

2/12 LECTURE 19: E, VAR & SUMMARY OF DATA.

X r.v. with pmf p_X . $\mu_X = E[X] = \sum_x x p_X(x)$ ($\mu = \mu_X$)
 $E[f(X)] = \sum_x f(x) p_X(x)$ $f: \mathbb{R} \rightarrow \mathbb{R}$ [Expectation/mean]
 $VAR[X] := E[(X - \mu)^2]$ (Variance of X).

$$= E[X^2 - 2\mu X + \mu^2]$$

$$= E[X^2] - \mu^2 \quad (\text{linearity})$$

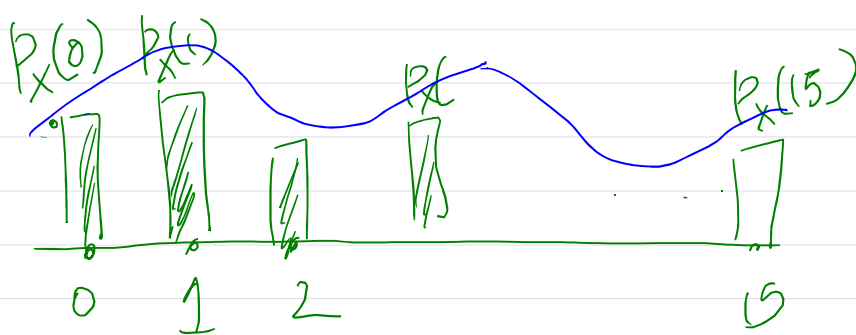
$$VAR[X] \geq 0$$

Alternative $E[(X - \mu)^2] = \sum_x (x - \mu)^2 p_X(x)$ (Take $f(x) = (x - \mu)^2$)

$$= \sum_x x^2 p_X(x) - \mu^2 = E[X^2] - \mu^2.$$

Ex 19.1 Assume X r.v. with pmf $p_X(x) = 0$ if $x \notin \{0, \dots, 15\}$.
 You are given $p_X(0), p_X(1), \dots, p_X(15)$.

Pictorially
 p_X is an histogram.



Can you summarise the histogram / r.v.?

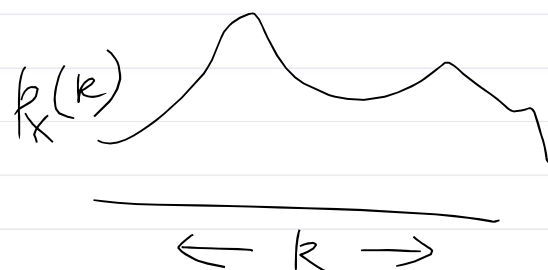
First summary is mean — $\mu_X = E[X] = \sum_x x p_X(x)$

Ex 19.2 Suppose I have scores of students in Quiz 1.

Define $p_X(k) = \frac{\# \text{ of scores} = k}{\# \text{ of students}}$. $\sum_{k=0}^{15} p_X(k) = 1$, $p_X(x) \geq 0$.

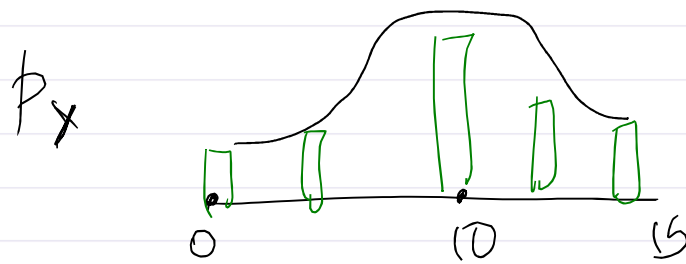
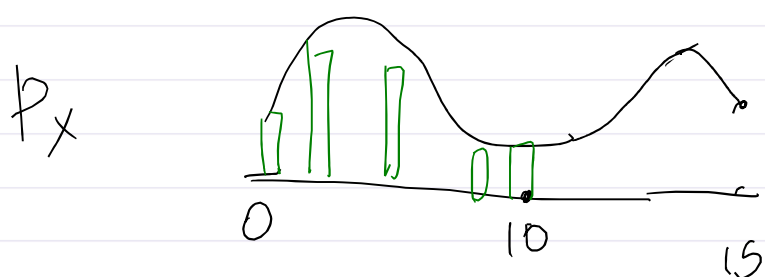
$\Rightarrow p_X$ is a pmf of some r.v. X .

(mean) $\mu_X = \sum_x x p_X(x)$ ($= 9.6\dots$)



Qn*: Find explicitly (Ω, \mathcal{P}) & $X: \Omega \rightarrow \mathbb{R}$ \exists $P(X=k) = p_X(k)$.

Qn: Suppose I have two pmfs p_x & p_y



Assume $\sum_x x p_x(x) = \mu_x = \mu_y = \sum_y y p_y(y)$.

$\text{VAR}[X] = \sum_x (x - \mu)^2 p_x(x) \Rightarrow$ larger variance \Rightarrow larger "disparity"

LEMMA 19.3 (i) $\text{VAR}[X] = 0$ iff $p_x(\mu) = 1$ (i.e., $X \equiv \mu$)

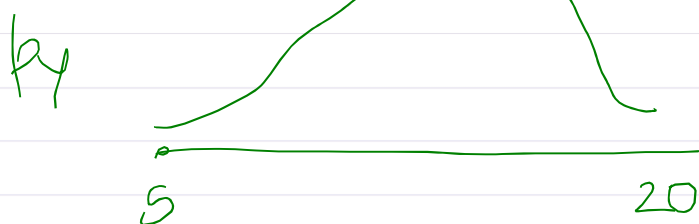
(2) $\text{VAR}[aX + b] = a^2 \text{VAR}[X] \quad \forall a, b \in \mathbb{R}$

(3) $\mu_{aX+b} = a\mu_x + b$. (recall)

[why not $\sum_x (x - \mu) p_x(x)$??]

Ex 19.4 (Again ex. of Quiz series). Suppose I add 5 points to all.

Then the new r.v. $Y = X + 5$. $p_y(x) = p_x(x - 5)$, $x \in \mathbb{R}$

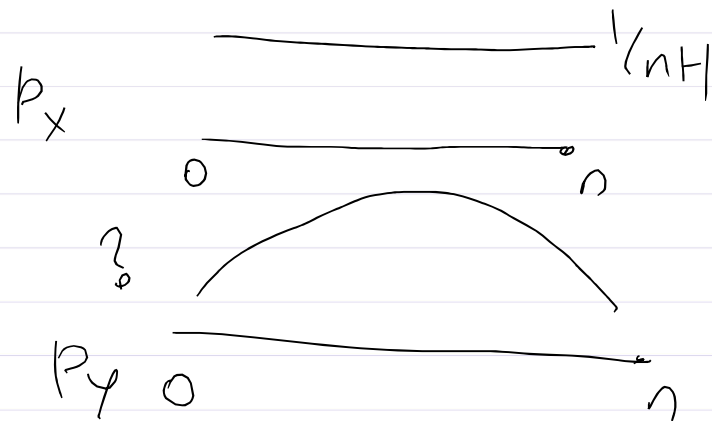


$\mu_y = \mu_x + 5$, $\text{VAR}[Y] = \text{VAR}[X]$

Ex 19.5 (A) $X \stackrel{d}{=} \text{Unif}\{0, 1, 2, \dots, n\}$. $\mathbb{E}[X] = \frac{n}{2}$

$Y \stackrel{d}{=} \text{Bin}(n, 1/2)$

$\mathbb{E}[Y] = \frac{n}{2}$



Compute $\text{VAR}(X)$, $\text{VAR}(Y)$.

[& what can you say about disparity?]

Further summaries - $E[X]$, $VAR[X]$, $E[X^k]$, $k \geq 1$.

For eg. if $X \stackrel{d}{=} \text{Ber}(p)$ then $X^k \stackrel{d}{=} \text{Ber}(p) \quad \forall k \geq 1$.

$$[P_X(1) = p = 1 - P_X(0)] \quad [\text{i.e., } P_{X^k}(1) = p = 1 - P_{X^k}(0)]$$

$$\Rightarrow E[X^k] = p \quad \forall k \geq 1.$$

EX
19.7

Let X be a r.v. $f: \mathbb{R} \rightarrow \mathbb{R}$ & define $Y = f(X)$.

Then Y is a r.v. with pmf $P_Y(\cdot) = ??$ (write it in terms of P_X).

EX
19.8

Sampling with replacement. Total population = N .

There are two types. # of type 1 = $Np \in \mathbb{N}$, $p \in [0, 1]$

of type 2 = $N(1-p) \in \mathbb{N}$.

Suppose you choose a sample of size n from the population at random & with replacement.

Let $X = \#$ of samples of type 1 in the n randomly chosen samples.

$$\begin{aligned} P_X(k) = P(X=k) &= \underset{(\text{check})}{\binom{n}{k}} \left(\frac{Np}{N}\right)^k \left(\frac{N(1-p)}{N}\right)^{n-k} \\ &= \binom{n}{k} p^k (1-p)^{n-k}. \end{aligned}$$

So X is a $\text{Bin}(n, p)$ r.v.

Qn. How to find p ?

Define $Y = \frac{X}{n}$. $P_Y\left(\frac{k}{n}\right) = \binom{n}{k} p^k (1-p)^{n-k} \quad k=0, \dots, n.$

$$E[Y] = \underset{(\text{linearity})^n}{E[X]} = p. \quad (\text{mean of } Y)$$

$$\underset{(\text{by LEMMA 19.6})}{VAR[Y]} = \underset{(A)}{\frac{VAR[X]}{n^2}} = \frac{np(1-p)}{n^2} = \frac{p(1-p)}{n}$$

$$\text{VAR}[Y] = \sum_x (y - p)^2 p_y(y) = \frac{p(1-p)}{n} \leq \frac{1}{4n}$$

$$(p(1-p) \leq \frac{1}{4})$$

$$\text{VAR}[Y] \leq \frac{1}{4n}$$

$$\text{If } n = 2500, \text{ VAR}[Y] \leq 10^{-4}$$

Intuitively, " $|Y - p| \leq 10^{-4}$ ";
 If p is 10^{-6} then this is bad; if " p is small or large", $\frac{p(1-p)}{n}$ is very small

Ex
19.4
/ (A)

Repeat the above exercise without replacement.

Find p_y , $E[Y]$ & $\text{VAR}[Y]$.

$$\mu_y = E[Y]; \quad \sigma_y = \text{SD}[Y] := \sqrt{\text{VAR}[Y]}$$

LEMMA
19.10

(Chebyshev's Ineq.)

X r.v.

$$P(|X - \mu| \geq t) \leq \frac{\sigma^2}{t^2} \quad \forall t > 0.$$

Proof:

$$|X - \mu| \geq t \Rightarrow (X - \mu)^2 \geq t^2 \quad \forall t > 0$$

$$\text{So } \{ |X - \mu| \geq t \} \subseteq \{ (X - \mu)^2 \geq t^2 \}$$

$$\Rightarrow P(|X - \mu| \geq t) \leq P((X - \mu)^2 \geq t^2)$$

$$\begin{aligned} & \text{(M. Ineq.)} \\ & \text{as } (X - \mu)^2 \geq 0 \end{aligned} \quad \leq \frac{E[(X - \mu)^2]}{t^2} = \frac{\sigma_X^2}{t^2}$$

So Chebyshev's ineq. \Rightarrow

$$P(|X - \mu| \geq k\sigma) \leq \frac{1}{k^2} \quad \forall k \geq 1$$

Apply to Eg. of sampling with replacement.

$$Y = \frac{X}{n}; \quad \mu_Y = p; \quad \sigma_Y = \sqrt{\frac{p(1-p)}{n}}$$

$$P(|Y - p| \geq K \sigma_Y) \leq \frac{1}{K^2}$$

"Prob. of Error in $|Y - p|$ is larger than $\frac{K \sqrt{p(1-p)}}{\sqrt{n}}$ is at most $\frac{1}{K^2}$!"

If $K=10$, then $|Y - p| \leq \frac{K \sqrt{p(1-p)}}{\sqrt{n}}$ with 99% accuracy.

Ex. Do this for sampling without replacement.

Find t s.t. $P(|Y - a| \geq t) \leq \frac{1}{100}$. ??