

Lecture 8 (Ch. 2)

Dealing with ambiguity (e.g., how wide is a curve)
Dealing with precise syntax (e.g., words and their order matters)
Population (the truth)
Sample (our observed data)
Random Variable
Types of data (well-defined, but ambiguous)
Histogram (its interpretation and uses)
Probability from histogram
Distribution (its interpretation and uses)
Probability from distribution
Named distributions
The random variable "template" associated with each distribution
Standardization (for Normal and other distributions)
Percentile/quantile (for sample and/or dist)
boxplots (uses and interpretation)
Derivation and application of Binomial and Poisson.

Summary of
what you have
learned, so far.

Time to quantify some of our qualitative ideas.

In prev. chapters we played with histograms of sample/data and distributions of (random) variables (cont. and discrete/categ.). Hists and dists are pillars of statistics. Hists describe the data/sample, while dists describe the population. One question we often ask is this: how likely is it that my data/hist came from, say, a normal dist with params $\mu = \dots$, $\sigma = \dots$?

One way to compare histograms with distributions is in terms of their summary measures. For example, we can compare the "location" of the Normal distr. (μ) with the "location" of the histogram.

Or the "width" of the former (σ) with the width of the histogram.

So, we need some measure of location and width for both histograms and distributions (for any dist, not just Normal).

	hist/sample	dist/pop
measure of location	?	?
measure of width	?	?

The first comparison between sample and pop. will happen soon (when we do qq-plots), and then more fully in the 2nd half of 390.

Measures of location for hist/sample:

- sample mean : $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$

The x for the i^{th} case

pros/cons

- sample median : $\tilde{x} = \text{middle of the ordered data.}$

data

Measures of spread for hist./sample:

- sample Range

standard deviation
(same units as \bar{x})

pros/cons

- sample variance $= \underbrace{S^2}_{\text{deviation.}} = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$

$S \sim$ "average" (typical) deviation.

$$\sum_{i=1}^n (x_i - \bar{x}) = \sum_{i=1}^n x_i - \sum_{i=1}^n \bar{x} = n\bar{x} - \bar{x} \underbrace{\sum_{i=1}^n 1}_n = 0.$$

Example : $x = c(1, 3, 8)$ cm

$$S = \sqrt{13} \text{ cm}$$

$$\bar{x} = \frac{1}{3} (1 + 3 + 8) = 4 \text{ cm}$$

$$S^2 = \frac{1}{3-1} [(1-4)^2 + (3-4)^2 + (8-4)^2] = \frac{1}{2} (9 + 1 + 16) = 13 \text{ cm}^2$$

Again : A lot of statistics is about explaining/understanding This variance. Recall, if There were no variance/change, we wouldn't say that we even have any data.

In short, we will use the following summary measures for location and spread of data:

Sample mean: $\bar{x} = \frac{1}{n} \sum_i x_i$

sample variance: $s^2 = \frac{1}{n-1} \sum_i (x_i - \bar{x})^2$ ← Because $\sum_i (x_i - \bar{x}) = 0$
"funny Average"

Then s will be another measure of spread, and it's even better than s^2 , because s has the same physical dimension as x itself. So, we can write things like $\bar{x} \pm s$ as a way of summarizing a histogram.

Important: Interpretation of \bar{x} is typical x
" " " s_x " typical deviation of x .
or s^2

In some problems where the $\frac{1}{n-1}$ is not important, one focuses on $S_{xx} \equiv \sum_i (x_i - \bar{x})^2$, i.e. just the numerator of s^2 .

Finally, note that all of these measures have the word "Sample," reminding you that they pertain to sample/data, not pop./distr.

(FYI)

For the mathematically-inclined: If you think of x_i as the components of an n -vector, then after you "center" each component (i.e. subtract \bar{x}), s is proportional to the magnitude of that vector.

Another way of computing s^2 .

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$$

$$= \frac{1}{n-1} \sum_{i=1}^n (x_i^2 - 2x_i\bar{x} + \bar{x}^2)$$

$$= \frac{1}{n-1} \left[\sum_{i=1}^n x_i^2 - 2 \underbrace{\bar{x}}_{\substack{\text{circled} \\ n\bar{x}}} \underbrace{\sum_{i=1}^n x_i}_{n\bar{x}} + (\bar{x})^2 \sum_{i=1}^n 1 \right]$$

$$= \frac{1}{n-1} \left[\sum_{i=1}^n x_i^2 - 2n(\bar{x})^2 + n(\bar{x})^2 \right]$$

$$= \frac{1}{n-1} \left[\sum_{i=1}^n x_i^2 - n(\bar{x})^2 \right]$$

$$= \frac{1}{n-1} \left[n \underbrace{\left(\frac{1}{n} \sum x_i^2 \right)}_{\overline{x^2}} - n(\bar{x})^2 \right] = \frac{n}{n-1} \left[\overline{x^2} - (\bar{x})^2 \right]$$

In summary:

$$s^2 = \frac{1}{n-1} \sum_i (x_i - \bar{x})^2$$

"Defining formula"

$$s^2 = \frac{n}{n-1} \left[\overline{x^2} - (\bar{x})^2 \right]$$

"Computational formula"
sometimes more useful

always faster (1 vs. 2 loops)
Not too important.

Example

$$x = c(1, 3, 8) \rightarrow x^2 = c(1, 9, 64) \rightarrow \overline{x^2} = \frac{74}{3}$$

$$s^2 = \frac{3}{2} \left[\frac{74}{3} - 16 \right] = \frac{3}{2} \frac{74-48}{3} = \frac{26}{2} = 13 \quad (\text{same as above}).$$

Keep the "big picture" in mind: We are looking for

		sample/hist.	pop./distr.
measure of	location	<p>Sample mean</p> $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$ <p>$\bar{x} \sim$ "typical x"</p>	?
	Spread	<p>Sample variance (s^2)</p> <p>Sample std. dev. (s)</p> $s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$ $= \frac{n}{n-1} (\bar{x}^2 - \bar{x}^2)$ <p>$s \sim$ "typical deviation in x"</p>	?

s, s^2 , and S_{xx} are all measures of spread.

Now, we need to come-up with corresponding things in the pop.

So, switch to distributions ($p(x), f(x)$). No Data/Sample!

different symbols
for the same thing.

1) Distribution mean
(or Expected Value)

$$= \mu_x = E[x] = \begin{cases} \sum_x x p(x) \\ \int x f(x) dx \end{cases}$$

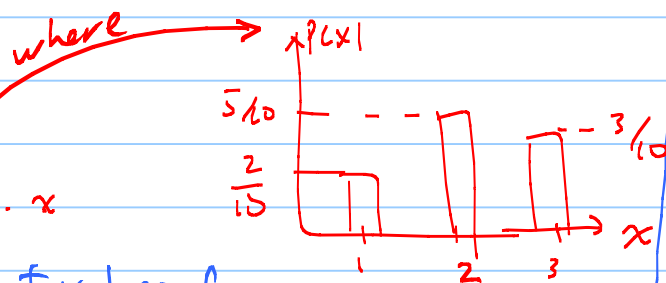
Motivation: Even though we are now in the realm of math ($p(x), f(x)$) not data, just to motivate the defn of $E[x]$, consider the following "population" $\{3, 2, 2, 1, 3, 2, 3, 1, 2, 2\}$.

$$\text{"avg."} = \frac{1}{10} [3 + 2 + 2 + \dots]$$

$$= \frac{1}{10} [3(3) + 5(2) + 2(1)]$$

$$= \frac{3}{10}(3) + \frac{5}{10}(2) + \frac{2}{10}(1) = \sum_{\text{distinct } x} p(x) \cdot x$$

distinct values of x



Compare :

Sample mean: $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$, \sim typical x (in dist. / pop.)

distr. mean (Expected Value): $\mu_x = E[x] = \sum_x x p(x), \int_{-\infty}^{\infty} x f(x) dx \sim$ typical deviation in x .

the book drops the x on μ_x , but then μ can be confused with the parameter of the Normal distr.

As examples, let's find the mean of Normal & Poisson, below.

$E[x]$ does not mean that E is a function of x . In fact, E is a \sum_x or an $\int_{-\infty}^{\infty} dx$, and so it is not a function of x .

$E[x]$ simply means that you need $p(x)$ or $f(x)$ to find it.

See examples, next lect; There is no x -dependence in $\mu_x = E[x]$.

FYI

The mean of $N(\mu, \sigma)$:

$$\begin{aligned}\mu_x \text{ (or } E[x]) &= \int_{-\infty}^{\infty} x f(x) dx = \frac{1}{\sqrt{2\pi}\sigma} \int_{-\infty}^{\infty} x e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2} dx \\ &= \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} (\mu + \sigma z) e^{-\frac{1}{2}z^2} dz \quad \left\{ \begin{array}{l} \text{Change of Var.} \\ \frac{x-\mu}{\sigma} = z \\ \frac{dx}{\sigma} = dz \end{array} \right. \\ &= \underbrace{\mu \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}z^2} dz}_{=1} + \underbrace{\sigma \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}} z e^{-\frac{1}{2}z^2} dz}_{=0}\end{aligned}$$

$\mu_x = \mu$

Now, you can see why the μ parameter of the normal distr. is called its mean.

Either look-up in Table of integrals, or note that z is odd, while the \int goes from $-\infty$ to $+\infty$.

You can also see why the subscript on μ_x is important!

Mean of Poisson (λ):

$$\mu_x = E[x] = \sum_{x=0}^{\infty} x \frac{e^{-\lambda} \lambda^x}{x!} = \dots = \lambda \sum_{x=0}^{\infty} \frac{e^{-\lambda} \lambda^{x-1}}{(x-1)!} = \lambda \quad \left\{ \begin{array}{l} \text{Change of Var.} \\ 1 = \sum_{x=0}^{\infty} p(x) \end{array} \right.$$

Now you can see why the λ param. of Poisson is called its mean.

Warning: Don't confuse \bar{x} , μ_x

Sample mean:

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$

distr. mean

→ For $N(\mu, \sigma)$, $\mu_x = \mu$

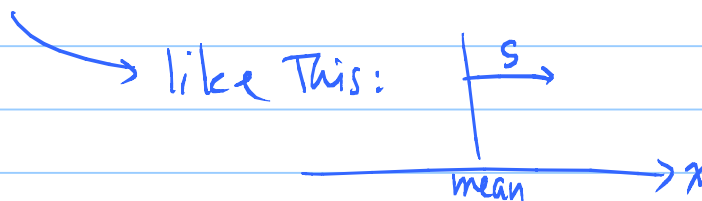
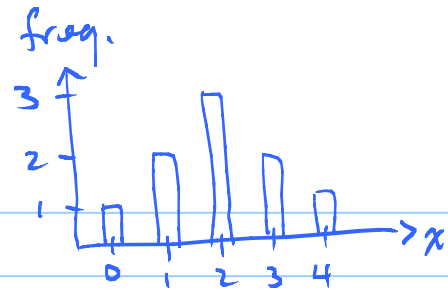
→ For Poisson (λ), $\mu_x = \lambda$

→ For other distributions, μ_x will be other things (next!)

hw lect8_1

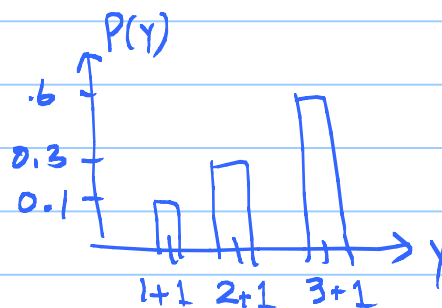
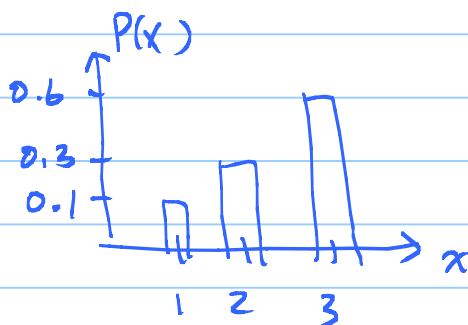
Consider the adjacent histogram.

- Compute/find the sample mean \bar{x} . Show work.
- Find The sample std. dev. s , using the defining formula.
- Draw the mean and the sample std. dev. on the histogram.



hw_lect8-2

- For the bottom/left distribution, find the distr. mean $\mu_x = E[x]$.
- For the bottom/right distribution, find the distr. mean $\mu_y = E[y]$. Note that $y = x + 1$.



hw_lect8_3:

To understand our formulas, it is often useful to see what they say if our data are just a bunch of 0's and 1's. So, suppose our data consists of n_0 zeros and n_1 ones. Note: $n = n_0 + n_1$.

- Find the sample mean
- Show that the sample variance is $(n_0 * n_1) / (n(n-1))$

hw-lect8-4

For the exponential distribution with parameter λ , find the mean. Hint: You may use this integral:

$$\int_0^{\infty} y e^{-\lambda y} dy = 1.$$