

CS70 LECTURE 29 : CONCENTRATION BOUNDS

Consider a nonnegative variable X .

For any a , let Y be a random variable $\text{let } Y = \begin{cases} 0 & \text{if } X < a \\ a & \text{if } X \geq a \end{cases}$

Then $Y \leq X$.

Since $X - Y \geq 0$, $\mathbb{E}[X - Y] \geq 0 \Rightarrow \mathbb{E}[X] - \mathbb{E}[Y] \geq 0 \Rightarrow \mathbb{E}[X] \geq \mathbb{E}[Y]$.

$$\mathbb{E}[Y] = 0 \cdot \mathbb{P}[X < a] + a \mathbb{P}[X \geq a] \leq \mathbb{E}[X]$$

$$\Rightarrow a \mathbb{P}[X \geq a] \leq \mathbb{E}[X] \Rightarrow \boxed{\mathbb{P}[X \geq a] \leq \frac{\mathbb{E}[X]}{a}}$$

Markov's Inequality

$$\begin{matrix} X \geq 0 \\ a > 0 \end{matrix} \Rightarrow \mathbb{P}[X \geq a] \leq \frac{\mathbb{E}[X]}{a}$$

Chebyshev's Inequality

Special application of Markov's Inequality

let X be a random variable (not necessarily nonnegative)

$$\text{let } \mu = \mathbb{E}[X]. \quad \sigma_x^2 \triangleq \mathbb{E}[(X - \mu)^2]$$

Chebyshev's inequality states that:

$$\mathbb{P}[|X - \mu| \geq c] \leq \frac{\sigma_x^2}{c^2}$$

Upperbound on how far X can be away from its mean.

i.e. if the variance is small, then the probability that X can be very much away from the mean is smaller.

Proof:

let $Y = |X - \mu|^2$. Then $Y \geq 0$.

let $a = c^2$.

Applying Markov's inequality, $\mathbb{P}[Y \geq c^2] \leq \frac{\mathbb{E}[Y]}{c^2}$

$$\Rightarrow \mathbb{P}[Y \geq c^2] = \mathbb{P}[|X - \mu|^2 \geq c^2] = \mathbb{P}[|X - \mu| \geq c] \leq \frac{\mathbb{E}[Y]}{c^2}$$

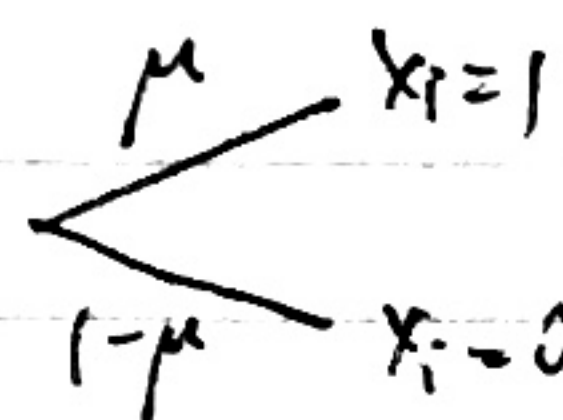
$$\therefore \mathbb{P}[|X - \mu| \geq c] \leq \frac{\mathbb{E}[Y]}{c^2} = \frac{\mathbb{E}[|X - \mu|^2]}{c^2} = \frac{\text{Var}[X]}{c^2}$$

Loose Interpretation for Markov's Inequality

If $\mathbb{E}[X]$ is small, then the probability that X takes on large values of a is small

Sample Mean

$$M_n = \frac{X_1 + X_2 + \dots + X_n}{n} \quad (\text{Sample mean})$$



μ is unknown, but fixed
Want to estimate μ .

X_i 's are independent, identically distributed (based on the same underlying Bernoulli)
(e.g. n independent flips of a coin)

The sample mean M_n is also a random variable in its own right.

$$E[M_n] = E\left[\frac{X_1 + X_2 + \dots + X_n}{n}\right] = \frac{1}{n} \sum_{i=1}^n E[X_i] = E[X]$$

M_n is an "estimator" for $E[X]$

all equal
share based off same
underlying experiment

$$\therefore \boxed{E[M_n] = E[X]}$$

M_n is an "unbiased estimator" for $E[X]$

Unbiased estimator: the mean of the estimator is the same as the mean of the unknown variable.

$$\text{Var}[M_n] = \text{Var}\left[\frac{X_1 + X_2 + \dots + X_n}{n}\right] = \left(\frac{1}{n}\right)^2 \text{Var}[X_1 + X_2 + \dots + X_n]$$

$$= \left(\frac{1}{n}\right)^2 (\text{Var}[X_1] + \dots + \text{Var}[X_n])$$

$$= \left(\frac{1}{n}\right)^2 \cdot n \text{Var}[X] = \frac{\text{Var}[X]}{n} \quad \therefore \boxed{\text{Var}[M_n] = \frac{\text{Var}[X]}{n}}$$

$$\lim_{n \rightarrow \infty} \text{Var}[M_n] = 0$$

As n increases, obtain a better and better estimate of the true mean.

Weak Law of Large Numbers

$$P(|M_n - E[M_n]| \geq \epsilon) \leq \frac{\text{Var}[M_n]}{\epsilon^2} \quad (\text{by Chebyshev})$$

$$\Rightarrow \boxed{P(|M_n - E[X]| \geq \epsilon) \leq \frac{\text{Var}[X]}{n\epsilon^2}}$$

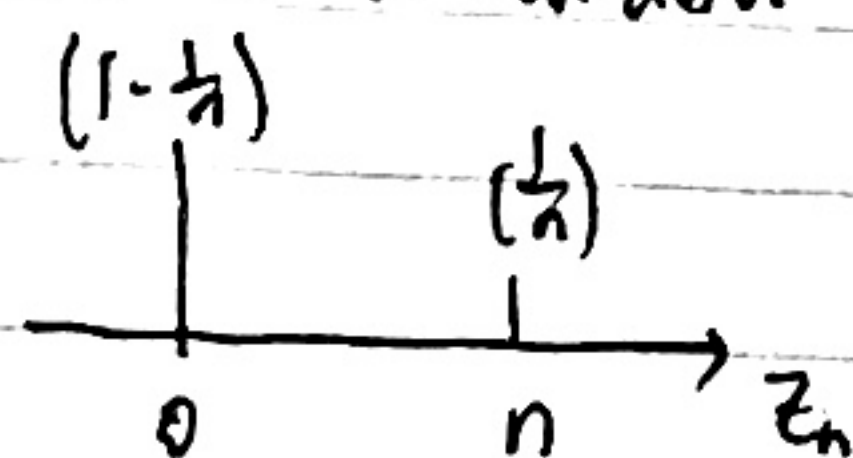
Weak means weak form of convergence

$$\lim_{n \rightarrow \infty} P(|M_n - E[X]| \geq \epsilon) \leq \lim_{n \rightarrow \infty} \frac{\text{Var}[X]}{n\epsilon^2} = 0$$

Since probability is non-negative to begin with, $\lim_{n \rightarrow \infty} P(|M_n - E[X]| \geq \epsilon) = 0$ for any arbitrary $\epsilon > 0$.

Sample Mean converges in probability to the true mean $E[X]$ as the number of trials (n) tends to ∞ .

Consider a random variable Z_n , with the following probability distribution.



$$E[Z_n] = 1$$

$$\text{Var}[Z_n] = \frac{1}{n} (n-1)^2 = \frac{(n-1)^2}{n}$$

$$\text{As } n \rightarrow \infty, \text{Var}[Z_n] = \lim_{n \rightarrow \infty} \frac{(n-1)^2}{n} = \infty$$

Pollster Problem

Want to estimate the true fraction M ($0 < M < 1$) of US voters who believe in some issue.

Determine the minimum number of voters n that we must poll so that we are at least 95% confident that our estimate M_n is within the range $(M - \epsilon, M + \epsilon)$

$$\begin{array}{l} \mu \swarrow X_i = 1 \quad (\text{vote A}) \\ \searrow 1 - \mu \quad X_i = 0 \quad (\text{vote B}) \end{array} \quad M_n = \frac{X_1 + X_2 + \dots + X_n}{n} \quad \text{our estimate } \mu = E[X].$$

$$\text{Want } \underbrace{P[|M_n - \mu| < \epsilon]}_{\text{accuracy}} \geq \underbrace{0.95}_{\text{confidence}}$$

Not in the form of Chebyshev's inequality. Just trivially message.

$$P[|M_n - \mu| < \epsilon] = 1 - P[|M_n - \mu| \geq \epsilon] \geq 0.95$$

$$\Leftrightarrow P[|M_n - \mu| \geq \epsilon] \leq 0.05$$

$$\text{Know } E[X_i] = \mu, \text{Var}[X] = \mu(1-\mu)$$

$$\Rightarrow \text{Var}[M_n] = \frac{\mu(1-\mu)}{n}$$

$$\text{Chebyshev's inequality says } P[|M_n - \mu| \geq \epsilon] \leq \frac{\text{Var}[M_n]}{\epsilon^2} = \frac{\text{Var}[X]}{n\epsilon^2} \leq 0.05$$

$$\text{Set } \epsilon = 0.01, \text{ so want } M_n \in [\mu - 0.01, \mu + 0.01].$$

$$\text{Worst case: maximize } \sigma_x^2 = \mu(1-\mu) \leq \frac{1}{4}.$$

$$\Rightarrow \frac{\sigma_x^2}{n(0.01)^2} \leq 0.05 \Rightarrow \boxed{n \geq 50000}$$

To allocate this large number - compromise on accuracy (ϵ)