

Lecture 2

Jan 11, 2024

- Sortedness Testing (Cont.)
- Concentration of random variables
(Markov, Chebyshev, Chernoff, Hoeffding)
- Running example: estimating coin bias.

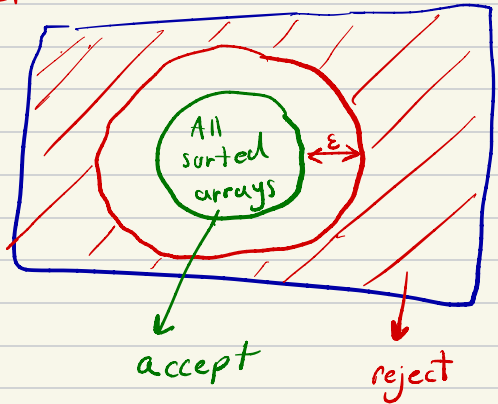
Testing sortedness (general case)

Recall def

Give an array A ,

Design an algorithm \mathcal{A} s.t. with prob. $1-\delta$

- if A is sorted, \mathcal{A} outputs **accept**
- if A is ϵ -far from being sorted, \mathcal{A} outputs **reject**



what does ϵ -far mean here?

distance between two array of size n :

entries we need to change to
change A to A'

$$\text{dist}(A, A') = \frac{\text{number of entries to change}}{n}$$

$P = \{ \text{all sorted arrays} \}$

$$\text{dist}(A, P) = \min_{A' \in P} \text{dist}(A, A')$$

ϵ -far from sortedness = We need to
change $\geq \epsilon \cdot n$ entries in A to get
a sorted array.

why randomly throwing darts in the
dark won't work in general case?

unable to detect local changes



New algorithm

Binary search base algorithm.

Sorted \Rightarrow binary search works.
" $\stackrel{?}{\Leftarrow}$ "

Assumption : WLOG, entries of A
are distinct.

Try s times

pick a random $i \in [n]$

$l \leftarrow \text{Binary search}(A, A[i])$

if $(l \neq i)$

return reject

return accept

* If A is sorted \Rightarrow

all calls to binary search work correctly

\Rightarrow the algorithm returns **accept** w. prob 1.

* If A is ϵ -far from sorted \Rightarrow ?

$p :=$ the probability of binary search fails
when we are ϵ -far

what we need:

$$\Pr[\text{outputting } \text{accept} \mid \epsilon\text{-far}] = (1-p)^s \leq \delta$$

$$\text{by setting } s = \frac{\log(1/\delta)}{p}$$

if $p < \epsilon \Rightarrow (1-\epsilon) \cdot n$ many entries
are nice

\Downarrow Lemma 1

$(1-\epsilon) \cdot n$ many entries are sorted

\Downarrow

A is not ϵ -far from being
sorted

$$\Rightarrow p \geq \epsilon \Rightarrow S = \frac{\log 1/\delta}{\epsilon}$$

would be enough.

Binary_search (array A, value x, indices h, t)

if (t < h)

return h

m ← $\lfloor \frac{h+t}{2} \rfloor$

if (A[m] = x)

return m

if (A[m] > x)

return binary_search (A, x, h, m-1)

if (A[m] < x)

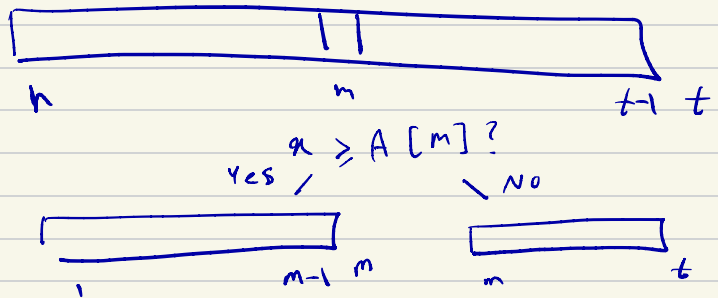
return binary_search (A, x, m+1, t)

Binary search on a sorted A
returns the smallest i such that

$$A[i] \geq x$$

Binary search $(A, x, 1, n+1)$

could return $n+1$.



$A[m]$ is called pivot

We say i is nice if the binary search on $x = A[i]$ returns i

Lemma 1 Suppose we have two nice indices i and $j \in [n]$. If $i < j$ then $A[i] < A[j]$

Proof
pivots of i

$m_1^{(i)}, m_2^{(i)}, \dots, m_{k_i}^{(i)} = i$

pivots of j

$m_1^{(j)}, m_2^{(j)}, \dots, m_{k_j}^{(j)} = j$

m^* = last mutual pivot



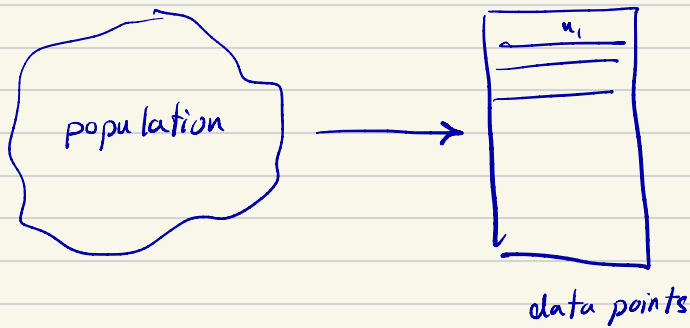
$$A[i] < A[m^*] < A[j]$$

Hypothesis testing (property testing of distions)

We used randomness to model the world.

data points are random samples from

an unknown data distribution



distribution p

x_1, \dots, x_m

$x_i \sim p$

Estimating coin bias

$$p = \Pr[\text{head}]$$

Testing a coin is fair:

- if $p = \frac{1}{2}$, output **accept** w. prob. $1 - \delta$.
- if $|p - \frac{1}{2}| > \epsilon$, output **reject** w. prob $1 - \delta$.

Algorithm

Flip a coin m times

$X \leftarrow$ # heads

$$\text{if } \left| \frac{X}{m} - \frac{1}{2} \right| \leq \epsilon$$

return **accept**

else

return **reject**

Question: How well $\frac{\bar{X}}{m}$ approximate p ?

what should be m ?

boils down \rightarrow How well $\frac{\bar{X}}{m}$ concentrate
around p ?

Concentration of random variables.

Questions: { Estimating average height of students
exit polls

n samples:

$$X_1, X_2, \dots, X_n \sim P$$

$$\bar{X}_n := \frac{1}{n} \sum_{i=1}^n X_i \rightarrow \mu := E_{X \sim P} [X]$$

Goal measure how much \bar{X}_n deviates from μ

Law of Large numbers

(weak) $\forall \varepsilon \quad \lim_{n \rightarrow \infty} \Pr[|\bar{X}_n - \mu| < \varepsilon] = 1$

(strong) $\Pr\left[\lim_{n \rightarrow \infty} \bar{X}_n = \mu\right] = 1$

Central Limit Theorem:

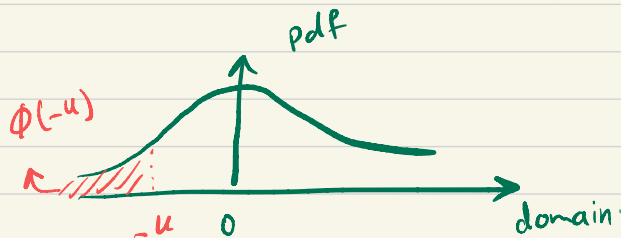
Var_{X~P}[X]

$$\sqrt{n} (\bar{X}_n - \mu) \rightarrow N(0, \sigma^2)$$

$$Z \sim N(0,1)$$

$$\Pr \left[\frac{\sqrt{n} |\bar{X}_n - \mu|}{\sigma} > u \right] \approx \Pr [|Z| > u] \\ = 2\Phi(-u)$$

where Φ is the cdf of the standard normal dist.



Look up table

$$u = 1.96 \rightarrow 2\Phi(-u) \approx 95\%$$

Hence: with prob. 0.95

$$\mu \in \left[\bar{X}_n - 1.96 \sigma / \sqrt{n}, \bar{X}_n + 1.96 \sigma / \sqrt{n} \right]$$

- Quality of Approximation varies depending on P .

These are asymptotic results. Very general, but

- work in the limit,

- Do not indicate the relationship among the parameters,

n, d, ϵ, δ ?

↓
dimension

↓
error

↘ confidence (in our example δ
has $1 - 0.95 = 0.05$)

what about finite sample setting?

Usefull tools to show concentration (tail bounds)

Markov's inequality:

For non-negative random variable X , and $a > 0$:

$$\Pr[X \geq a] \leq \frac{\mathbb{E}[X]}{a}$$

proof.

$$\mathbb{E}[X] = \int_0^{\infty} x \Pr[X=x] dx$$

pdf ↙

$$= \int_0^a x \Pr[X=x] dx + \int_a^{\infty} x \Pr[X=x] dx$$

$$\geq 0 + \int_a^{\infty} \Pr[X=x] dx$$

$$\geq a \Pr[X \geq a]$$

$$\Rightarrow \Pr[X \geq a] \leq \frac{\mathbb{E}[X]}{a} \quad \square$$

→ back to coin example

works well for small p

if $p \leq 0.01$

$$\Pr\left[\frac{X}{m} > 0.1\right] \leq \frac{E[X]}{0.1} \leq 0.1$$

not very meaningful when $p = \frac{1}{2}$

Chebyshev's inequality

For a random variable with finite mean and variance, and $k > 0$:

$$\Pr [|X - \mathbb{E}[X]| \geq k \sigma] \leq \frac{1}{k^2}$$

proof: ↙ standard deviation of X

$$\Pr [|X - \mathbb{E}[X]| \geq k \sigma]$$

$$= \Pr [(X - \mathbb{E}[X])^2 \geq k^2 \sigma^2]$$

$$\leq \frac{\mathbb{E} [(X - \mathbb{E}[X])^2]}{k^2 \sigma^2} = \frac{\sigma^2}{k^2 \sigma^2} = \frac{1}{k^2} \quad \square$$

↙ Markov

→ back to coin example

$$E\left[\frac{X}{m}\right] = p$$

$$\text{Var}\left[\frac{X}{m}\right] = \frac{p(1-p)}{m}$$

$$\Pr\left[\left|\frac{X}{m} - p\right| > \varepsilon\right] \leq \frac{\text{Var}\left[\frac{X}{m}\right]}{\varepsilon^2} \leq \frac{1}{m\varepsilon^2} \leq \delta$$

$$m = \frac{1}{\delta \cdot \varepsilon^2}$$

right dependencies to ε

but not δ

Chernoff bound:

m Bernoulli random variable: X_1, X_2, \dots, X_m

$$X_i \sim \text{Ber}(p_i) \quad X_i = \begin{cases} 1 & \text{with prob } p_i \\ 0 & \text{with prob } 1-p_i \end{cases}$$

empirical mean $X := \frac{1}{m} \sum_{i=1}^m X_i$

and true mean $\mu := \frac{1}{m} \sum_{i=1}^m p_i$

$$\Pr [X - \mu > \epsilon \mu] \leq e^{-m p \epsilon^2 / 3}$$

$$\Pr [\mu - X > \epsilon \mu] \leq e^{-m p \epsilon^2 / 2}$$

general structure of the proof:

(can be applied to any random variable)

For all $\varepsilon > 0$, $t > 0$:

$$\Pr[X > \varepsilon] = \Pr[e^{tX} > e^{t\varepsilon}]$$

$$\leq \frac{E[e^{tX}]}{e^{t\varepsilon}} = e^{-t\varepsilon} M_X(t)$$

↓
Markov

↓
moment generating func

since the bound holds for any t , we can conclude:

$$\Pr[X \geq \varepsilon] \leq \inf_{t > 0} e^{-t\varepsilon} M_X(t) \quad \square$$

→ back to coin example

$$Pr[|X - p| \geq \epsilon] \leq 2 \exp(-mp\epsilon^2)$$

$$m = \frac{1}{p\epsilon^2} \log \frac{2}{\delta} \leq \delta$$

works when $p = 1/2$. ✓

Hoeffding bound:

$$\Pr [X - \mu > \varepsilon] < e^{-2m\varepsilon^2}$$

$$\Pr [\mu - X < \varepsilon] < e^{-2m\varepsilon^2}$$

→ back to coin example

$$m = \frac{\log(2/\delta)}{2\varepsilon^2} \Rightarrow$$

$$\Pr [|X - \mu| > \varepsilon] \leq \delta$$