

Lecture 3

Jan 23, 2025

- Concentration of random variables

(Markov, Chebyshev, Chernoff, Hoeffding)

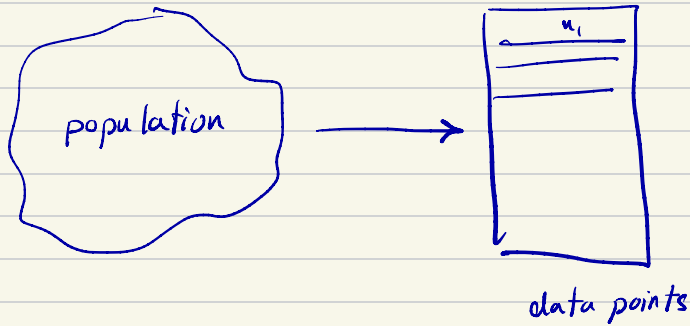
- Running example: estimating coin bias.

Hypothesis testing (property testing of distributions)

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

if $|\frac{X}{m} - \frac{1}{2}| \leq ?$

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:

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

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

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:

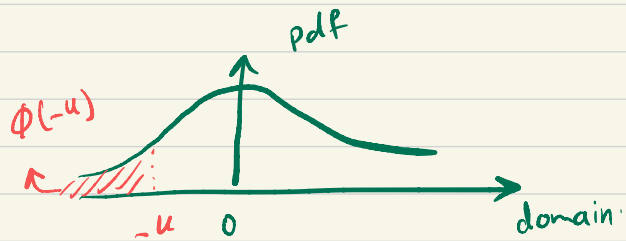
Var_{X~P}[X]

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

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

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

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 ?

\downarrow dimension \downarrow error \rightarrow 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 ↙

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

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

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

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

→ 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 $p = 1/2$. ✓

(or when p is a constant)

Hoeffding bound:

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

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

→ back to coin example

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

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