

Τροπογραμμένα ισότιμα

Markov's Inequality

Αν $X \geq \mu$ αρνητική η μεταβλητή $\Pr[X \geq t] \leq \frac{E[X]}{t}$ για $t > 0$

Chebyshov's Inequality

Αν X είναι τυχαία μεταβλητή με $\text{Var}[X] = \sigma^2$

$$\Pr[|X - E[X]| \geq \kappa \sigma] \leq \frac{1}{\kappa^2} \quad \text{για } \kappa \geq 1$$

Hoeffding Inequality

Αν $X = \sum X_i$ με X_i ανεξάρτητες τηλ. στο $\{a_i, b_i\}$ τότε για κάθε $t > 0$

$$\Pr[|X - E[X]| \geq t] \leq 2 e^{-\frac{2t^2}{\sum_i (b_i - a_i)^2}}$$

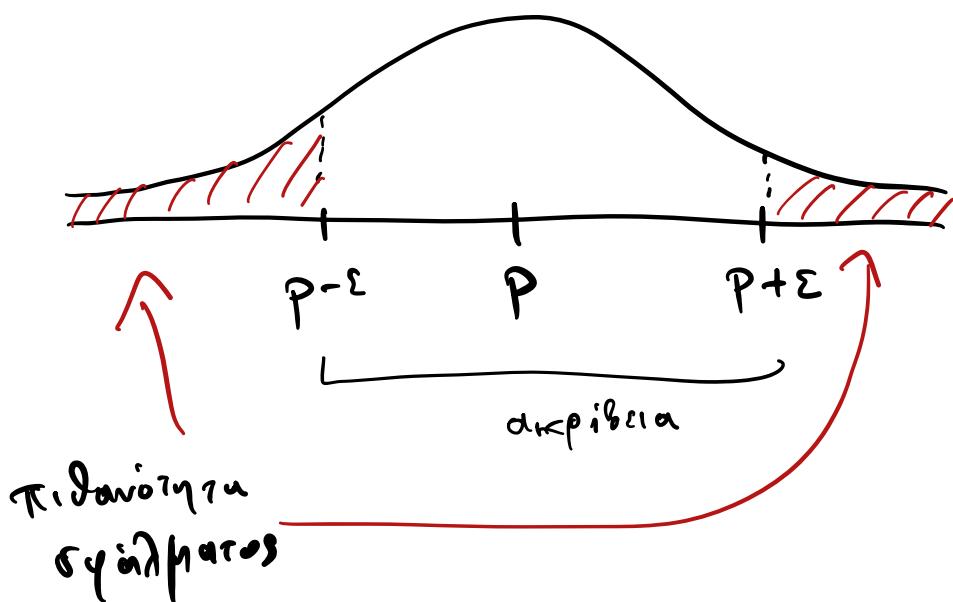
Πιθανότητα π ευσκέψιας p

$$\text{Av } \hat{p} = \frac{\sum_{i=1}^m x_i}{m}$$

$$\Pr[|\hat{p} - p| > \varepsilon] \leq e^{-\Omega(\varepsilon^2 m)}$$

$$M_\varepsilon \quad m = \Theta\left(\frac{1}{\varepsilon^2} \log(1/\delta)\right) \text{ διεργασία}$$

$$|\hat{p} - p| < \varepsilon \quad \mu_\varepsilon \text{ πιθανότητα } 1 - \delta$$



Εκτίμηση

Πλήρους

Διαφορετικών

Στοιχείων

Τρόβλημα: Μεγάλο stream από δεδομένα

x_1, \dots, x_N

Πότε είναι διαφορετικά;

Έστω τα διαφορετικά στοιχεία έχουν πλήρος η

Με $O(n)$ μήχανη

μπορώ να τα υπολογίσω
ακριβώς.

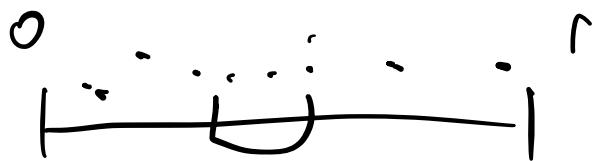
Τεχνική: Hash κάθε στοιχείου

$x_i \rightarrow h(x_i) \sim U[0, 1]$

$h(x_1) = h(x_8) \quad \text{αν}$

τότε $x_1 = x_8$

$h(x) = md5(x)$



$$Y = \min_{i=1}^n h(x_i) = \min_{i=1}^n s_i$$

or
 $s_i \sim U[0, 1]$

Claim: $E[Y] = \frac{1}{n+1}$

Anoferen

$$\begin{aligned} E[Y] &= \int_0^1 \Pr[Y=t] \cdot t \, dt \\ &= \int_0^1 (-\Pr[Y \geq t])' \cdot t \, dt \\ &= \int_0^1 \Pr[Y \geq t] \, dt \\ &= \int_0^1 (1-t)^n \, dt = \int_0^1 t^n \, dt \\ &= \frac{1}{n+1} \end{aligned}$$

D

$$n = \frac{1}{y} - 1$$

$$\text{Claim : } \text{Var}[Y] \leq \frac{1}{(n+1)^2}$$

$$\text{Var}[Y] = E[Y^2] - E[Y]^2$$

Ano δc13y .

$$\begin{aligned} E[Y] &= \int_0^1 \Pr[Y=t] t^2 dt \\ &= \int_0^1 n(1-t)^{n-1} t^2 dt \\ &= \frac{2}{(n+1)(n+2)} \leq \frac{2}{(n+1)^2} \end{aligned}$$

Ano Chebyshev

$$\Pr\left[\left| Y - \frac{1}{n+1} \right| \geq \frac{\kappa}{n+1} \right] \leq \frac{1}{\kappa^2}$$

Measuring the Variance

Εναριθμήστε το σειράς της γραφής.

$$Y_1, Y_2, \dots, Y_t$$

$$\bar{Y} = \frac{\sum_{i=1}^t Y_i}{t} \rightarrow \left\{ \begin{array}{l} E[\bar{Y}] = \frac{1}{t} \sum_{i=1}^t E[Y_i] \\ = \frac{1}{n+1} \\ \text{Var}[\bar{Y}] = \frac{1}{t^2} \text{Var}[\sum Y_i] \\ = \frac{\sum \text{Var}[Y_i]}{t^2} \\ = \frac{1}{t(n+1)^2} \end{array} \right.$$

Ano Chebyshev

$$\Pr \left[\left| \bar{Y} - \frac{1}{n+1} \right| \geq \frac{\kappa}{\sqrt{t(n+1)}} \right] \leq \frac{1}{\kappa^2}$$

A_v $\kappa = 2$ $\frac{1}{\sqrt{6}}$ $t = 1$

$$t = \frac{4}{\varepsilon^2}$$

$$\frac{1}{\delta \varepsilon^2}$$

$$\Pr \left[\left| \bar{Y} - \frac{1}{n+1} \right| \geq \varepsilon \frac{1}{(n+1)} \right] \leq \frac{1}{4}$$

'Appa μ_C $n, \text{davotyta}$ $\frac{3}{4}$

$$\frac{1}{n+1} (1-\varepsilon) \leq \bar{Y} \leq (1+\varepsilon) \frac{1}{n+1}$$

$$\hat{n} = \frac{1}{\bar{Y}} - 1$$

$$\frac{n+1}{1+\varepsilon} \leq \frac{1}{\bar{Y}} \leq \frac{n+1}{1-\varepsilon}$$

[Flajolet - Martin algorithm]

πιδανοίτα

1 - δ

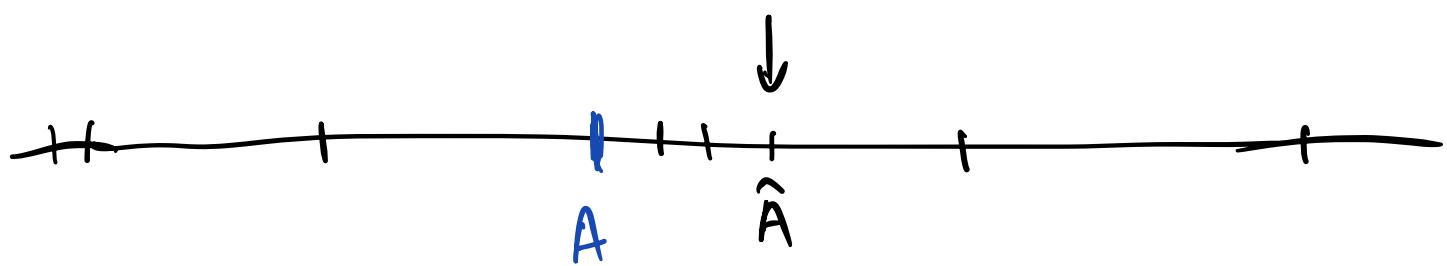
Γενική τεχνική για boosting
τη πιδανοίτας συστήματος.
εξηγήσου.

Έστω οτικός εξώς μεριδιας A
εξηγήσου \hat{A} τον ϵ -
wore $|\hat{A} - A| \leq \epsilon$ με πιδανοίτα $\frac{3}{4}$

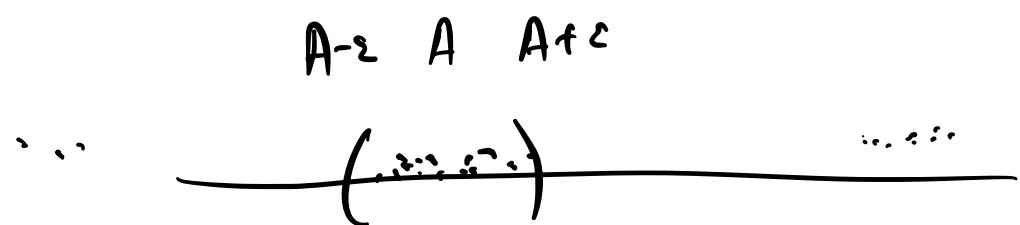
Μηρών να αντίστοιχου την οιδανότητα
με αντίστοιχη περιπάτωση του εκτιμήσι

$$\hat{A}_1, \hat{A}_2, \dots, \hat{A}_m$$

$$\text{median}(\hat{A}_1, \dots, \hat{A}_m)$$



Claim: Αν ην ρανού ανο $\frac{m}{2}$
των \hat{A}_i μερονού $|\hat{A}_i - A| \leq \varepsilon$
τότε $|\text{median}(\hat{A}_1, \dots, \hat{A}_m) - A| \leq \varepsilon$



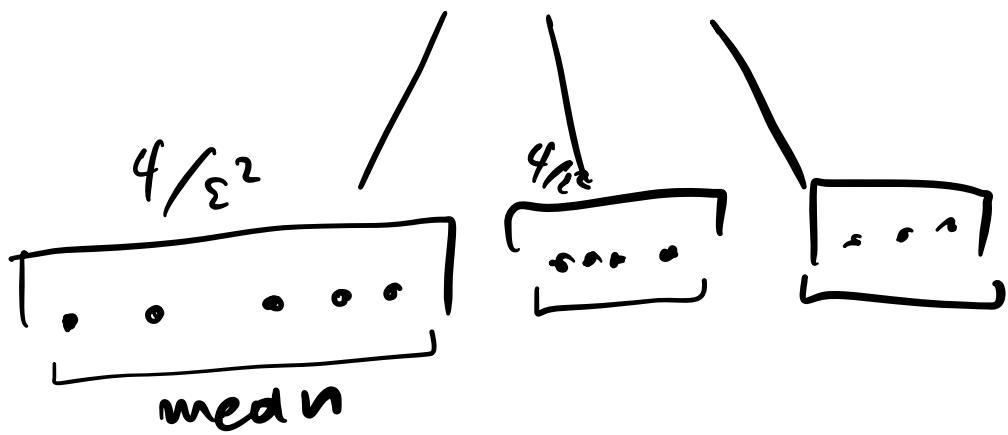
Έστω Z_i οιναίη και δεκτριά τη.
 του γεγονότος ότι το \hat{A}_i είναι
 κανό. Δηλαδή:

$$Z_i = \begin{cases} 1 & \text{αν } |\hat{A}_i - A| > \varepsilon \\ 0 & \text{αλλιώς} \end{cases}$$

$$E[Z_i] \leq \frac{1}{4}$$

$$\begin{aligned} & \Pr \left[\sum_{i=1}^m Z_i \geq \frac{m}{2} \right] \\ &= \Pr \left[\left| \sum_{i=1}^m Z_i - \frac{m}{4} \right| \geq \frac{m}{4} \right] \\ &\leq e^{-\Theta(m)} \leq \delta \end{aligned}$$

$$m = \Theta(\log^{1/\delta})$$

$\log(1/\delta)$
median $O\left(\frac{\log(1/\delta)}{\epsilon^2}\right)$ 

Median - of - Means