

Concentration Inequalities

Markov's inequality:

If X is a non-negative random variable, then for any $t > 0$

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

Proof:

$$Z = \begin{cases} 0 & \text{if } X < t \\ t & \text{if } X \geq t \end{cases}$$

Clearly, with probability one,

$$Z \leq X$$

Take expectation of both

... - ...
sides.

$$E[Z] \leq E[X]$$

The left-hand side is

$$\begin{aligned} E[Z] &= 0 \cdot \Pr[X < t] + t \cdot \Pr[X \geq t] \\ &= t \cdot \Pr[X \geq t] \end{aligned}$$

Thus

$$t \cdot \Pr[X \geq t] \leq E[X]$$

Divide both sides by t . ~~□~~

Chebyshev's inequality

11 \forall ... \leq 1

If X is any random variable then for any $t > 0$,

$$\Pr[|X - E[X]| \geq t] \leq \frac{\text{Var}(X)}{t^2}$$

Proof:

Let

$$Z = (X - E[X])^2$$

Then Z is non-negative.

Markov implies that

$$\Pr[Z \geq t^2] \leq \frac{E[Z]}{t^2} \quad (*)$$

The left-hand side of (*)

$$\Pr[Z \geq t^2] = \Pr[(X - E[X])^2 \geq t^2]$$

$$= \Pr[|X - E[X]| \geq t]$$

The expectation on the right side of (*)
is :

$$\begin{aligned} E[Z] &= E[(X - E(X))^2] \\ &= \text{Var}(X) \end{aligned}$$

a Very common case is
that we have a random
variable that is a sum
or an average of many independent
random variables.

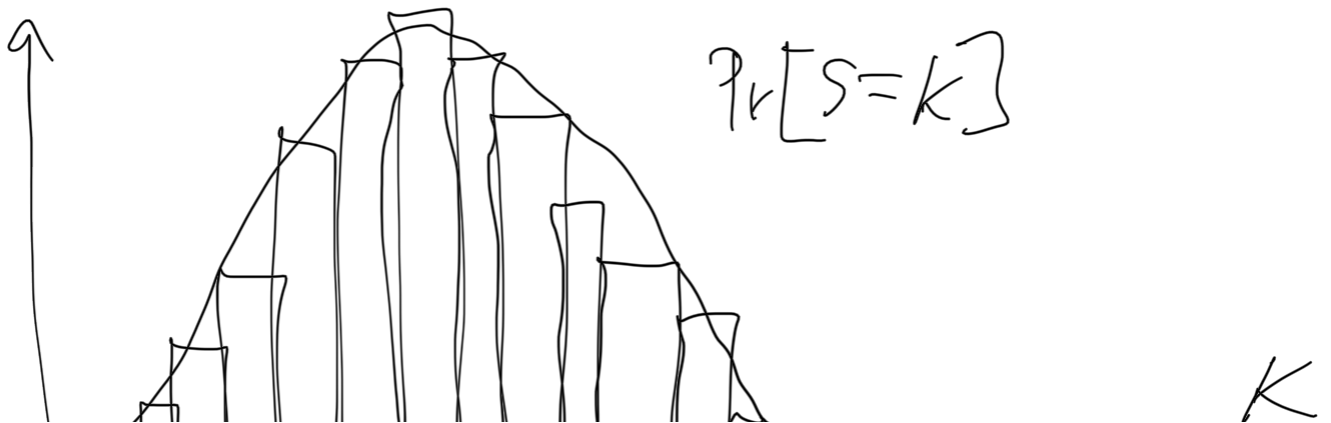
• The simplest example is

X_1, X_2, \dots, X_n are
i.i.d. Bernoulli variables.

$$X_i = \begin{cases} 1 & \text{with probability } p \\ 0 & \text{with probability } 1-p. \end{cases}$$

$$S = \sum_{i=1}^n X_i$$

S has binomial distribution





Let's use Chebyshev's inequality for the binomial case

$$Z = \frac{1}{n} \sum_{i=1}^n X_i$$

$$E[Z] = p$$

$$\text{Var}(Z) = \frac{p(1-p)}{n}$$

$$\Pr\left[\left|\frac{1}{n} \sum_{i=1}^n X_i - p\right| \geq t\right] \leq \frac{p(1-p)}{nt}$$

Central limit theorem

If X_1, X_2, \dots are i.i.d.

with $E[X_i] = \mu$ and $\text{Var}(X_i) = \sigma^2$

— 1

n

(, n, n)

Then for every $t \in \mathbb{R}$,

$$\lim_{n \rightarrow \infty} \Pr \left[\frac{\left(\sum_{i=1}^n x_i \right) - n\mu}{\sigma \sqrt{n}} \geq t \right] = \frac{1}{\sqrt{2\pi}} \int_t^{\infty} e^{-x^2/2} dx$$

$$Z = \frac{\left(\sum_{i=1}^n x_i \right) - n\mu}{\sigma \sqrt{n}}$$

Then Z converges in distribution to $N(0,1)$.

For $t \geq 1$,

$$\int_t^{\infty} e^{-x^2/2} dx \leq e^{-t^2/2}$$

So we expect

$$\Pr \left[\sum_{i=1}^n x_i - n\mu \geq t \right] \leq e^{-t^2/2}$$

$$P\left[\frac{\sum_{i=1}^n x_i}{\sqrt{np(1-p)}} \leq t\right]$$

$$P\left[\left(\frac{1}{n} \sum x_i\right) - p \geq t\right]$$

$$= P\left[\left(\frac{1}{n} \sum x_i - p\right) \sqrt{\frac{n}{p(1-p)}} \geq t \sqrt{\frac{n}{p(1-p)}}\right]$$

$$= P\left[\frac{\sum_{i=1}^n x_i - np}{\sqrt{np(1-p)}} \geq \dots\right]$$

$$\sim e^{-\frac{t^2 n}{2p(1-p)}}$$

This is exponentially better than Chebyshev's inequality.

We need a better bound.

Hoerl's bound ?