

Εκφράσεις

Διακρίσεων

Κατανομών

Έχουμε n στοιχεία

$$[n] = \{1, 2, \dots, n\}$$

Μια κατανομή $p \in [0, 1]^n$ με $\sum_{i=1}^n p_i = 1$

Ο χώρος όλων των κατανομών

$$D = \left\{ q \in [0, 1]^n : \sum_{i=1}^n q_i = 1 \right\}$$

Total Variation Distance

(ή ℓ_1 - distance)

$$d_{TV}(p, q) = \frac{1}{2} \sum_{i=1}^n |p_i - q_i| = \frac{1}{2} \|p - q\|_1$$

$$n=2$$

Κορώνα 1
Ροδάκισσα 2

$$d_{TV} \left((p, 1-p), (q, 1-q) \right)$$

$$= \frac{1}{2} (|p-q| + |1-p-(1-q)|)$$

$$= |p-q|$$

Τώρα δειγμάτωμα χρειαζόμαστε από μια κατανομή p για να μάθουμε \hat{p}

με $d_{TV}(\hat{p}, p) \leq \varepsilon$ με πιθανότητα $1-\delta$

Εκτιμητής: M με m δείγματα X_1, \dots, X_m

$$\hat{p}_i = \frac{\sum_{j=1}^m \mathbb{1}\{X_j = i\}}{m}$$

$$E[d_{TV}(\hat{p}, p)] = ?$$

$$E[\hat{p}_i] = E\left[\frac{\sum_{j=1}^m 1\{x_j=i\}}{m}\right] = \frac{1}{m} m p_i = p_i$$

$$\text{Var}[\hat{p}_i] = \frac{1}{m^2} m \text{Var}[1\{x_j=i\}]$$

$$= \frac{1}{m} p_i (1 - p_i) \leq \frac{1}{m} p_i$$

$$E[|\hat{p}_i - p_i|] \leq \sqrt{E[(\hat{p}_i - p_i)^2]} = \sqrt{\text{Var}[\hat{p}_i]} \leq \sqrt{\frac{p_i}{m}}$$

$$E[d_{TV}(\hat{p}, p)]$$

$$= \frac{1}{2} \sum_{i=1}^n E[|\hat{p}_i - p_i|]$$

$$\leq \frac{\sum_{i=1}^n \sqrt{p_i}}{2\sqrt{m}} \leq \frac{\sqrt{n}}{2\sqrt{m}}$$

Av $m = 4 \frac{n}{\epsilon^2}$ τότε

$$E[d_{TV}(\hat{p}, p)] \leq \frac{\epsilon}{4}$$

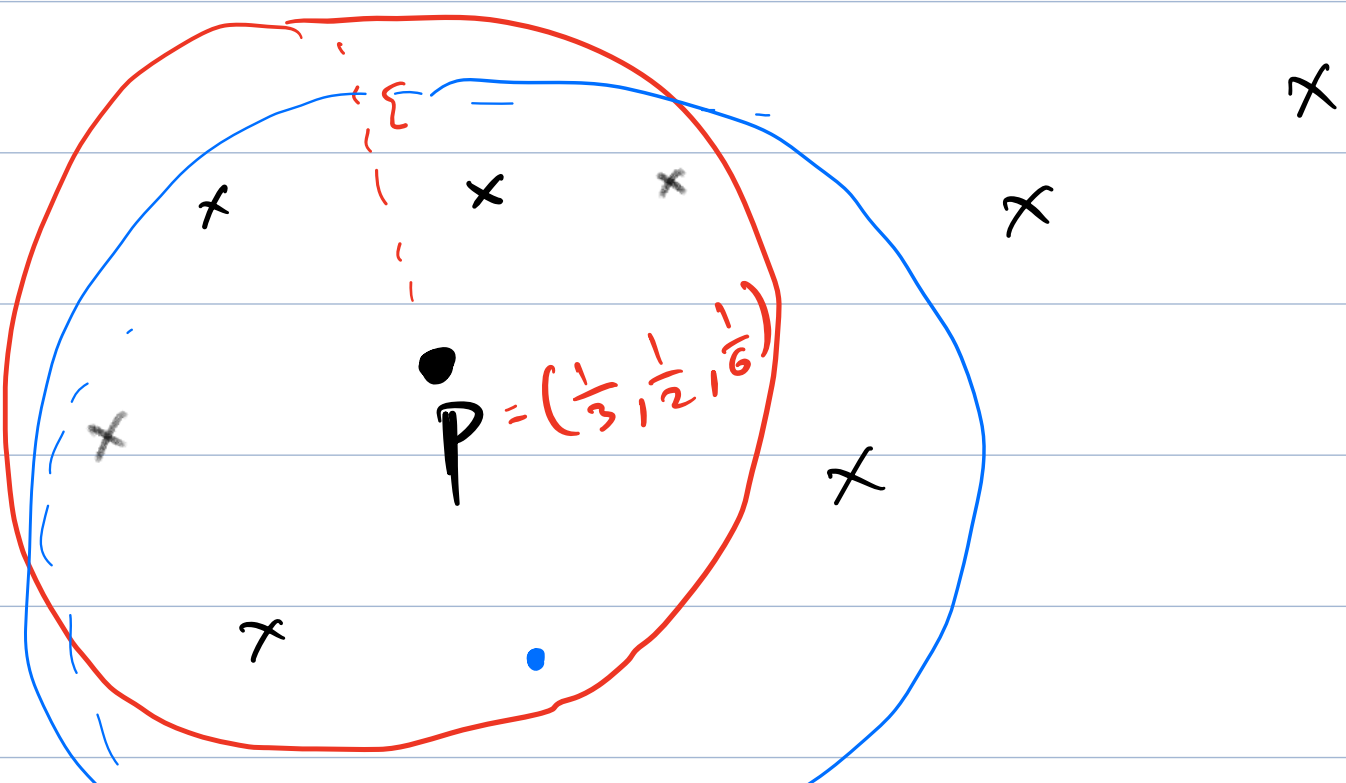
άρα

$$\Pr[d_{TV}(\hat{p}, p) > \epsilon] \leq \frac{1}{4}$$

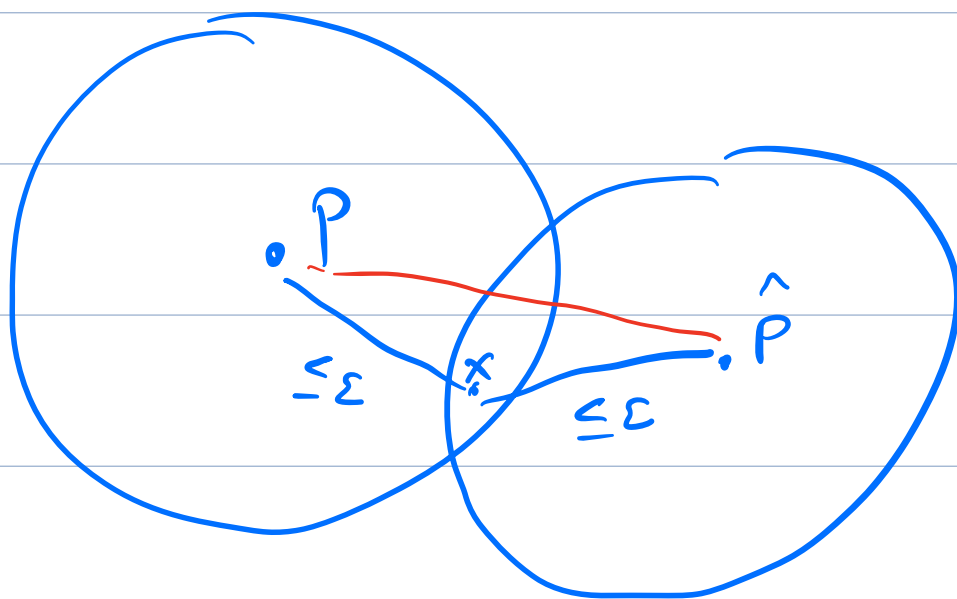
από Markov

Συνολικά με $m = \Theta\left(\frac{n}{\epsilon^2}\right)$ δείγματα
 βρίσκουμε \hat{p} τέτοιο ώστε
 $d_{TV}(\hat{p}, p) \leq \epsilon$ με πιθαν. $\frac{3}{4}$

Πως μπορούμε να αυξήσουμε την
 πιθανότητα;



Θ α βγαίνω \hat{p} η οποία
 είναι πιο κοντά στις περιπτώσεις

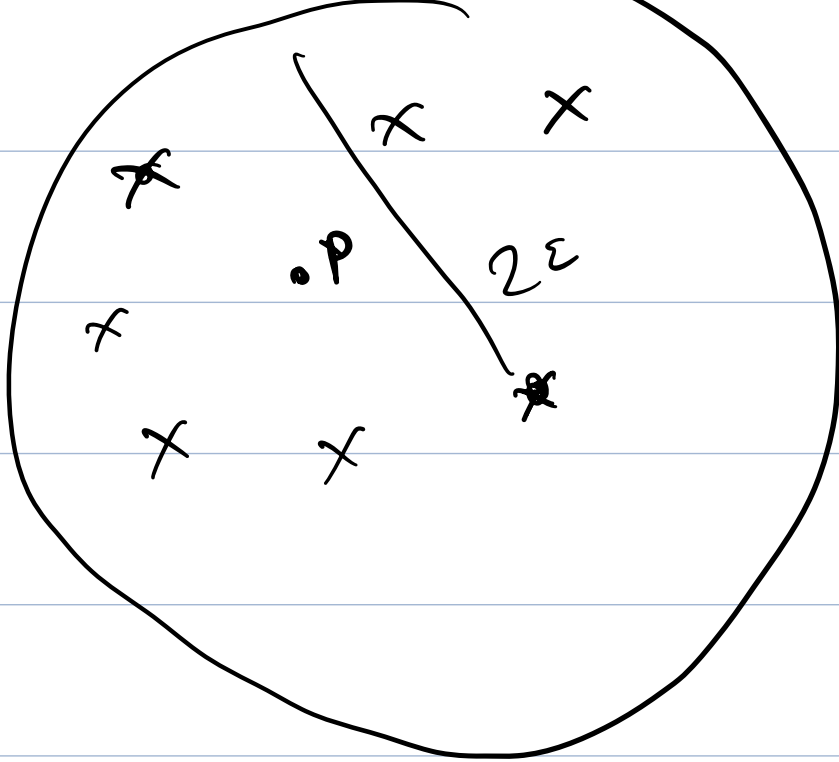


με $\Theta(\log(1/\delta))$ επαναλήψεις

με πιθανότητα δ πετυχαίνω

πάνω από τα μισά:

$$d_{TV}(\hat{p}, p) \leq 2\varepsilon$$



$$\Theta \left(\frac{n + \log(1/\epsilon)}{\epsilon^2} \right) \quad \text{δείγματα}$$

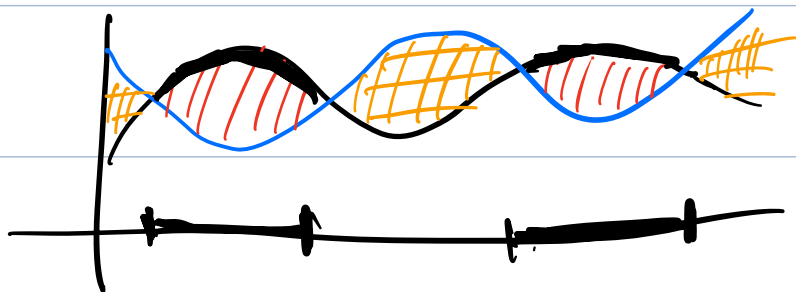
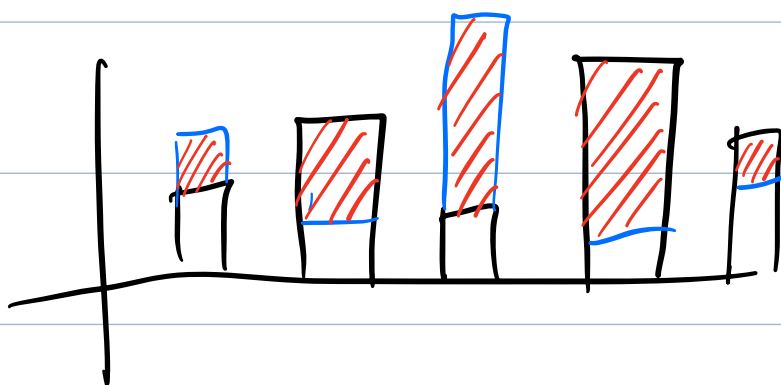
αρκούν για να μάθω π

TV - distance ϵ με πιθανότητα $1-\delta$.

Μετρικές στο χώρο των κατανομήων

- TV - distance p, q

$$d_{TV}(p, q) = \frac{1}{2} \sum |p_i - q_i| = \sum_{i=1}^n (p_i - q_i)_+ \quad \leftarrow$$
$$= \sum_{i=1}^n (q_i - p_i)_+$$



$$d_{TV}(p, q) = \max_{S \subseteq \mathcal{X}} |p(S) - q(S)|$$

$$\left| \Pr_{i \sim p} [i \in S] - \Pr_{i \sim q} [i \in S] \right|$$

Για να μάθω \hat{p} τ.ω.
 $d_{TV}(\hat{p}, p) \leq \epsilon$ με πιθανότητα $1 - \delta$
 χρειαζόμαστε $O\left(\frac{n + \log(1/\delta)}{\epsilon^2}\right)$ δείγματα

$$\hat{p}_i = \frac{\sum_{j=1}^m \mathbb{1}_{\{x_j = i\}}}{m}$$

ϕ_i άρνη ένα $S \subseteq [n]$

$$\Pr \left[\left| \hat{p}(S) - p(S) \right| \geq \epsilon \right] \leq e^{-\Omega(\epsilon^2 m)}$$

$$\text{Goal} \quad \forall \epsilon > 0 \quad \exists n \quad \text{d.t.v.} \quad (P, \hat{P}) \leq \epsilon$$

$$\text{Principle} \quad \forall S \subseteq [n]$$

$$|P(S) - \hat{P}(S)| \leq \epsilon$$

$$\text{Pr} [\exists S \subseteq [n] : |P(S) - \hat{P}(S)| > \epsilon]$$

$$\leq \sum_{S \subseteq [n]} \text{Pr} [|P(S) - \hat{P}(S)| > \epsilon]$$

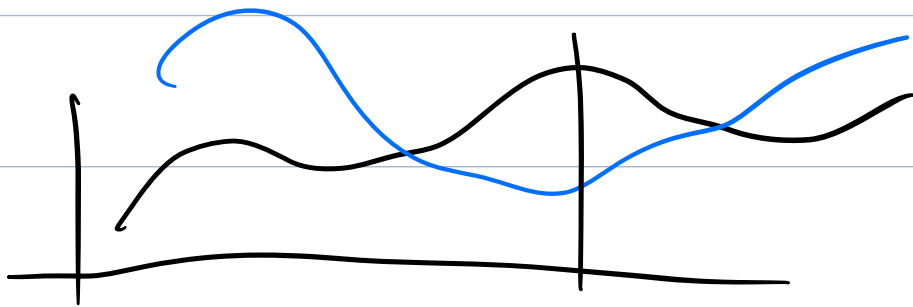
Union bound

$$\leq 2^n e^{-\Omega(\epsilon^2 m)}$$

$$\leq \delta$$

$$\text{As a} \quad m = O\left(\frac{n + \log(1/\delta)}{\epsilon^2}\right)$$

- Kolmogorov Distance



$$d_K(p, q) = \max_t \left| \Pr_{x \sim p}[x \leq t] - \Pr_{x \sim q}[x \leq t] \right|$$

- KL-divergence

$$D_{KL}(p \parallel q) = \sum_i p_i \log \frac{p_i}{q_i}$$

$$D_{KL}(p \parallel q) \geq 0$$
$$= \sum p_i \log p_i - \sum p_i \log q_i$$
$$= C - E_{i \sim p} [\log(q_i)]$$

Pinsker's Inequality

$$d_{TV}(p, q) \leq \sqrt{\frac{1}{2} D_{KL}(p \parallel q)}$$

Εκτίμηση ημιδιατάξιμου εκφρατός

Lower Bound

Χρειάζεται $\frac{1}{\epsilon^2}$;

Έστω $p = \text{Ber}(\frac{1}{2})$

$q = \text{Ber}(\frac{1}{2} + \epsilon)$

Θ.δ.ο. για ημιδιατάξιμα n φορές

$\frac{1}{4}$ χρειάζονται $\Omega(\frac{1}{\epsilon^2})$

δείγματα

Le Cam's method

Έστω 2 κατανομές p, q
σε ένα χώρο X

Στόχος: Να βρούμε test

$$T: X \rightarrow \{p, q\}$$

with το error probability

$$\max_x \left\{ \Pr_{x \sim p} [T(x) = q], \Pr_{x \sim q} [T(x) = p] \right\}$$

$$\left[\text{Lemma: Error Prob} \geq \frac{1}{2} (1 - d_{TV}(p, q)) \right]$$

Στατιστική ανάλυση

$$P = \text{Ber}\left(\frac{1}{2}\right)^{\otimes m}$$

$$Q = \text{Ber}\left(\frac{1}{2} + \varepsilon\right)^{\otimes m}$$

$$\text{Prob of Error} \geq \frac{1}{2} \left(1 - d_{TV}(P, Q)\right)$$

$$\left[\text{Claim: } d_{TV}(P, Q) \leq O(\sqrt{m} \varepsilon) \right]$$

$$\geq \frac{1}{2} \left(1 - O(\sqrt{m} \varepsilon)\right)$$

$$\geq \frac{1}{4}$$

εκτός αν

$$m = \Omega\left(\frac{1}{\varepsilon^2}\right)$$

$$[\text{Claim : } d_{TV}(P, Q) \leq O(\sqrt{m} \varepsilon)]$$

$$d_{KL}(P, Q) = \sum_{x \in \{0,1\}^m} P_x \log \frac{P_x}{Q_x}$$

$$= \sum_{x \in \{0,1\}^m} \prod_{i=1}^m p(x_i) \log \frac{\prod_{i=1}^m p(x_i)}{\prod_{i=1}^m q(x_i)}$$

$$= \sum_{i=1}^m \underbrace{\sum_{x_i \in \{0,1\}^{m-1}} \prod_{j \neq i} p(x_j)}_1 \underbrace{\sum_{x_i \in \{0,1\}} p(x_i) \log \frac{p(x_i)}{q(x_i)}}_{d_{KL}(p, q)}$$

$$= m d_{KL}(p, q)$$

$$= m d_{KL}(\text{Ber}(\frac{1}{2}), \text{Ber}(\frac{1}{2} + \varepsilon))$$

$$d_{KL}(\text{Ber}(\frac{1}{2}), \text{Ber}(\frac{1}{2} + \varepsilon))$$

$$= \frac{1}{2} \log \frac{\frac{1}{2}}{\frac{1}{2} + \varepsilon} + \frac{1}{2} \log \frac{\frac{1}{2}}{\frac{1}{2} - \varepsilon}$$

$$= \frac{1}{2} \log \frac{1}{1 + 2\varepsilon} + \frac{1}{2} \log \frac{1}{1 - 2\varepsilon}$$

$$= \frac{1}{2} \log \frac{1}{(1 + 2\varepsilon)(1 - 2\varepsilon)} = \frac{1}{2} \log \frac{1}{1 - 4\varepsilon^2}$$

$$= -\frac{1}{2} \log (1 - 4\varepsilon^2)$$

$$= -\frac{1}{2} (-4\varepsilon^2 + \mathcal{O}(\varepsilon^4))$$

$$= 2\varepsilon^2 + \mathcal{O}(\varepsilon^4) = \mathcal{O}(\varepsilon^2)$$

$$d_{KL}(P, Q) = O(m \varepsilon^2)$$

$$d_{TV}(P, Q) = O(\sqrt{m} \varepsilon)$$

Ανάλυση του Λήμματος του Le Cam

Error Prob =

$$\max \left\{ \Pr_{x \sim P} [T(x) = Q], \Pr_{x \sim Q} [T(x) = P] \right\}$$

$$\geq \frac{1}{2} \left(\Pr_{x \sim P} [T(x) = Q] + \Pr_{x \sim Q} [T(x) = P] \right)$$

$$= \frac{1}{2} \left(\sum_x P(x) \mathbb{1}[T(x)=Q] + \sum_x Q(x) \mathbb{1}[T(x)=P] \right)$$

$$= \frac{1}{2} \left(\sum_x P(x) (1 - \mathbb{1}[T(x)=P]) + \sum_x Q(x) \mathbb{1}[T(x)=P] \right)$$

$$= \frac{1}{2} \left(1 - \sum_x (P(x) - Q(x)) \mathbb{1}[T(x)=P] \right)$$

$$\geq \frac{1}{2} \left(1 - \sum_x (P(x) - Q(x))_+ \right)$$

$$= \frac{1}{2} (1 - d_{TV}(P, Q)) \quad \square$$