# Lecture 13:

## Testing Distributions

- Uniformity (cont.)

- Monotonicity

Turning to a new model :

Probability distributions - get samples of distribution



Domain $D$, $|D| = n$ ← known

$P_i = \Pr[p \text{ outputs } i]$ ←unknown

outputs iid samples

↖ this is all we can learn from

Examples:

Lottery data

Shopping choices

experimental outcomes

⋮

What do we want to know?

is it uniform? eg. lottery
is it high entropy?
large support? (many distinct elements have >0 probability
is p monotone increasing, K-modal, monotone hazard rate...?

how can we do it?

$\chi^2$ test

plug in estimate

learn distribution, Maximum likelihood estimates

Goal: sample complexity SUBLINEAR in $n$

## Testing Uniformity

The goal:

↙ Uniform dist on D

- if $P \equiv U_D$ then tester outputs PASS ← with prob $\geq 3/4$

- if $dist(P, U_D) > \varepsilon$ then tester outputs FAIL

which measure of distance?

$\ell_1, \ell_2$, KL-divergence, Earth mover, Jensen-Shannon

↑ today's focus

## Distances

$\ell_1$-distance : $\|p-q\|_1 = \sum_{i \in D} |p_i - q_i|$

$\ell_2$-distance : $\|p-q\|_2 = \sqrt{\sum_{i \in b} (p_i - q_i)^2}$

$$\|p-q\|_2 \leq \|p-q\|_1 \leq n^{1/2} \|p-q\|_2$$

examples:

① $p = (1, 0, 0, \dots 0)$



$q = (\frac{1}{n}, \frac{1}{n}, \dots \frac{1}{n})$



$\ell_1$ distance:

$\|p-q\|_1 = (\frac{n-1}{n}) + (n-1) \cdot \frac{1}{n}$

$\approx 2$

$\ell_2$-distance:

$\|p-q\|_2^2 = (1 - \frac{1}{n})^2 + (n-1)(\frac{1}{n})^2$

$\approx 1$

② $p = (\frac{2}{n}, \frac{2}{n}, \dots \frac{2}{n}, 0, 0, \dots 0)$

$q = (0, 0, \dots 0, \frac{2}{n}, \frac{2}{n}, \dots \frac{2}{n})$





$\ell_1$ distance:

$\|p-q\|_1 = n \cdot (\frac{2}{n}) = 2$

$\ell_2$-distance: $\|p-q\|_2^2 = n \cdot (\frac{2}{n})^2 = \frac{4}{n}$

$\|p-q\|_2 = \frac{2}{\sqrt{n}}$

## "Plug-in" Estimate :

Algorithm :

- take $m$ samples from $p$

- estimate $p(x)$ $\forall x$ via

$$\hat{p}(x) = \frac{\# \text{ times } x \text{ occurs in sample}}{m}$$

- if $\sum_x \left| \hat{p}(x) - \frac{1}{n} \right| > \varepsilon$   reject

   else   accept.

Analysis : ( better analyses exist )

pick $m$ s.t. $\forall x, \left| \hat{p}(x) - p(x) \right| < \frac{\varepsilon}{n}$ $\Rightarrow$ $\| \hat{p} - p \|_1 < \varepsilon$

$\boxed{\text{So, if } p = U_n \text{ then } p \text{ passes}}$

by $\Delta \ddagger$, if $\| p - \hat{p} \|_1 < \varepsilon$  $+$  $\| \hat{p} - U \|_1 < \varepsilon$

then   $\| p - U \|_1 < 2\varepsilon$.

$\boxed{\text{So, if } \| p - U_n \|_1 > 2\varepsilon \text{ this test is likely to fail}}$

how many samples? $\Omega\left(\frac{n}{\varepsilon}\right)$ maybe even worse ...

$\boxed{\text{for each } x, \text{ need to see it at least once in order to give non zero estimate.}}$

$\Theta(n)$? Can we do better?

Better analysis:

Claim $E[\|\hat{p}-p\|_{1,\|}] \leq \sqrt{\frac{n}{m}}$

Pf

$E[\|\hat{p}-p\|_1] = \sum_x E[|\hat{p}(x)-p(x)|]$ ⟵ note: $E[\hat{p}(x)] = \frac{1}{m}E\left[\sum 1_{i^{th}\text{ sample is } x}\right]$

$\leq \sum_x \sqrt{E[(\hat{p}(x)-p(x))^2]}$ ⟵ (Jensen's ≠)

$= \frac{1}{m}\sum_{i=1}^{m} E[1_{i^{th}\text{ sample is } x}]$

$= \sum_x \sqrt{Var(\hat{p}(x))}$

$= \frac{m \cdot p(x)}{m} = p(x)$

$Var(\hat{p}(x)) = \frac{1}{m^2} m \, p(x)(1-p(x))$

$\leq \sum_x \sqrt{\frac{p(x)}{m}}$

$\leq \frac{p(x)}{m}$

$\leq \frac{1}{\sqrt{m}} \cdot \sqrt{n}$ ⟵ since $\max_{p \in \text{prob dist over domain of size } n} \sum \sqrt{p(x)}$ is $\sqrt{n}$

So picking $m = \Omega\left(\frac{n}{\varepsilon^2}\right)$ gives

$E[\|\hat{p}-p\|_1] \leq \frac{\varepsilon}{2}$

by Markov's ≠ : with prob $1-\frac{1}{2}$, $\|\hat{p}-p\|_1 \leq \varepsilon$

Note, this says can "learn" (approximate) any dist w.r.t. $L_1$ distance in $\theta(n/\varepsilon^2)$ samples

# $L_2$ - Distance (squared):

$$\|p - u_{[n]}\|_2^2 = \sum_{i \in [n]} (p_i - \tfrac{1}{n})^2$$

$$= \sum p_i^2 - \tfrac{2}{n} \underbrace{\sum p_i}_{=1} + \underbrace{\sum (\tfrac{1}{n})^2}_{=\tfrac{1}{n}}$$

$$= \underbrace{\sum p_i^2}_{} - \tfrac{1}{n}$$

Collision probability of $p$:

$$\|p\|_2^2 \equiv \Pr_{s, t \sim p} [ s = t ] = \sum p_i^2$$

$$\text{for} \quad p = u, \quad \|p\|_2^2 = \tfrac{1}{n}$$
$$\text{for} \quad p \neq u, \quad \|p\|_2^2 > \tfrac{1}{n}$$

$$= \underbrace{\|p\|_2^2}_{\substack{\text{we can} \\ \text{estimate} \\ \text{this}}} - \underbrace{\|u_{(n)}\|_2^2}_{\substack{\text{we know this} \\ \text{since we know } n}}$$

# Algorithm

1. take $s$ samples from $p$   ① how many samples?
2. let $\hat{c} \leftarrow$ estimate of $\|p\|_2^2$ from sample   ② how?
3. if $\hat{c} < \tfrac{1}{n} + \delta$ pass   ③ what should $\delta$ be?
   else fail

First:
How to estimate $\|p\|_2^2$ ?

Naive idea:

take two new samples:
$$X_i \leftarrow \begin{cases} 1 & \text{if samples are equal} \\ 0 & \text{o.w} \end{cases}$$

" gives $\theta(k)$ samples of collision probability
from $k$ samples of $p$ "

Better idea: recycle - use _all_ pairs in sample

" gives $\theta(k^2)$ samples of collision probability
from $k$ samples of $p$ "

Estimate by recycling:

- Take $s$ samples from $p$: $X_1 \cdots X_s$
- for each $1 \le i < j \le s$
$$\delta_{ij} \leftarrow \begin{cases} 1 & \text{if } X_i = X_j \\ 0 & \text{if } X_i \ne X_j \end{cases}$$
- Output $\hat{c} \leftarrow \dfrac{\sum_{i<j} \delta_{ij}}{\binom{s}{2}}$

$\delta_{ij}$'s not independent so can't use Chernoff

Analysis: $E[\hat{c}] = \dfrac{1}{\binom{s}{2}} \cdot \binom{s}{2} \cdot E[\delta_{ij}]$
$$= \|p\|_2^2$$

How well do we need to estimate $\|p\|_2^2$ ?

Assumption ✴: $\quad |\hat{c} - \|p\|_2^2| < \Delta$

↰ this is our parameter that determines whether our approximation is good. Spoiler: will set $\Delta = \frac{\varepsilon^2}{2}$

will take enough samples so that this holds with prob $\geq 3/4$

What happens if ✴ holds with $\Delta = \frac{\varepsilon^2}{2}$ ?

**Correct behavior!**

• if $p = U_{[n]}$ then $\hat{c} \leq \|U_{[n]}\|_2^2 + \Delta = \frac{1}{n} + \frac{\varepsilon^2}{2}$

so test will PASS

• if $\|p - U_{[n]}\|_2 > \varepsilon$ then $\|p - U_{[n]}\|_2^2 > \varepsilon^2$

but $\|p\|_2^2 = \|p - U_{[n]}\|_2^2 + \frac{1}{n}$ ← see p.6

$\qquad\qquad > \varepsilon^2 + \frac{1}{n}$

† $\qquad \hat{c} > \|p\|_2^2 - \Delta$ ← ✴

$\qquad\qquad \geq \varepsilon^2 + \frac{1}{n} - \Delta = \varepsilon^2 + \frac{1}{n} - \frac{\varepsilon^2}{2} = \frac{\varepsilon^2}{2} + \frac{1}{n}$

so test will FAIL

**Remaining Question:**

How many samples do we need to estimate $\hat{c}$ to within $\Delta$ ?

## Analysis

$$E[\mathbb{6}_{ij}] = Pr[\mathbb{6}_{ij} = 1]$$
$$= \|p\|_2^2$$

$$E[\hat{c}] = \frac{1}{\binom{s}{2}} \binom{s}{2} E[\mathbb{6}_{ij}] = \|p\|_2^2$$

$$Pr\left[\,|\hat{c} - \|p\|_2^2| > \rho\,\right] \le \frac{Var[\hat{c}]}{\rho^2}$$

*(red, right margin)* recall:
$$Var[X] = E[(X - E[X])^2]$$

*(magenta)* Chebyshev ≠

**Fact** $Var[aX] = a^2 Var[X]$

So $\quad Var[\hat{c}] = Var\left[\frac{1}{\binom{s}{2}} \cdot \sum_{i<j} \mathbb{6}_{ij}\right]$

$$= \frac{1}{\binom{s}{2}^2} Var\left[\sum_{i<j} \mathbb{6}_{ij}\right]$$

**Lemma** $\quad Var\left[\sum \mathbb{6}_{ij}\right] \le 4\left(\binom{s}{2}\|p\|_2^2\right)^{3/2}$

**Why?** (proof...)

$\quad$ **def.** $\quad \overline{\mathbb{6}}_{ij} = \mathbb{6}_{ij} - E[\mathbb{6}_{ij}]$

$\quad$ So $\quad E[\overline{\mathbb{6}}_{ij}] = 0$

$\quad$ Also ∴ $E[\overline{\mathbb{6}}_{ij}\overline{\mathbb{6}}_{kl}] \le E[\mathbb{6}_{ij}\mathbb{6}_{kl}]$

*(magenta, left brace) Verify at home? (or trust...)*
$$\left(\sum p(x)^3\right)^{1/3} \le \left(\sum p(x)^2\right)^{1/2}$$
$$s^2 \le 3\binom{s}{2}$$
$$\binom{s}{3} \le s^3/6$$

*(blue box, right)* Fact $\Rightarrow$
$$Var[\hat{c}] \le \frac{4 \cdot \left(\binom{s}{2}\|p\|_2^2\right)^{3/2}}{\binom{s}{2}^2}$$
$$\le \theta\left(\|p\|_2^3 / s\right)$$

*(purple box, right)* ← trick – will rewrite variance as $E[\overline{\mathbb{6}}_{ij}]^2$.
why?
$$Var\left[\sum \overline{\mathbb{6}}_{ij}\right] = E\left[\left(\sum \overline{\mathbb{6}}_{ij} - E[\sum \overline{\mathbb{6}}_{ij}]^{\,0}\right)^2\right]$$
$$= E\left[\left(\sum \mathbb{6}_{ij} - E[\sum \mathbb{6}_{ij}]\right)^2\right]$$
$$= Var\left[\sum \mathbb{6}_{ij}\right]$$

*(purple)* e.g. $(a^3+b^3)^2 \le (a^2+b^2)^3$
$$a^6 + 2a^3b^3 + b^6 \le a^6 + b^6$$
$$+ 3a^4b^2 + 3a^2b^4$$

So

$$Var\left[\sum_{i<j}\overline{\delta_{ij}}\right] = E\left[\left(\sum_{i<j}\overline{\delta_{ij}} - E\left[\sum_{i<j}\overline{\delta_{ij}}\right]\right)^2\right]$$

$$= E\left[\left(\sum_{i<j}\overline{\delta_{ij}}\right)^2\right]$$

$$= E\left[\sum_{i<j}\overline{\delta_{ij}}^2 + \sum_{\substack{i<j\\k<l\\i,j,k,l\ distinct}}\overline{\delta_{ij}}\,\overline{\delta_{kl}} + \sum_{\substack{i<j\\k<l\\i,j,l\ distinct}}\overline{\delta_{ij}}\,\overline{\delta_{il}} + \sum_{\substack{i<j\\k<j\\i,j,k\ distinct}}\overline{\delta_{ij}}\,\overline{\delta_{kj}}\right]$$

$$+ \sum \overline{\delta_{ij}}\,\overline{\delta_{jl}} \quad ⑤$$
$$+ \sum \overline{\delta_{ij}}\,\overline{\delta_{ki}} \quad ⑥$$

① ② ③ ④

---

① $\quad E\left[\sum_{i<j}\overline{\delta_{ij}}^2\right] \leq E\left[\sum \delta_{ij}^2\right] = \binom{s}{2}\|p\|_2^2$

$E[\delta_{ij}] = E[\delta_{ij}^2]$ since $\delta_{ij}$ is indicator var

**Trick helps here:** → gets rid of lots of terms

② $\quad$ independent

$$E\left[\sum_{\substack{i<j\\k<l\\all\ 4\ distinct}}\overline{\delta_{ij}}\,\overline{\delta_{kl}}\right] \leq \sum E[\overline{\delta_{ij}}]\,E[\overline{\delta_{kl}}] = 0$$

③ $\quad E\left[\sum_{\substack{i,j,l\\distinct}}\overline{\delta_{ij}}\,\overline{\delta_{il}}\right] \leq E\left[\sum \delta_{ij}\cdot\delta_{il}\right] = \sum_{\substack{i,j,l\\distinct}} pr[x_i = x_j = x_l]$

expected # 3-way collisions

$$\leq \binom{s}{3}\sum_x p(x)^3$$

$$\leq \frac{s^3}{6}\left(\sum_x p(x)^2\right)^{3/2}$$

$$\leq \frac{\sqrt{3}}{2}\binom{s}{2}^{3/2}\left(\|p\|_2^2\right)^{3/2} \quad \text{by the facts}$$

$$\frac{1}{6}(s^2)^{3/2} < \frac{\left(3\binom{s}{2}\right)^{3/2}}{6} = \frac{\sqrt{3}}{2}\binom{s}{2}^{3/2}$$

④ same as 3
⑤
⑥

In total:

$$\text{Var}\left[\sum_{i<j} b_{ij}\right] = \text{Var}\left[\sum_{i<j} \overline{b}_{ij}\right]$$

$$\leq \binom{5}{2}\|p\|_2^2 + 0 + 4 \cdot \frac{\sqrt{3}}{2}\left(\binom{5}{2}\|p\|_2^2\right)^{3/2}$$

$$\leq 4\left[\binom{5}{2}\|p\|_2^2\right]^{3/2}$$

∎

Putting lemma into Chebyshev:

use $p = \frac{\varepsilon^2}{2}$

$$\Pr\left[\left|\hat{c} - \|p\|_2^2\right| > \frac{\varepsilon^2}{2}\right] \leq \frac{Var[\hat{c}]}{\varepsilon^4} \cdot 4$$

note $\frac{1}{\left(\frac{s}{2}\right)^{1/2}} \leq \frac{1}{\sqrt{\frac{s^2}{2}}}$
$\leq \frac{2}{s}$

recall this from const. in front of $\hat{c}$ →

$$\leq \frac{4\left[\left(\frac{s}{2}\right)\|p\|_2^2\right]^{3/2}}{\left(\frac{s}{2}\right)^2 \, \varepsilon^4} \cdot 4 \qquad \leq \frac{32}{\varepsilon^4} \cdot \frac{1}{s} \cdot \|p\|_2^3$$

also want this to be $\leq 1$ $\qquad \leq 1$

So Pick $s \geq \Omega\left(\frac{1}{\varepsilon^4}\right)$

Note: Can get better bnd

1) Testing closeness to any known
   distribution — reduce to uniform case!

2) lower bound

How to estimate $\|p - u\|_1$ ?

1) $\|p-u\|_1 = 0 \iff \|p-u\|_2^2 = 0 \iff \|p\|_2^2 = \frac{1}{n}$

2) if $\|p-u\|_1 > \varepsilon \implies \|p-u\|_2 > \frac{\varepsilon}{\sqrt{n}}$

$\implies \|p-u\|_2^2 > \frac{\varepsilon^2}{n}$

$\implies \|p\|_2^2 \geq \frac{1}{n} + \frac{\varepsilon^2}{n}$

either additive estimate with error $\leq \frac{\varepsilon^2}{2n}$

or mult error $\leq (1 \pm \frac{\varepsilon^2}{3})$

suffices

would have this if have additive error $\leq \frac{\varepsilon^2}{3n} \cdot \|p\|_2^2$

to get additive error $\leq \frac{\varepsilon^2}{3n} \|p\|_2^2$

suffices to have

$s \geq \frac{const \cdot \sqrt{n}}{\varepsilon^2}$ samples

since $\Pr\left[ |\hat{c} - \|p\|_2^2| \geq \gamma \|p\|_2^2 \right] \leq \frac{k \cdot \|p\|_2^3}{s \cdot \gamma^2 (\|p\|_2^2)^2} \leq \frac{k}{s \cdot \gamma^2 \cdot \|p\|_2}$

$\overset{const}{\downarrow}$

$\left[ \text{note} \quad \|p\|_2^2 > \frac{1}{n} \text{ so } \|p\|_2 > \frac{1}{\sqrt{n}} \text{ so } \frac{1}{\|p\|_2} < \sqrt{n} \right]$

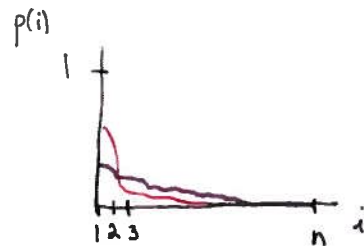$\leq \frac{k \cdot \sqrt{n}}{s \cdot \gamma^2}$ $\left[ \text{note: we need } \gamma \approx \frac{\varepsilon^2}{3} \right]$

So picking $s >> \frac{\sqrt{n}}{\varepsilon^4}$ will give small probability of error $\implies \approx \frac{k \cdot \sqrt{n}}{s} \cdot \frac{1}{\varepsilon^4}$

# Testing & Learning Monotone Distributions (over totally ordered domain)

<u>Def</u>. $p$ over $[n]$ is "monotone decreasing"

if $\forall\ i \in [n-1]$ $\quad p(i) \geq p(i+1)$



Monotonicity Tester:

· if $p$ monotone increasing, Pass with prob $\geq 3/4$

· if $p$ ε-far in $L_1$ dist from mon increasing, Fail with prob $\geq 3/4$

Useful tool: "Birge Decomposition"

(note: this is a different decomposition than in homework (upcoming)
in particular, it is <u>oblivious</u>!)

decompose domain $1..n$ into $\ell = \Theta\left(\frac{\log \varepsilon n}{\varepsilon}\right) \approx \Theta\left(\frac{\log n}{\varepsilon}\right)$ intervals

$$I_1^{\varepsilon},\ I_2^{\varepsilon},\dots I_{\ell}^{\varepsilon} \quad \text{s.t.}$$

← will drop ε
in notation
once it is fixed

$$\left|I_{k+1}^{\varepsilon}\right| = \left\lfloor (1+\varepsilon)^{i} \right\rfloor$$

$$\left|I_1^{\varepsilon}\right| = \left|I_2^{\varepsilon}\right| = \dots = 1$$
$$\left|I_a^{\varepsilon}\right| = \left|I_{a+1}^{\varepsilon}\right| = \dots = 2$$
$$\vdots$$

but then at some point the sizes grow
exponentially

define "flattened distribution"

$$\forall \; 1 \leq j \leq \ell$$
$$\forall \; i \in I_j \qquad \tilde{q}_\varepsilon (i) = \frac{q(I_j)}{|I_j|}$$

← assign all elements in same interval the same probability

note: $q(I_j) = \tilde{q}_\varepsilon(I_j)$

__Birgé's Thm__ if $q$ mon decreasing then $\|\tilde{q}_\varepsilon - q\|_1 < \varepsilon$

__Corr__ if $q$ $\varepsilon$-close to mon decreasing then $\|\tilde{q}_\varepsilon - q\|_1 < O(\varepsilon)$

__Testing Algorithm :__

how can we do this? $\tilde{q}$ isn't even if $q$ monotone, exactly uniform. See problem from next hw set.

Take samples of $q$
do uniformity test for each partition (using samples that fell in it)
    (if not enough samples then pass)    fail if any partition fails

$w_j \leftarrow$ # samples that fell in partition $j$
use LP to verify $w$ close to monotone

↖ note this is LP on $O(\log n)$ vars

__How many samples?__
for each partition with enough weight, say $\frac{\varepsilon}{\log n}$, need $\frac{\sqrt{n}}{\varepsilon^2}$ samples

$\approx \; O(\sqrt{n} \; polylog \, n \cdot poly \frac{1}{\varepsilon})$

↖ need $\frac{\sqrt{n} \cdot \log n}{\varepsilon^3}$ for each one
need another $\log\log n$ for union bound

(note: this can be improved !!)

## Last step:

difficulty

Sampling error might make $w_j$'s <u>look</u> non monotone

purple is <u>not</u> monotone
but is <u>close</u>

good thing: only $\frac{\log n}{\sqrt{\varepsilon}}$ variables!

Can be solved via brute force
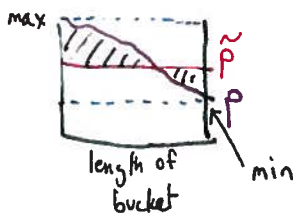LP (actually quite efficient)
$\vdots$

so: monotone $p$ likely to pass
$\varepsilon$-far from monotone $p$: either (1) non uniform in buckets
or (2) $w$ far from monotone

Slightly changing perspective...

What if we <u>know</u> dist $q$ is monotone, can we <u>learn</u> it?

Yes! use sampling to estimate $\tilde{q}_\varepsilon(I_j)$'s

<u>Birge's Thm</u> ⇒ Can learn monotone distributions to w/in $\varepsilon$ $L_1$ error

in $\Theta\left(\frac{1}{\varepsilon^3} \log n\right)$ samples.

# Proof of Birge's Thm :

Error in bucket



gross upper bound on error:
$$\leq (\text{max} - \text{min}) \cdot \text{bucket length}$$

## Partition of Intervals:

- Size 1 Intervals $|I_j| = 1$
- Short Intervals $|I_j| < 1/\varepsilon$
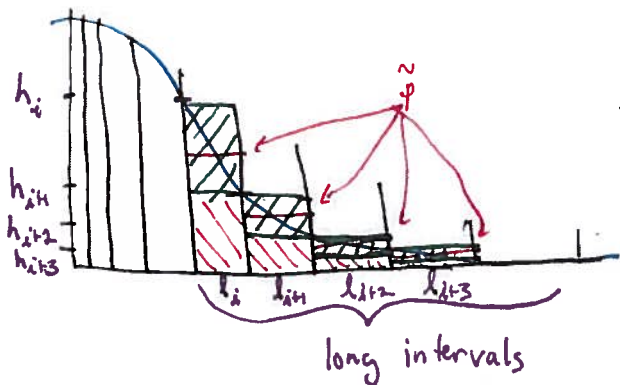- Long Intervals $|I_j| \geq 1/\varepsilon$

← if we have any short intervals, there are $\Omega(1/\varepsilon)$ of these
if not, we can learn the distribution

⊃ if we have these then
max prob $\leq \varepsilon$ (since # size 1 intervals is $\Omega(1/\varepsilon)$)

↑ therefore min size 1 interval has prob $\leq \varepsilon$ which upper bounds later probabilities too since $p$ is monotone

$$\text{total error} \leq \sum_{j=1}^{\ell} |I_j| \cdot (\text{max prob in } I_j - \text{min prob in } I_j)$$

$$= \underbrace{\sum 1 \cdot 0}_{\substack{\text{size 1} \\ \text{intervals}}} + \underbrace{\sum |I_j| (\text{max} - \text{min})}_{\substack{\text{short} \\ \text{intervals}}} + \underbrace{\sum |I_j| (\text{max} - \text{min})}_{\substack{\text{long} \\ \text{intervals}}}$$

$\underbrace{\phantom{xxxxx}}$
0
since no difference

$\underbrace{\phantom{xxxxx}}$
Omitted: idea is bound similarly to the long intervals but need to group together intervals of same size

$\underbrace{\phantom{xxxxx}}$
see below

## Picture for long intervals:



long intervals

green rectangles = upper bnd on error

$$\text{error} \leq (h_i - h_{i+1}) \ell_i + (h_{i+1} - h_{i+2}) \ell_{i+1} + (h_{i+2} - h_{i+3}) \ell_{i+2} + \ldots$$

$$= h_i \ell_i + h_{i+1} (\ell_{i+1} - \ell_i) + h_{i+2} (\ell_{i+2} - \ell_{i+1}) + h_{i+3} (\ell_{i+3} - \ell_{i+2})$$

all $h$'s in this area are $< \varepsilon$!

positive, $+ \approx \varepsilon \cdot \ell_{i+1}$ by way that we partitioned

$$\leq \varepsilon \left[ \ell_i + \sum h_i \ell_{i-1} \right]$$

get rid of this when bounding short intervals ↗

this is area of red rectangles, which is upper bounded by $p$ so sum is $\leq 1$