

# Distribution Testing

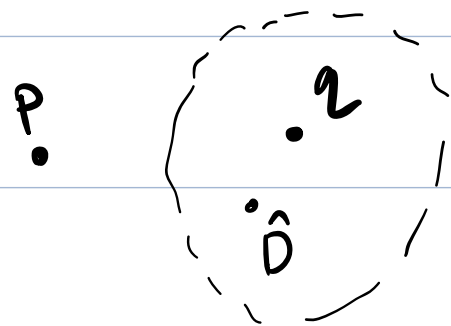
Δίνονται  $\Sigma$  κατανομές  $p, q$

και sample access από μια άγνωστη  $D$

Ξέρουμε ότι είτε  $D = p$

είτε  $D = q$

Πόσα δείγματα χρειαζόμαστε για να διακρίνουμε  
αυτές τις περιπτώσεις;



- Approach 1

Μαθαίνω  $\hat{D}$  τω  $d_{TV}(D, \hat{D}) < \frac{\epsilon}{2}$

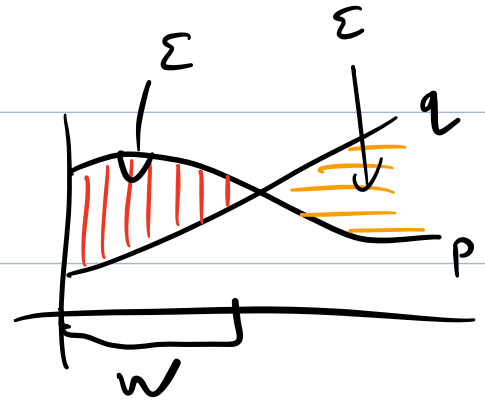
με πιθανότητα  $1 - \delta$

Χρειαζόμαστε  $m = \Omega\left(\frac{n + \log(1/\delta)}{\epsilon^2}\right)$

- Approach 2

$$d_{TV}(p, q) = \varepsilon$$

||



$$\sum_i (p_i - q_i)_+$$

$$W = \{ i \in [n] : p_i > q_i \}$$

$$p(w) = q(w) + \varepsilon$$

$$D(w) = p(w)$$

∩

$$D(w) = q(w)$$

Μετρών το ποσοστό δειγμάτων  
της  $D$  που έπεσαν στο  $w$

$$|\hat{d} - D(w)| < \frac{\varepsilon}{2}$$

Χρειάζονται μόνο  $O(\frac{1}{\varepsilon^2})$   
δειγμάτα

$$\text{Αν } p = D$$

$$|\hat{d} - p(w)| < \frac{\varepsilon}{2}$$

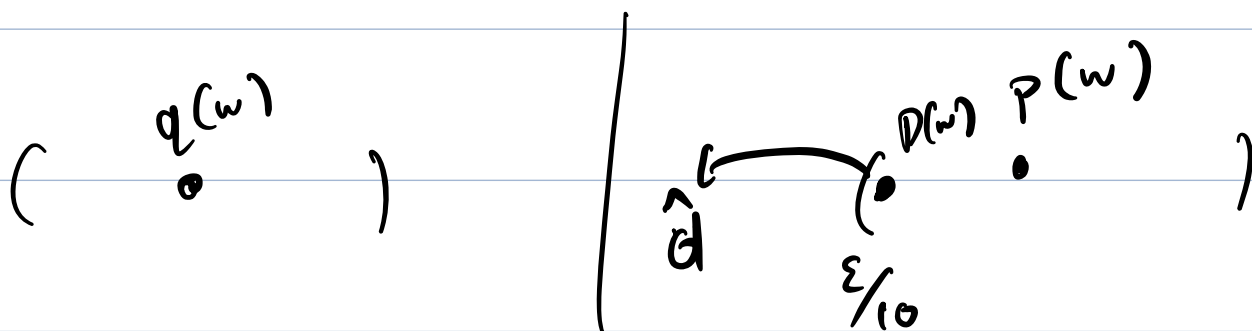
$$\begin{aligned} |\hat{d} - q(w)| &= |\hat{d} - p(w) - \varepsilon| \\ &\leq \varepsilon - |\hat{d} - p(w)| \\ &> \varepsilon/2 \end{aligned}$$

Έστω ότι  $\eta \in D$   
 είναι μόνο  $\frac{\varepsilon}{10}$   
 κοντά στην  $p$  ή  $q$   
 δλδ

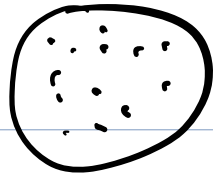
$$d_{TV}(D, p) < \frac{\varepsilon}{10} \quad \text{ή} \quad d_{TV}(D, q) < \frac{\varepsilon}{10}$$

$$d_{TV}(p, q) = \varepsilon$$

$$|p(w) - D(w)| < \frac{\varepsilon}{10}$$



# Open Testing σε Learning



Οικογένεια κατανομών  $\mathcal{D}$

$\epsilon$ -cover : Για μια οικογένεια  
κατανομών  $\mathcal{D}$  ένα  
υποσύνολο  $\mathcal{D}_\epsilon \subseteq \mathcal{D}$

λέγεται  $\epsilon$ -cover αν για  
κάθε  $p \in \mathcal{D} \exists p' \in \mathcal{D}_\epsilon$   
ώστε  $d_{TV}(p, p') \leq \epsilon$

$$\pi.X. \quad D = \{ \text{Ber}(p) : p \in [0, 1] \}$$

$$D_\varepsilon = \{ \text{Ber}(p) : p \in \{ \varepsilon, 2\varepsilon, 3\varepsilon, \dots \} \}$$

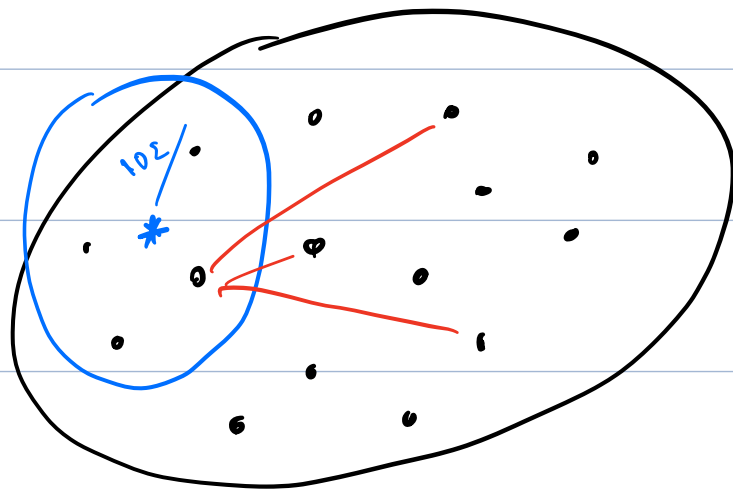
$$|D_\varepsilon| = \frac{1}{\varepsilon}$$

$$\pi.X. \quad D = \{ p \in [0, 1]^n : \sum_{i=1}^n p_i = 1 \}$$

$$D_\varepsilon = \{ p \in \{ \frac{\varepsilon}{n}, 2\frac{\varepsilon}{n}, 3\frac{\varepsilon}{n}, \dots \}^n : \sum p_i = 1 \}$$

$$|D_\varepsilon| \leq \left( \frac{n}{\varepsilon} \right)^n \cdot \underbrace{1 \cdot 1 \cdot 1 \cdot 1 \cdot \dots}_{n \text{ times}}$$

$$\left( \frac{n}{\varepsilon} \right)^n \approx \left( \frac{1}{\varepsilon} \right)^n$$



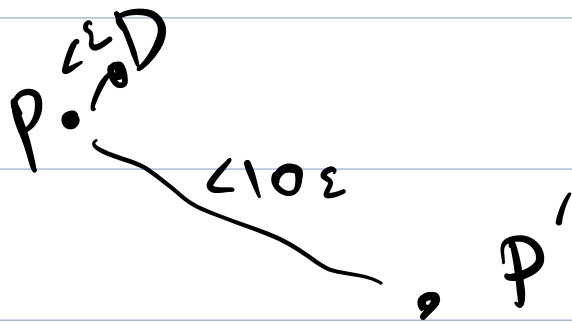
Tournament :

Τρίγων του testing αλγόριθμο  
για κάθε ζεύγος κατανομιών

- Νικητής είναι όποιος κερδίζει  
όλες τις κατανομές σε απόσταση 10 ε

Έστω

$$p \text{ που έχει}$$
$$d_{TV}(p, D) < \varepsilon$$



Αν δαχίσω έναν τέτοιο νικητή  $p'$

$$d_{TV}(p', D) < 11\varepsilon$$

Η πιθανότητα σφάλματος  $\leq |D_\varepsilon|^2 \delta$

Αν δώσω συνολικά  $\leq \delta'$

$$\text{δίνω} \quad \delta = \frac{\delta'}{|D_\varepsilon|^2}$$



$$|D_{\varepsilon}|$$

Συνολικά το πλήθος των  
δειγμάτων που χρειάζονται είναι

$$\Theta\left(\frac{1}{\varepsilon^2} \log(1/\delta)\right) \text{ δηλαδή}$$

$$\Theta\left(\frac{1}{\varepsilon^2} \log(|D_{\varepsilon}|^2 / \delta')\right)$$

$$\Theta\left(\frac{\log |D_{\varepsilon}|}{\varepsilon^2} + \frac{\log(1/\delta')}{\varepsilon^2}\right)$$

Για διακριτές κατανομές σε  $n$  στοιχεία

$$|D_\varepsilon| \leq \left(\frac{1}{\varepsilon}\right)^n$$

$$\Theta\left(\frac{\log|D_\varepsilon|}{\varepsilon^2} + \frac{\log(1/\delta)}{\varepsilon^2}\right)$$

=

$$\Theta\left(\frac{n \log(1/\varepsilon)}{\varepsilon^2} + \frac{\log(1/\delta)}{\varepsilon^2}\right)$$

=

$$\tilde{\Theta}\left(\frac{n}{\varepsilon^2} + \frac{\log(1/\delta)}{\varepsilon^2}\right)$$

$$\tilde{\Theta}(f(n, m)) \equiv \Theta(f(n, m) \log^{o(1)} f(n, m))$$

# Distribution Testing

Έστω ότι μου δίνεται μια  $q$   
και έχω δείγματα από μια  $p$

Πρέπει να διακρίνω

- $p = q$

- $d_{TV}(p, q) > \epsilon$

Γενικά χρειαζόμαστε  $O\left(\frac{\sqrt{n}}{\epsilon^2}\right)$

δείγματα για να διακρίνουμε το πρόβλημα

# Uniformity Testing

$$p = \left(\frac{1}{n}, \dots, \frac{1}{n}\right)$$

$$d_{TV}(p, (\frac{1}{n}, \dots, \frac{1}{n})) > \epsilon$$

Γενικότατα

$$\frac{\|p\|_{2/3}}{\epsilon^2} = \left( \sum p_i^{2/3} \right)^{3/2} \leq \sqrt{n}$$

δειχνάτα για γενικό p

### Αλγόριθμος

Με  $m$  δείγματα από την  $p$

$C \leftarrow \# \text{ collisions}$

οπότε  $C = \sum_{i < j} \mathbb{1}\{x_i = x_j\}$

Αν  $C < \frac{1 + \epsilon^2/2}{n} \binom{m}{2}$

τότε "Uniform"

Αλλάζω "OXI Uniform"

- Αναμένω  $n_j$  collisions

$$E[C] = \sum_{i < j} \Pr[x_i = x_j]$$

$$= \sum_{i < j} \sum_{k=1}^n \Pr[x_i = k \wedge x_j = k]$$

$$= \sum_{i < j} \sum_{k=1}^n p_k^2$$

$$= \sum_{k=1}^n p_k^2 \cdot \binom{n}{2}$$

$$= \|p\|_2^2 \binom{n}{2}$$

$$\|p\|_2^2 \geq \frac{\|p\|_1^2}{n} = \frac{1}{n}$$

( =  $\sigma_{\text{TV}}$  uniform )

Εστω  $\delta > 0$  το  $p$  έχει

$$d_{\text{TV}}(p, \mathcal{U}_n) > \varepsilon$$

||

$$\|p - \mathcal{U}_n\|_1 > \varepsilon$$

$\Leftrightarrow$

$$\|p - \mathcal{U}_n\|_2^2 > \frac{\varepsilon^2}{n}$$

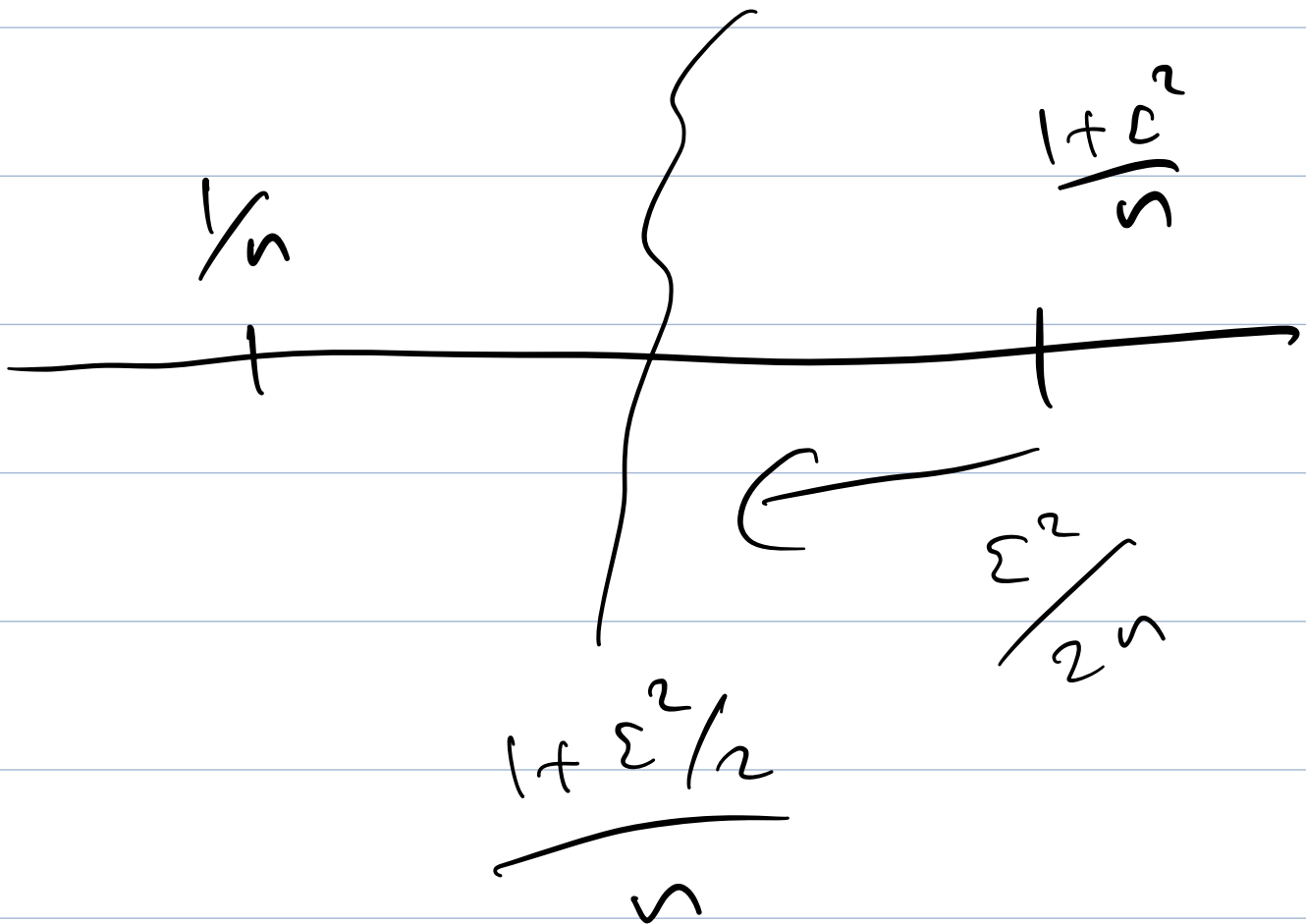
$$\|p\|_2^2 = \|\mathcal{U}_n\|_2^2 + \|p - \mathcal{U}_n\|_2^2$$

$$\geq \frac{1}{n} + \frac{\varepsilon^2}{n}$$

$P_{\text{uniform}}$   $d_{TV}(P, U_n) > \epsilon$

$$E[C] > \binom{m}{2} \frac{1 + \epsilon^2}{n}$$

$$E[C] = \binom{m}{2} \frac{1}{n}$$



Ano Chebyshev

$$\Pr [ |C - E[C]| \geq 2 \sqrt{\text{Var}[C]} ] \leq \frac{1}{4}$$

$$\sqrt{\text{Var} \left[ \frac{C}{\binom{m}{2}} \right]} \leq \frac{\varepsilon^2}{4n}$$

$$\text{Var}[C] \leq 2 m^3 \|p\|_2^3$$

$$\text{Var} \left[ \frac{C}{\binom{m}{2}} \right] \leq O \left( \frac{\|p\|_2^3}{m} \right)$$

$$\|p\|_2^3 < \frac{\varepsilon^4}{m}$$



$$\frac{1}{m} \quad \frac{1}{16n^2}$$

$$m \gg 0 \left( \frac{n^2 \|p\|_2^3}{\varepsilon^4} \right)$$

$$12 \gg 0 \left( \frac{\sqrt{n}}{\varepsilon^4} \right)$$