

Today:

- wrap up learning arbitrary discrete distributions (see previous notes)
 - learning of monotone discrete distributions
-

Monotone distributions:

D on $[n]$ where $p_i =$ probability of i
is monotone if $p_1 \geq p_2 \geq p_3 \geq \dots \geq p_n$

How many samples from D are needed
to output D' such that $d_{TV}(D, D') \leq \epsilon$
with probability $4/5$?

We will show $O(\epsilon^{-3} \log n)$

Partition $[n]$ into buckets $\{a_i, a_{i+1}, \dots, b_i\}$
 $\{1, 2, 3, \dots, n\}$

$$a_1 = b_1 = 1$$

then

$$a_{i+1} = b_i + 1$$

$$b_{i+1} = \lceil b_i(1 + \varepsilon) \rceil$$

(truncate when n is reached)

Algorithm:

- collect $t = C \cdot \varepsilon^{-3} \log n$ independent samples from D

- output D' such that probability of $i \in \{a_j, \dots, b_j\}$ equals $(\# \text{samples in } \{a_j, \dots, b_j\}) / (b_j - a_j + 1)$

Intuition: replace each bucket with the uniform distribution on the bucket of the same total probability

Claim: monotone distributions close to distributions uniform

on each bucket $\{a_i, a_i+1, a_i+2, \dots, b_i\}$

D - monotone distribution on $[n]$

p_i = probability of i in D

D_* - distribution created by making each bucket of D uniform

q_i = probability of i in D_*

$$\text{For } i \in \{a_j, \dots, b_j\}, q_i = \frac{\sum_{k=a_j}^{b_j} p_k}{a_j - b_j + 1}$$

$$d_{TV}(D, D_*) \leq \epsilon$$

Homework 2

Proof: Let's bound $\|p - q\|_1 \stackrel{\leftarrow}{=} 2 d_{TV}(D, D_*)$

where

$$p = (p_1, \dots, p_n) \text{ and } q = (q_1, \dots, q_n)$$

Consider bucket j : $\{a_j, a_j+1, \dots, b_j\}$

If $a_j = b_j$: D equals D_* on this bucket

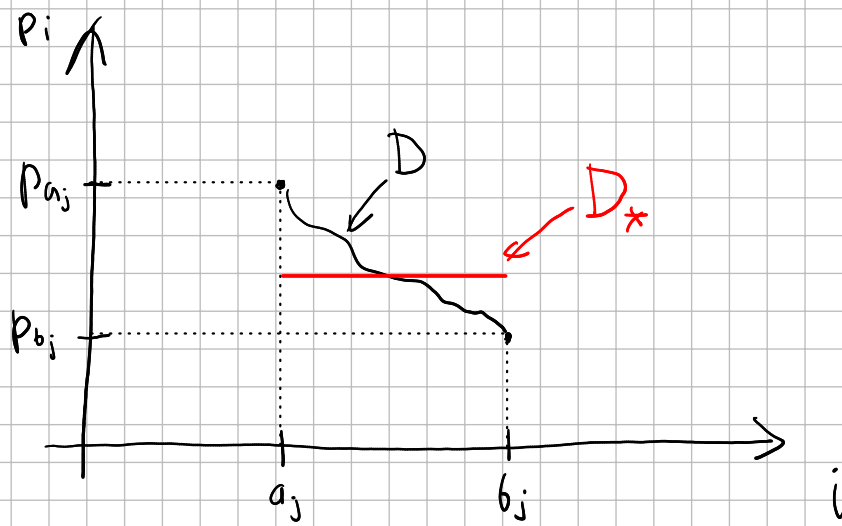
If $a_j < b_j$:

$$b_j = \lceil b_{j-1} (1 + \epsilon) \rceil > a_j = b_{j-1} + 1$$

\Downarrow

$$b_{j-1} \epsilon > 1$$

$\boxed{16-3}$



$$\begin{aligned}
 \sum_{k=a_j}^{b_j} |p_k - q_k| &\leq (b_j - a_j + 1) \cdot (p_{a_j} - p_{b_j}) \\
 &\leq \lceil \varepsilon b_{j-1} \rceil \cdot (p_{a_j} - p_{b_j}) \\
 &\leq (\varepsilon b_{j-1} + 1) (p_{a_j} - p_{b_j}) \\
 &\leq 2\varepsilon b_{j-1} (p_{a_j} - p_{b_j})
 \end{aligned}$$

Let $B =$ number of buckets

$$\|p - q\|_1 \leq \sum_{j=2}^B 2\varepsilon b_{j-1} (p_{a_j} - p_{b_j})$$

$$= 2\varepsilon \sum_{j=2}^B \left((p_{a_j} - p_{b_j}) \sum_{k=1}^{j-1} (b_k - a_k + 1) \right)$$

$$= 2\varepsilon \sum_{k=1}^{B-1} \left((b_k - a_k + 1) \sum_{j=k+1}^B (p_{a_j} - p_{b_j}) \right)$$

$$\leq 2\varepsilon \sum_{k=1}^{B-1} (b_k - a_k + 1) p_{a_{k+1}}$$

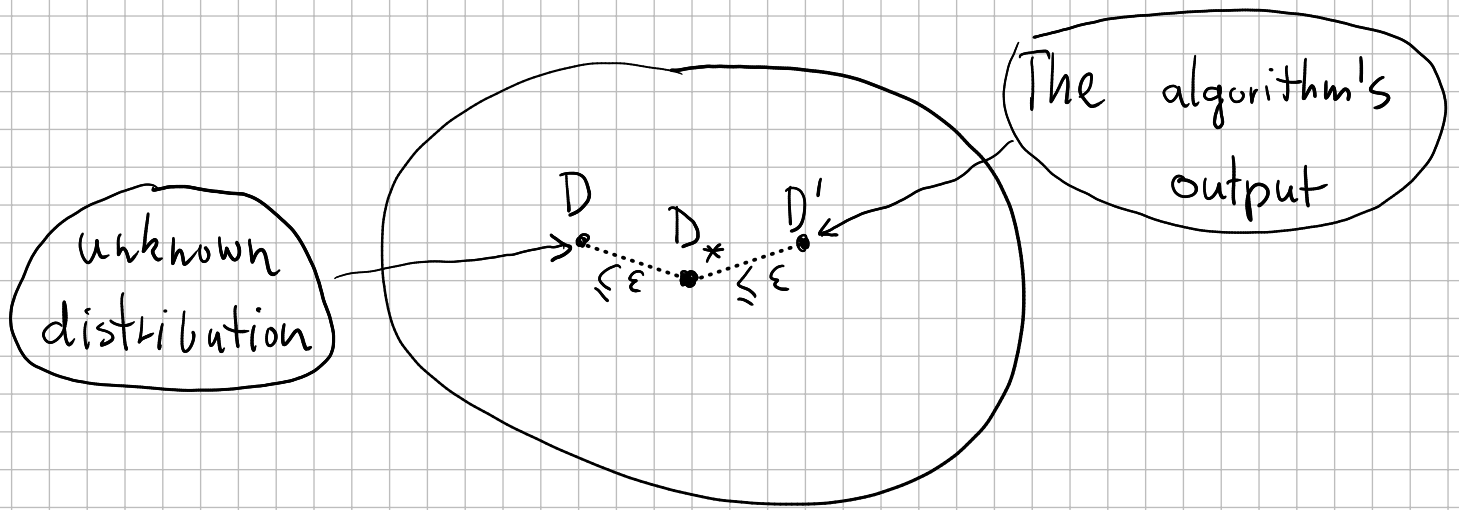
$$\leq 2\varepsilon \sum_{k=1}^{B-1} \sum_{i=a_k}^{b_k} p_i \leq 2\varepsilon \sum_{i=1}^n p_i = 2\varepsilon$$

116-4

This implies $d_{TV}(D, D_*) \leq \epsilon$ (via Homework 2)



Why the algorithm works:



$$d_{TV}(D, D') \leq 2\epsilon \text{ with probability } 99/100$$

because

$$d_{TV}(D_*, D') \leq \text{with probability } 99/100$$

Why this holds?

- D_* and D' are both uniform on each bucket

- L_1 -distance between D_* and D'

$$= \sum_{i=1}^B |(\text{probability of bucket } i \text{ in } D_*)$$

$$- (\text{probability of bucket } i \text{ in } D')|$$

- hence it suffices to estimate the distribution

[16-5] on $B = O\left(\frac{1}{\epsilon} \log n\right)$ buckets up to $d_{TV}(\dots) \leq \epsilon$

Sketch of the general idea

- Sample $O(1/\epsilon)$ elements
- Divide $[n]$ into $O(1/\epsilon)$ ranges, each containing $O(1)$ samples
- The ranges correspond roughly to $\Theta(\epsilon)$ mass of the distribution
- One of the ranges contains the peak
- For each range as the "peak candidate", create $O(\frac{1}{\epsilon} \log n)$ buckets for ranges before and after as in learning monotone distributions (one increasing, one decreasing)
- Total of $O(\epsilon^{-2} \log n)$ intersecting buckets
- Take their intersections to create $O(\epsilon^{-2} \log n)$ non-intersecting buckets

Example:



- Learn the distribution on the resulting $O(\varepsilon^{-2} \log n)$ buckets up to ε in total variation distance with $O\left(\frac{O(\varepsilon^{-2} \log n)}{\varepsilon^2}\right) = O(\varepsilon^{-4} \log n)$ samples
- Output distribution uniform on each bucket with the learned probabilities of buckets