1, \dots, q-1\} \rightarrow \{0, 1, \dots, q-1\}$$

$$\# \text{ πλήθος πολυωνύμων βαθμού } \leq d \text{ στο } \mathbb{Z}_q = q^{d+1} = n$$

Για να κατασκευάσω την πρώτη στήλη του πίνακα  $A$ , αποτιμώ το  $g_1$  πάνω σε όλα τα στοιχεία του  $\mathbb{Z}_q$ , παίρνοντας ένα  $q$ -διάστατο διάνυσμα. Εν συνεχεία, αντικαθιστώ κάθε στοιχείο (έστω το  $g_1(\ell)$ ) αυτού του  $q$ -διάστατου διανύσματος με ένα  $q$ -διάστατο διάνυσμα το οποίο έχει 1 στη θέση  $g_1(\ell)$ . Έτσι παίρνω ένα διάνυσμα με  $q \cdot q = q^2$  στοιχεία. Μπορώ να κάνω το ίδιο για όλα τα  $i \in [n]$ , αν  $n = q^{d+1}$ . Προσέξτε ότι δε χρειάζεται να γράψω τον πίνακα κάτω. Μπορώ να υπολογίσω την  $i$ -οστή στήλη ή ένα στοιχείο του  $A$  αποτιμώντας απλά το  $g_i$  πάνω στο  $\mathbb{Z}_q$ .

$$\begin{bmatrix} g_1(0) \\ g_1(1) \\ g_1(2) \\ \vdots \\ g_1(q-1) \end{bmatrix} = \begin{bmatrix} 3 \\ 9 \\ \vdots \\ q-6 \\ 37 \end{bmatrix} \Rightarrow \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \\ - \\ 0 \\ 0 \\ 0 \\ 1 \\ 0 \\ - \\ 0 \\ 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}$$

Κάθε στήλη αντιστοιχεί σε ένα πολυώνυμο. Το που είναι το μη μηδενικό της στοιχείο στο  $r$ -οστό πίνακα εξαρτάται από το  $g_r(l)$

Από τα παραπάνω έχουμε

$$\langle A_i, A_j \rangle = \frac{\#r : g_i(r) = g_j(r)}{t} \leq \frac{d}{t}$$

Για να πάρουμε έναν  $\varepsilon$ -(μη συνοχικό πίνακα) θα πρέπει  $\frac{d}{t} \leq \varepsilon$ . Συνολικά, έχουμε

- $q = t$
- $\varepsilon = \frac{d}{q}$
- $m = q^2$
- $n = q^{d+1}$ ,

από όπου θέτοντας  $q = \Theta\left(\varepsilon^{-1} \frac{\log n}{\log \log n + \log(1/\varepsilon)}\right)$  και  $d = \Theta\left(\frac{\log n}{\log q}\right)$ , παίρνουμε μία άπειρη ακολουθία  $\varepsilon$ -(μη συνοχικών) πινάκων με

$$m = \Theta\left(\varepsilon^{-2} \left(\frac{\log n}{\log \log n + \log(1/\varepsilon)}\right)^2\right).$$

### Count Sketch για $\ell_2$ $\varepsilon$ -βαρέα στοιχεία

**Πρόβλημα των  $\ell_2$  βαρέων στοιχείων :** Ένα  $i \in [n]$  λέγεται ( $\ell_2$ )  $\varepsilon$ -βαρέο αν  $x_i^2 \geq \varepsilon \cdot \|x\|_2^2 = \varepsilon \cdot \sum_{i=1}^n x_i^2$

\* Για να λύσω το πρόβλημα των  $\ell_2$  βαρέων στοιχείων να βρω για κάθε  $i \in [n]$  ένα  $x'_i$  ώστε  $x'_i = x_i \pm \sqrt{\frac{\varepsilon}{2}} \|x\|_2$ , ή ισοδύναμα

$|x_i - x'_i| \leq \sqrt{\frac{\varepsilon}{2}} \|x\|_2$ . Αν έχω βρει ένα τέτοιο  $x'_i$  θα λέω ότι το  $i$  εκτιμήθηκε σωστά.

Ιδέα: Συνδυάζουμε το CountMin με το AMS. Τοποθετούμε ένα τυχαίο πρόσημο σε κάθε  $x_i$  πριν το στείλουμε στο μετρητή.

### Πολυπλοκότητα Αλγορίθμου Count Sketch:

- Χώρος:  $O(\varepsilon^{-1} \lg n)$  λέξεις
- Χρόνος Ανανέωσης:  $O(\lg n)$
- Χρόνος Ερώτησης:  $O(n \lg n)$
- Πιθανότητα λάθους:  $n^{-10}$

Το γραμμικό του σκιαγράφημα μοιάζει δηλαδή κάπως έτσι:

$$R = \Theta(\lg n) \begin{bmatrix} \Theta(\varepsilon^{-1}) & \begin{bmatrix} \cdot & \cdot & \cdot & 0 & \pm 1 & \cdot & \pm 1 & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \pm 1 & 0 & \cdot & 0 & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & 0 & 0 & \cdot & 0 & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & 0 & \pm 1 & \cdot & 0 & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & 0 & 0 & \cdot & \pm 1 & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \pm 1 & 0 & \cdot & 0 & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & 0 & \pm 1 & \cdot & 0 & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \pm 1 & 0 & \cdot & \pm 1 & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & 0 & 0 & \cdot & 0 & \cdot & \cdot & \cdot \end{bmatrix} \end{bmatrix}$$

Αν το δούμε σαν  $R = \Theta(\log n)$  σειρές από  $\Theta(\varepsilon^{-1})$  μετρητές, έχουμε το εξής.

$$\begin{bmatrix} C_{1,1} & C_{1,2} & \dots & C_{1,\varepsilon^{-1}} \\ C_{2,1} & C_{2,2} & \dots & C_{2,\varepsilon^{-1}} \\ \vdots & \vdots & \vdots & \vdots \\ C_{R,1} & C_{R,2} & \dots & C_{R,\varepsilon^{-1}} \end{bmatrix}$$

$h_1 : [n] \rightarrow [\Theta(\varepsilon^{-1})]$  2-ανεξάρτητες

$\sigma_r : [n] \rightarrow \{-1, +1\}$  4-ανεξάρτητες

Όταν βλέπω  $(i, \Delta) \quad \forall r \in [R] : C_{r,h_r(i)} \leftarrow C_{r,h(i)} + \sigma_r(i)\Delta$

Για  $r \in [R]$  :

$$C_{r,h_r(i)} = \sum_{j:h_r(j)=h_r(i)} \sigma_r(j)x_j \quad \rightarrow \quad C_{r,h(i)} \cdot \sigma_r(i) = x_i + \sum_{j \neq i: h_r(j)=h_r(i)} x_j \cdot \sigma_r(j) \cdot \sigma_r(i)$$

Για κάθε  $i \in [n]$  θα πάρω μια εκτίμηση για το  $q_i$  από την  $r$ -οστή επανάληψη, και στη συνέχεια θα πάρω τη διάμεσο των  $R$  αυτών εκτιμήσεων. Για να είναι ο εκτιμητής στην  $r$ -οστή επανάληψη σωστός με πιθανότητα  $9/10$  (δηλαδή το σφάλμα της τάξης  $\sqrt{\varepsilon/2}\|x\|_2$ ) αρκεί

$$\mathbb{P} \left[ \left| \sum_{j \neq i} \sigma_r(i)\sigma_r(j)x_j \right|^2 \geq \frac{\varepsilon}{2}\|x\|_2^2 \right] \leq \frac{1}{10}$$

Υπολογίζω την αναμενόμενη τιμή του τετραγώνου του σφάλματος και κάνω Markov.

$$\mathbb{E} \left[ \sum_{j \neq i: h_r(i)=h_r(j)} \sigma_r(i)\sigma_r(j)x_j \right]^2 = \mathbb{E} \left[ \left( \sum \delta_j \eta_j x_j \right)^2 \right] = \sum_{j \neq i} \mathbb{E} [\delta_j, \delta_{j'}, \eta_j, \eta_{j'}] x_j x_{j'} \stackrel{\text{why?}}{=} \sum_{j \neq i} x_j^2 \Theta(\varepsilon),$$



όπου

$$\delta_j = 1 \text{ αν } h_r(j) = h_r(i)$$

$$\eta_j = \sigma_r(j) \cdot \sigma_r(i)$$

Όπως είπαμε, εφόσον έβγαλα την αναμενόμενη τιμή μπορώ να κάνω Markov για να δείξω ότι το λάθος δε θα ξεφύγει πολύ.

Ένας δεύτερος περιγραφικός τρόπος να δεις κανείς ότι η εκτίμηση είναι σωστή είναι ο ακόλουθος. Ο μετρητής  $i \in [n]$  αναμένει να 'δεχτεί'  $O(\|x\|_2^2) \ell_2^2$  μάζα από τα υπόλοιπα  $j \neq i$ , άρα από Markov θα έχει πιθανότητα 9/10 το πολύ  $O(\|x\|_2^2) \ell_2^2$  μάζα (ας την πούμε  $M$  για τη συνέχεια) από τα υπόλοιπα  $j \neq i$ . Δεδομένου ότι αυτό θα συμβεί, ποια η πιθανότητα τα τυχαία πρόσημα να ανεβάσουν πολύ το λάθος, πολύ πάνω από  $M$ ; Εφόσον αναμένουμε τα τυχαία πρόσημα να μας δώσουν ακριβώς  $M$  (όπως στην απόδειξη του AMS), μπορώ πάλι να κάνω Markov για να δείξω ότι το λάθος δε θα ξεφύγει πολύ. Η διαφορά της απόδειξης αυτής από την προηγούμενη είναι ότι πρώτα υπολογίζω τη συνεισφορά από τα  $h$ , και μετά από τα  $\sigma$ , ενώ στην άλλη τα κάνω και τα δύο μαζί.

Συνεχίζοντας εκεί που μείναμε, με πιθανότητα 9/10 ο εκτιμητής του  $x_i$  στην  $r$ -οστή επανάληψη είναι σωστός. Εφόσον επαναλαμβάνουμε  $\Theta(\log n)$  φορές, μπορούμε να ριζούμε σε  $n^{-101}$  την πιθανότητα λάθους<sup>1</sup> για ένα  $i$ . Άρα η πιθανότητα να υπάρχει ένα  $i$  που εκτιμήθηκε λάθος είναι το πολύ  $n \cdot n^{-101} = n^{-100}$ .

	Χώρος	Χρόνος Ερωτήματος	Χρονος Ανανέωσης	Πιθανότητα Λαθους
Count Sketch	$O(\varepsilon^{-1} \lg n)$	$O(n \lg n)$	$O(\lg n)$	$n^{-100}$
Count Sketch + Code	$O(\varepsilon^{-1} \lg n \lg \frac{1}{\varepsilon})$	$O(\varepsilon^{-1} \lg n \lg \frac{1}{\varepsilon})$	$O(\lg n \cdot \log(1/\varepsilon))$	$\varepsilon^{100}$

Θα χρησιμοποιήσουμε το εξής εργαλείο για να επιταχύνουμε το CountSketch, στην επόμενη δι-άλεξη.

### **Κώδικας του Spielman, 1996:**

Υπάρχει μια συνάρτηση από το  $\text{enc} : [n] \rightarrow \{0, 1\}^{20 \lg n}$  για την οποία ισχύει:

Αν έχω το  $\text{enc}(x)$  για κάποιο  $x \in [n]$  αλλά με το πολύ  $\frac{1}{3}(20 \lg n)$  ψηφία του  $\text{enc}(x)$  λάθος, μπορώ να βρώ το  $x$  σε  $O(\lg n)$  χρόνο.

---

<sup>1</sup>Η σταθερά στον έκθετη έρχεται με ένα πολλαπλασιαστικό κόστος στη σταθερά μέσα στο  $\Theta(\log n)$ . Για να κάνω την πιθανότητα λάθους  $n^{-200}$  πρέπει να διπλασιάσω τη σταθερά μέσα στο  $\Theta(\log n)$ , άρα και τον αριθμό των επαναλήψεων.