

Lecture 16:

Hypothesis Testing

Some Problems: (Given samples of p)

is $p = q$ (e.g. $q = U_D$)
or ϵ -far from q

Complexity (in terms of $n = |D|$)

$$\sqrt{n}$$

is p ϵ -close to q
or 2ϵ -far from q

$$\frac{n}{\log n}$$

(Given samples of q) is $p = q$
or p ϵ -far from q

$$n^{2/3}$$

(Given samples of q) is p ϵ -close to q
or 2ϵ -far from q

$$\frac{n}{\log n}$$

is p monotone
or ϵ -far from monotone

$$\sqrt{n}$$

is p ϵ -close to monotone
or ϵ -far from monotone

$$\frac{n}{\log n}$$

ignoring ϵ 's
 $\log n$'s

Other problems considered:

estimate entropy, support size

Independence?

represented well via k-histogram?

monotone hazard rate

-
-
-

A useful tool:

Given: (1) collection of distributions (via complete description) \mathcal{H}

(2) Samples of p such that $\exists q \in \mathcal{H}$ for which $\text{dist}(p, q)$ is small

\mathcal{H} contains good approximation to p

Goal: Output $h \in \mathcal{H}$ s.t. $\text{dist}(p, h)$ small

↑↑
strong assumption

Question:

How many samples needed in terms of $|\mathcal{H}|$ + domain size?

Is this the same as testing closeness, uniformity?

Do lower bounds apply?

NO! guaranteed p close to
some $q \in \mathcal{H}$

What we want:

Given h_1, h_2 explicit distributions
 p via samples

procedure that outputs h_i that is closer to p

What if both are roughly same distance?

maybe either one is ok?

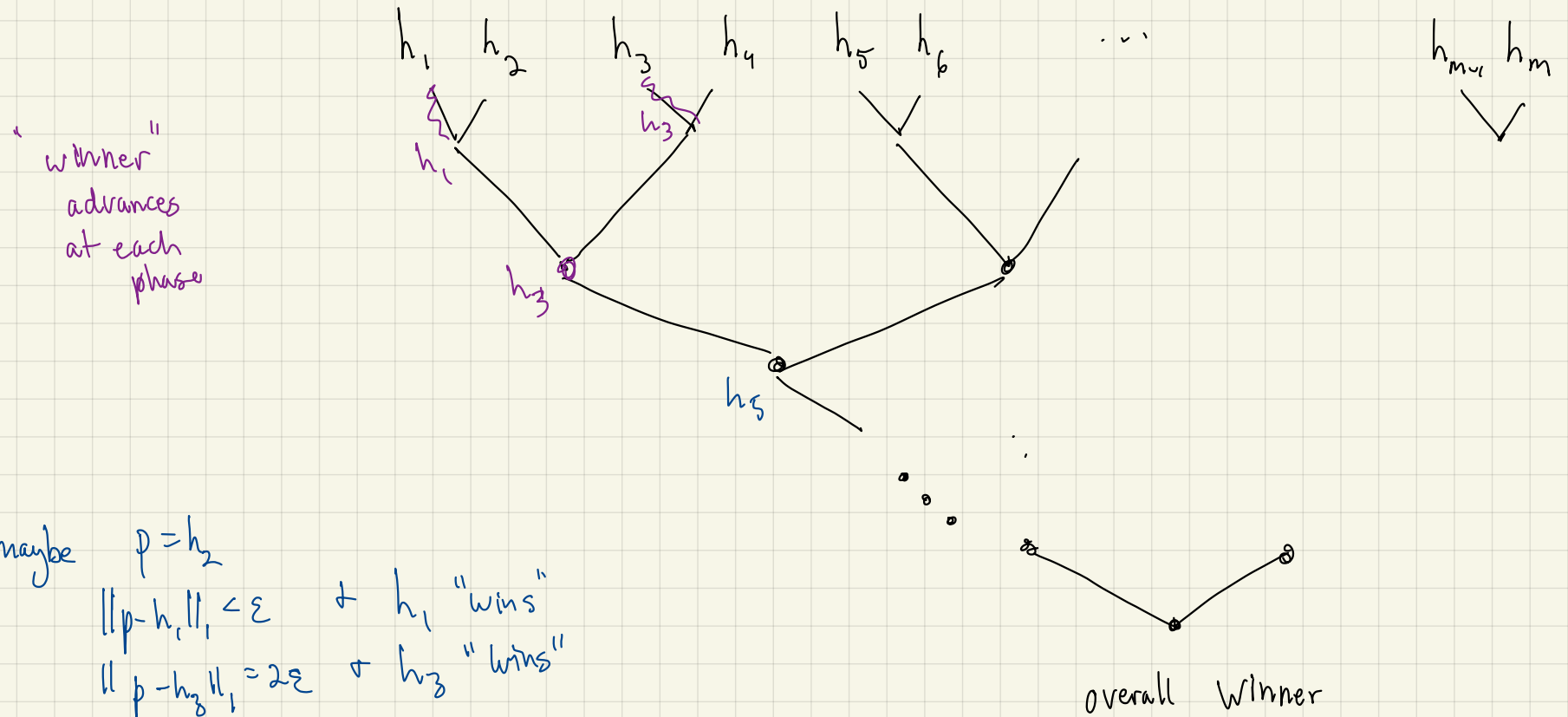
or maybe not...

More general Goal:

Given set of hypotheses \mathcal{H}
 \dagger p via samples

find $h \in \mathcal{H}$ closest to p
or pretty close to best

Find best hypothesis via "tournament"?



"winner" advances at each phase

maybe $p = h_2$
 $\|p - h_1\|_1 < \epsilon$ if h_1 "wins"
 $\|p - h_3\|_1 = 2\epsilon$ if h_3 "wins"
 $\|p - h_5\|_1 = 3\epsilon$ if h_5 "wins"
 \vdots
 \uparrow

overall winner $\leq O(\log n \cdot \epsilon)$ distance from best hypothesis?

- won't use simple tournament

← instead compare every pair

- will add notion of "tie"

Output hypothesis that wins or ties
every match

- hopefully there is one

- " it is close to p

A "subtool" for comparing two hypotheses:

Thm given (1) sample access to p
(2) h_1, h_2 hypothesis distributions (fully known to algorithm)
(3) accuracy parameter ϵ' , confidence parameter δ'

then Algorithm "choose" takes $O(\log(\frac{1}{\delta'}) / (\epsilon')^2)$ samples + outputs
 $h \in \{h_1, h_2\}$ satisfying:

if one of h_1, h_2 has $\|h_i - p\|_1 < \epsilon'$

then with prob $\geq 1 - \delta'$, output h_j has $\|h_j - p\|_1 < 12\epsilon'$

i.e. if both h_1, h_2 far no guarantees

if one ϵ' -close + one is really far ($> 12\epsilon'$) we will output close one

if one is ϵ' close \rightarrow output either one?
other is $10\epsilon'$ close but both $\leq 10\epsilon'$ close so not too bad

getting kind of complicated just to specify?

Actually a bit stronger: (focus on case where ≥ 1 close)

Thm p given via samples
 h_1, h_2 fully known + p is ϵ' -close to at least one of h_1, h_2
 ϵ', δ' given

Algorithm "choose" takes $O((\log \frac{1}{\delta'}) (\frac{1}{\epsilon'})^2)$ samples + outputs $h \in \{h_1, h_2\}$ such that:

(1) If h_i more than $\underline{12\epsilon'}$ -far from p , unlikely to output h_i as winner
very far $2e^{-m(\epsilon')^2/2}$ or tie

(2) If h_i more than $\underline{10\epsilon'}$ -far from p , unlikely to output h_i as winner
pretty far but not too bad \uparrow BUT COULD TIE

Proof of subtool:

idea:

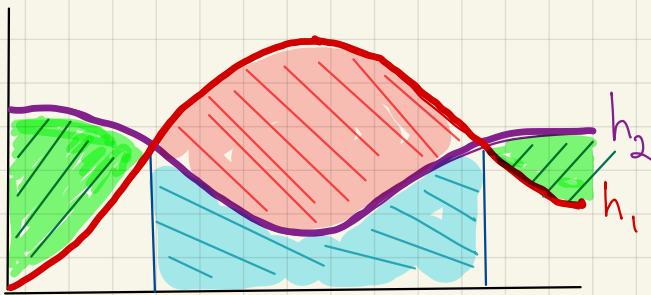
wlog h_1 is ε' -close to p

if h_2 is close to h_1 , then h_2 is pretty close to p
+ can win or tie

else h_2 is far from h_1 , then also far from p

red + green areas are big

distance of $h_1 + h_2$



estimate $\text{pr}[x \in A]$
 $x \neq p$

see if it is more
like h_1 or
more like h_2

$A =$ part of domain where $h_1(x) > h_2(x)$

Algorithm Choose: Input p, h_1, h_2

\swarrow samples
 \nwarrow explicit description

First some definitions:

$$A = \{x \mid h_1(x) > h_2(x)\}$$

$$a_1 = h_1(A) \quad \leftarrow \text{red + blue areas}$$

$$a_2 = h_2(A) \quad \leftarrow \text{blue area}$$

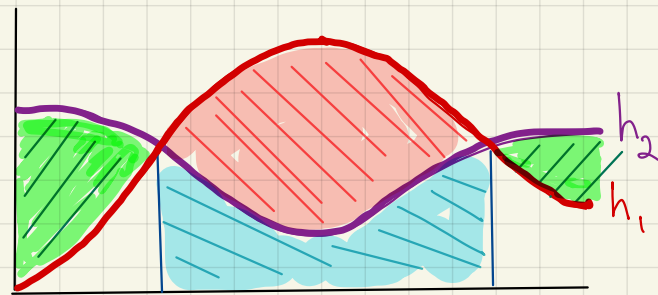
$$\text{note } \|h_1 - h_2\|_1 = 2(a_1 - a_2)$$

1. if $a_1 - a_2 \leq 5\epsilon'$ declare "tie" + return
(no samples needed)

2. draw $m = 2 \log \frac{1}{\delta'}$ samples S_1, \dots, S_m from p
 $(\epsilon')^2$

3. $\alpha \leftarrow \frac{1}{m} |\{i \mid S_i \in A\}|$ \leftarrow estimates $\Pr_{x \in p}[x \in A]$

4. if $\alpha > a_1 - \frac{3}{2}\epsilon'$ return h_1
else if $\alpha < a_2 + \frac{3}{2}\epsilon'$ return h_2
else declare "tie"



$$\text{green area} = \text{red area} = a_1 - a_2$$

$$L_1 \text{ dist} = \text{green} + \text{red} = 2 \cdot \text{red} = 2(a_1 - a_2)$$

$$C = \sum_x \min(h_1(x), h_2(x))$$

$$\text{green} = \sum h_2(x) - C$$

$$= 1 - C$$

$$\text{red} = \sum h_1(x) - C$$

$$= 1 - C$$

Why does it work?

- h_1 or h_2 is ε' -close to A (given)

- If "tie" in step 1:

$h_1 + h_2$ are $\underbrace{10\varepsilon'}_{=2(a_1 - a_2)}$ -close to each other

\Rightarrow both $\leq 11\varepsilon'$ -close to $p \Rightarrow$ "tie" is OK output

- otherwise reach step 2:

$$\|h_1 - h_2\|_1 > 10\varepsilon' \quad (a_1 - a_2 > 5\varepsilon')$$

Algorithm Choose:

$$A = \{x \mid h_1(x) > h_2(x)\}$$

$$a_1 = h_1(A)$$

$$a_2 = h_2(A)$$

$$\text{note } \|h_1 - h_2\|_1 = 2(a_1 - a_2)$$

1. if $\underbrace{a_1 - a_2}_{\frac{1}{2}L_1 \text{ dist}} \leq 5\varepsilon'$ declare "tie" + return h (no samples needed)

2. draw $m = 2 \frac{\log \frac{1}{\delta'}}{(\varepsilon')^2}$ samples S_1, \dots, S_m from p

$$3. \alpha \leftarrow \frac{1}{m} |\{i \mid S_i \in A\}|$$

4. if $\alpha > a_1 - \frac{3}{2}\varepsilon'$ return h_1
 else if $\alpha < a_2 + \frac{3}{2}\varepsilon'$ return h_2
 else declare "tie" + return h_1

$\left\{ \begin{array}{l} \text{if } p = h_1, E[\alpha] = a_1 \\ \text{if } p = h_2, E[\alpha] = a_2 \end{array} \right.$



green area = red area = $a_1 - a_2$
 $L_1 \text{ dist} = \text{green} + \text{red}$
 blue area = a_2
 blue + red area = a_1

Why does it work?

- h_1 or h_2 is ϵ' -close to A (given)
- If "tie" in step 1, algorithm does right thing
- Otherwise reach step 2: $\|h_1 - h_2\|_1 > 10\epsilon'$ ($a_1 - a_2 > 5\epsilon'$)

$$E[\alpha] = \Pr_{x \in p} [x \in A] = p(A)$$

assume (Chernoff) that with high prob $|\alpha - E[\alpha]| \leq \frac{\epsilon'}{2}$

h_1 assigns a_1 weight to A
 h_2 " " a_2 " " A

if p is ϵ' -close to h_1 , assigns $\geq a_1 - \epsilon'$ wt to A
 $\alpha \geq a_1 - \epsilon' - \frac{\epsilon'}{2}$ return h_1 whp

if p is ϵ' -close to h_2 , assign $\leq a_2 + \epsilon'$ wt to A
 $\alpha \leq a_2 + \epsilon' + \frac{\epsilon'}{2}$ return h_2 whp

Algorithm Choose:

$$A = \{x \mid h_1(x) > h_2(x)\}$$

$$a_1 = h_1(A)$$

$$a_2 = h_2(A)$$

$$\text{note } \|h_1 - h_2\|_1 = 2(a_1 - a_2)$$

1. if $a_1 - a_2 \leq 5\epsilon'$ declare "tie" + return h (no samples needed)
2. draw $m = 2 \frac{\log \frac{1}{\delta'}}{\epsilon'^2}$ samples S_1, \dots, S_m from p
3. $\alpha \leftarrow \frac{1}{m} |\{i \mid S_i \in A\}|$
4. if $\alpha > a_1 - \frac{3}{2}\epsilon'$ return h_1
 else if $\alpha < a_2 + \frac{3}{2}\epsilon'$ return h_2
 else declare "tie" + return h_1

$\left\{ \begin{array}{l} \text{if } p = h_1, E[\alpha] = a_1 \\ \text{if } p = h_2, E[\alpha] = a_2 \end{array} \right.$



green area = red area = $a_1 - a_2$
 L_1 dist = green + red
 blue area = a_2
 blue + red area = a_1

The cover method - a method for learning distributions

def. \mathcal{C} is an " ϵ -cover" of \mathcal{D} if $\forall p \in \mathcal{D}$
 $\exists q \in \mathcal{C}$ s.t. $\|p - q\|_1 \leq \epsilon$

smaller set of distributions (pointing to \mathcal{C})
big set of distributions (pointing to \mathcal{D})

Why useful? hopefully \mathcal{C} much smaller than \mathcal{D} it allows to approx \mathcal{D}
note \mathcal{C} not unique

Thm \exists algorithm, given $p \in \mathcal{D}$, which takes
 $O\left(\frac{1}{\epsilon^2} \log |\mathcal{C}|\right)$ samples of p + outputs $h \in \mathcal{C}$
s.t. $\|h - p\|_1 \leq 6\epsilon$ with prob $\geq \frac{9}{10}$

union bnd over $|\mathcal{C}|$ rather than $|\mathcal{D}|$ →

Thm \exists algorithm, given $p \in \mathcal{D}$, which takes
 $O\left(\frac{1}{\epsilon^2} \log |\mathcal{C}|\right)$ samples of p + outputs $h \in \mathcal{C}$
 s.t. $\|h - p\|_1 \leq 6\epsilon$ with prob $\geq \frac{9}{10}$

Pf.

Since $p \in \mathcal{D}$, $\exists q_{\text{opt}} \in \mathcal{C}$ s.t. $\|p - q_{\text{opt}}\|_1 \leq \epsilon$

run "Choose" on p with every pair $q_1, q_2 \in \mathcal{C}$

q_{opt} : win or tie all games

if q' is $\geq 6\epsilon$ -far from p , then $\geq 5\epsilon$ -far from q_{opt}
 \Rightarrow loses to q_{opt}

equivalently: if q' wins or ties all games
 $\Rightarrow \leq 5\epsilon$ far from q_{opt}
 $\leq 6\epsilon$ far from p

need all matches to give correct output (estimate of α)
 union bound on $\binom{|\mathcal{C}|}{2}$ many matches \square

Applications:

Example 1: learning distribution of a coin

domain = $\{0, 1\}$

need to learn bias

Here $\mathcal{D} = [0, 1]$

if use $\mathcal{C} = \left\{ 0, \frac{1}{k}, \frac{2}{k}, \frac{3}{k}, \dots, \frac{k}{k} \right\}$

\forall bias p let $\frac{i}{k} < p \leq \frac{i+1}{k}$

so pick

$\tilde{p} \leftarrow$ closest $\frac{i}{k}$ then $|p - \tilde{p}| = \frac{1}{k} \leq \epsilon$

picking $k = \Theta\left(\frac{1}{\epsilon}\right)$

$|\mathcal{C}| = k+1 = \Theta\left(\frac{1}{\epsilon}\right)$

need $O\left(\frac{1}{\epsilon^2} \log \frac{1}{\epsilon}\right)$ samples

using cover method

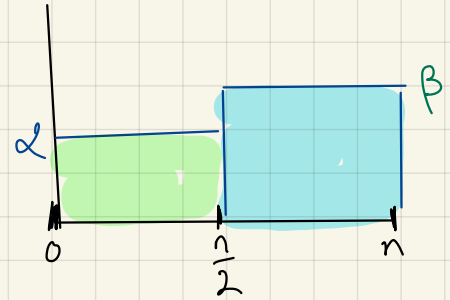
Example 2: 2-bucket distributions

now need to specify α and β

$$\text{so } \mathcal{C} = \left\{ \left(\frac{i}{k}, \frac{j}{k} \right) \mid i, j \in \{0, \dots, k\} \right\}$$

$$|\mathcal{C}| = \Theta\left(\left(\frac{1}{\varepsilon}\right)^2\right)$$

$$\# \text{ samples is } O\left(\frac{1}{\varepsilon^2} \log \frac{1}{\varepsilon}\right)$$



Example 3: monotone distributions

$$\text{Birge} \Rightarrow \mathcal{C} = \left\{ \left(\frac{i_1}{k}, \dots, \frac{i_{\lceil \log n / \varepsilon \rceil}}{k} \right) \mid i_1, i_2, \dots \in \{0, \dots, k\} \right\}$$

$$|\mathcal{C}| = \Theta\left(\frac{1}{\varepsilon^{\lceil \log n / \varepsilon \rceil}}\right) \Rightarrow \# \text{ samples is } O\left(\frac{1}{\varepsilon^3} \log n \log \frac{1}{\varepsilon}\right)$$