

Sums of Random Variables

- We are often interested in the behavior of a sum of random variables $S_n = X_1 + X_2 + \dots + X_n$.
- For instance, the sample mean $\hat{\mu}_x$ or M_n

$$\hat{\mu}_x = M_n = \frac{1}{n} \sum_{i=1}^n X_i$$

is frequently used to estimate the mean from data.

- How many samples n are needed to obtain a good estimate?
- Ideally, we could answer this question precisely using the PMF (or PDF) of a sum of random variables.
- Unfortunately, calculating the exact PMF (or PDF) is **difficult**, even in simple scenarios.

- As a starting point, we will see how to calculate the mean and variance of a sum of random variables

$$S_n = X_1 + X_2 + \dots + X_n.$$

$$\mathbb{E}[S_n] = \mathbb{E}\left[\sum_{i=1}^n X_i\right] = \sum_{i=1}^n \mathbb{E}[X_i] \quad \text{using Linearity of Expectation}$$

$$\text{Var}[S_n] = \mathbb{E}\left[(S_n - \mathbb{E}[S_n])^2\right]$$

$$= \mathbb{E}\left[\left(\sum_{i=1}^n X_i - \sum_{i=1}^n \mathbb{E}[X_i]\right)^2\right]$$

$$= \mathbb{E}\left[\left(\sum_{i=1}^n (X_i - \mathbb{E}[X_i])\right)^2\right]$$

$$\left(\sum_{i=1}^n a_i\right)^2 = \sum_{i=1}^n \sum_{j=1}^n a_i a_j$$

$$= \mathbb{E}\left[\sum_{i=1}^n \sum_{j=1}^n (X_i - \mathbb{E}[X_i])(X_j - \mathbb{E}[X_j])\right]$$

Linearity of Expectation

$$= \sum_{i=1}^n \sum_{j=1}^n \mathbb{E}[(X_i - \mathbb{E}[X_i])(X_j - \mathbb{E}[X_j])]$$

$$= \sum_{i=1}^n \sum_{j=1}^n \text{Cov}[X_i, X_j]$$

Only requires pairwise second-order statistics.

- To further simplify our calculations, we sometimes make additional assumptions on X_1, X_2, \dots, X_n .
- The random variables X_1, X_2, \dots, X_n are said to be **independent and identically distributed (i.i.d.)** if they are independent and have the same underlying marginal PMF, $P_{x_i}(x_i) = P_x(x_i)$, or PDF, $f_{x_i}(x_i) = f_x(x_i)$ for $i = 1, 2, \dots, n$,

Discrete Case: $P_{x_1, x_2, \dots, x_n}(x_1, x_2, \dots, x_n) = P_x(x_1) P_x(x_2) \cdots P_x(x_n) = \prod_{i=1}^n P_x(x_i)$

Continuous Case: $f_{x_1, x_2, \dots, x_n}(x_1, x_2, \dots, x_n) = f_x(x_1) f_x(x_2) \cdots f_x(x_n) = \prod_{i=1}^n f_x(x_i)$

- Ex: X_1, X_2, \dots, X_n are i.i.d. Bernoulli(p), $P_x(x) = p^x (1-p)^{1-x}$

$$P_{x_1, x_2, \dots, x_n}(x_1, x_2, \dots, x_n) = \prod_{i=1}^n p^{x_i} (1-p)^{1-x_i} = p^{\sum_{i=1}^n x_i} (1-p)^{n - \sum_{i=1}^n x_i}$$

- Ex: X_1, X_2, \dots, X_n are i.i.d. Gaussian(μ, σ^2), $f_x(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(x-\mu)^2\right)$

$$f_{x_1, x_2, \dots, x_n}(x_1, x_2, \dots, x_n) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(x_i-\mu)^2\right) = \frac{1}{(2\pi\sigma^2)^{\frac{n}{2}}} \exp\left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i-\mu)^2\right)$$

- The mean and variance of a sum $S_n = \sum_{i=1}^n X_i$ of i.i.d. random variables X_1, \dots, X_n are

$$\mathbb{E}[S_n] = \mathbb{E}\left[\sum_{i=1}^n X_i\right] = n \mathbb{E}[X] \quad \text{Var}[S_n] = \text{Var}\left[\sum_{i=1}^n X_i\right] = n \text{Var}[X]$$

Compute using the marginal PMF $P_X(x)$ or PDF $f_X(x)$.

Why? $\mathbb{E}\left[\sum_{i=1}^n X_i\right] = \sum_{i=1}^n \mathbb{E}[X_i]$
Identically Distributed
 $= \sum_{i=1}^n \mathbb{E}[X]$
 $= n \mathbb{E}[X]$

$\text{Var}\left[\sum_{i=1}^n X_i\right] = \sum_{i=1}^n \sum_{j=1}^n \text{Cov}[X_i, X_j]$
 $= \sum_{i=1}^n (\text{Cov}[X_i, X_i] + \sum_{j \neq i} \text{Cov}[X_i, X_j])$
Cov[X,X] = Var[X]
 $= \sum_{i=1}^n \text{Var}[X_i] = n \text{Var}[X]$
independent
Identically Distributed

- The mean and variance of the sample mean $M_n = \frac{1}{n} \sum_{i=1}^n X_i$ of i.i.d. random variables X_1, \dots, X_n are

$$\mathbb{E}[M_n] = \mathbb{E}\left[\frac{1}{n} \sum_{i=1}^n X_i\right] = \mathbb{E}[X] \quad \text{Var}[M_n] = \text{Var}\left[\frac{1}{n} \sum_{i=1}^n X_i\right] = \frac{1}{n} \text{Var}[X]$$

Why? $\mathbb{E}\left[\frac{1}{n} \sum_{i=1}^n X_i\right] = \frac{1}{n} \mathbb{E}\left[\sum_{i=1}^n X_i\right]$
 $= \frac{1}{n} \cdot n \mathbb{E}[X]$

$\text{Var}\left[\frac{1}{n} \sum_{i=1}^n X_i\right] = \frac{1}{n^2} \text{Var}\left[\sum_{i=1}^n X_i\right]$
Var[aX] = a^2 Var[X]
 $= \frac{1}{n^2} \cdot n \text{Var}[X]$
Goes to 0 as n increases.

- The **sample variance** $\hat{\theta}_x^2$ or V_n is frequently used to estimate the variance from data

$$\hat{\theta}_x^2 = V_n = \frac{1}{n-1} \sum_{i=1}^n (x_i - M_n)^2 \quad \text{where } M_n = \frac{1}{n} \sum_{i=1}^n x_i$$

- Why do we multiply by $\frac{1}{n-1}$ instead of $\frac{1}{n}$?

→ Consider $U_n = \frac{1}{n} \sum_{i=1}^n (x_i - M_n)^2 \stackrel{\text{can be shown}}{=} \frac{1}{n} \sum_{i=1}^n x_i^2 - M_n^2$.

- Calculate $\mathbb{E}[U_n]$. Is it $\text{Var}[X]$?

$$\begin{aligned} \mathbb{E}[U_n] &= \mathbb{E}\left[\frac{1}{n} \sum_{i=1}^n x_i^2 - M_n^2\right] \\ &\stackrel{\text{Linearity of Expectation}}{=} \frac{1}{n} \sum_{i=1}^n \mathbb{E}[x_i^2] - \mathbb{E}[M_n^2] \\ &= \frac{1}{n} \sum_{i=1}^n (\text{Var}[X] + (\mathbb{E}[X])^2) - \text{Var}[M_n] - (\mathbb{E}[M_n])^2 \\ &= \text{Var}[X] + \cancel{(\mathbb{E}[X])^2} - \frac{1}{n} \text{Var}[X] - \cancel{(\mathbb{E}[X])^2} \\ &= \frac{n-1}{n} \text{Var}[X] \quad \text{This is a biased estimator.} \end{aligned}$$

- Using a $\frac{1}{n-1}$ factor corrects this issue.