# Lecture 9 (Ch. 2)

So far, we have introduced the following measures of location and spread of

| Histogram/Sample/Data | Distribution/population |

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

↑ sample mean

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

↑ Sample Variance, $s = $ sample std. dev.

*Important*

~ typical deviation

$$\mu_x = E[x] = \sum_x x\, p(x) \;,\; \int_{-\infty}^{\infty} x f(x)\, dx$$

dist./pop. mean

?

next.

---

Now, a popular measure of distr. spread is (distr.) Variance:

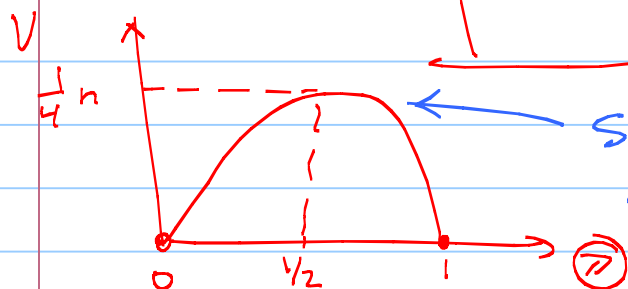$$V[x] = \sigma_x^2 = \begin{cases} \sum_{\forall x} (x - E[x])^2 p(x) \\[2mm] \int_{-\infty}^{\infty} (x - E[x])^2 f(x)\, dx \end{cases}$$

$E[x] = \sum_x x\, p(x)$

The book drops the $x$, but then this $\sigma$ can be confused with the $\sigma$ of the Normal distr.

As with sample std. dev., $\sigma_x = $ distr. std. dev. ← Same units as $x$.

Let's compute $V[x] = \sigma_x^2$ for some of our special distr:

$$\text{Binomial}(n,\pi): \boxed{V[x] \equiv \sigma_x^2} = \overset{E[x]=n\pi}{\underset{\text{binomial}}{\sum (x-n\pi)^2 \, p(x)}}$$

$$= \cdots = n\pi(1-\pi)$$



$$V$$
$$\tfrac{1}{4}n$$

So the largest possible standard deviation is

$$\sqrt{\tfrac{1}{4}n} = \tfrac{1}{2}\sqrt{n} \, , \text{ at } \pi=\tfrac{1}{2}.$$

I.e. If I'm trying to estimate the number of heads out of $n$ tosses, the largest "uncertainty" is $\tfrac{1}{2}\sqrt{n}$ ( $\boxed{\text{at } \pi=\tfrac{1}{2}}$ ).

---

$$\underline{\text{Poisson}(\lambda)}: V[x]^{=\sigma_x^2} = \sum (x-\lambda)^2 \, \overset{\text{poisson}}{p(x)} = \cdots = \lambda$$

$$\text{Recall } E[x]=\lambda \longleftarrow \text{same}$$

---

$$\underline{\text{Normal}(\mu,\sigma)}: V[x]^{=\sigma_x^2} = \int (x-\mu)^2 \, \overset{\text{Normal}}{f(x)} \, dx = \cdots = \sigma^2$$
$$E[x]=\mu$$

which is why the param. $\sigma^2$ is called variance.

---

Q1: Let $f(x)=2x, \ 0<x<1$. In the last qz, we found $E[x]=\mu_x=\tfrac{2}{3}$. Now, find $V[x]\equiv\sigma_x^2$?

A) 1     B) $\tfrac{1}{2}$     C) $\boxed{\tfrac{1}{18}}$     D) This $f(x)$ has no variance.

$$V[x]\equiv\sigma_x^2 = \int_{-\infty}^{\infty}(x-\mu_x)^2 f(x)\,dx = \int_0^1 \left(x-\tfrac{2}{3}\right)^2 2x\,dx = 2\int_0^1\left(x^3-\tfrac{4}{3}x^2+\tfrac{4}{9}x\right)dx$$

$$= 2\left[\tfrac{1}{4}x^4 - \tfrac{4}{3}\tfrac{1}{3}x^3 + \tfrac{4}{9}\tfrac{1}{2}x^2\right]_0^1 = 2\left[\tfrac{1}{4}-\tfrac{4}{9}+\tfrac{2}{9}\right] = 2\left(\tfrac{1}{4}-\tfrac{2}{9}\right) = \tfrac{2}{36} = \tfrac{1}{18}$$

By now, you should be familiar with the meaning of

histograms vs. distributions

Sample mean $\bar{x}$ vs. distv. mean $E[x] = \mu_x$

" Variance $s^2$ vs. " Variance $V[x] = \sigma_x^2$

" std. dev. $s$ vs. " std. dev. $\sigma_x$

Finally, given that we can compute all of the above quantities, you can then compute the proportion of times $x$ is expected to be within some std. dev. of its mean, for ANY distv. $\uparrow$ 1, 1.96, 2, ---

For examples, for the normal dist. we can now say that 68% of $x$'s fall within 1 std. dev. of the mean.

But now we can say things like that for any distv. even skewed ones:



e.g. Poisson ---

$\mu_x - \sigma_x \quad \mu_x \quad \mu_x + \sigma_x$

Computing areas like this will eventually enable us to provide some measure of confidence when we try to estimate a population parameter, later,

Single summary of
histogram location

Sample mean:

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

~ typical x/obs.

Single summary of
histogram spread

Sample variance:

recall computational
formula, too. →

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

Sample std. dev. = s.
~ typical deviation/spread

Recall why $\frac{1}{n} \sum (x_i - \bar{x})'$ will not do!

Single Summary of
distribution/population location

dist./pop mean, or $E[x]$

$$\mu_x \equiv E[x] = \sum_x x \, p(x) \, , \, \int_{-\infty}^{\infty} x f(x) \, dx$$

Single summary of
dist/pop. spread

dist/pop. Var. or $V[X]$

$$\sigma_x^2 \equiv V[x] = \sum_x (x - \mu_x)^2 p(x)$$

$$\int_{-\infty}^{\infty} (x - \mu_x)^2 f(x) \, dx$$

Eg. binomial $(n, \pi)$: $\mu_x = n\pi$ $\quad\quad \sigma_x^2 = n\pi(1-\pi)$

poisson $(\lambda)$: $\mu_x = \lambda$ $\quad\quad \sigma_x^2 = \lambda$

Normal $(\mu, \sigma)$: $\mu_x = \mu$ $\quad\quad \sigma_x^2 = \sigma^2$

uniform $(a, b)$: $\mu_x = \frac{a+b}{2}$ $\quad\quad \sigma_x^2 = \frac{(b-a)^2}{12}$

Exponential $(\lambda)$: $\mu_x = \frac{1}{\lambda}$ $\quad\quad \sigma_x^2 = \left(\frac{1}{\lambda}\right)^2$
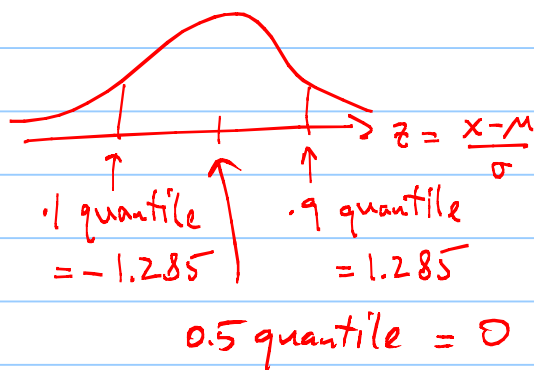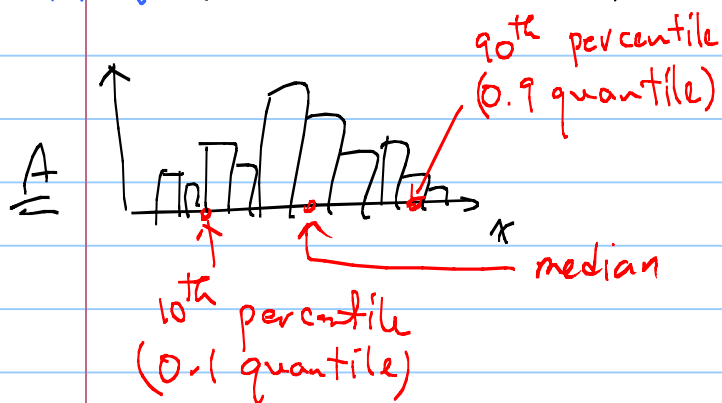
# One more Thing in Ch 2.

This business of estimating pop. parameters refers to
← dist.

any parameter. Specifically, $\bar{x}$ and $s$ provide
point estimates for $\mu, \sigma$, respectively, of The Normal
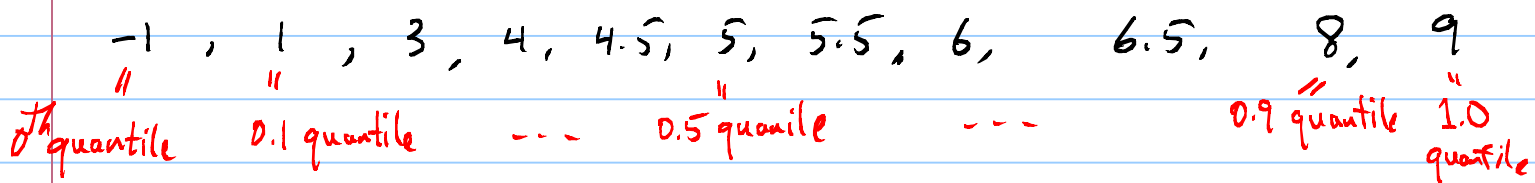distr. if The data come from a Normal distr. to begin with.

Q: But how do we know if our data come from a Normal?

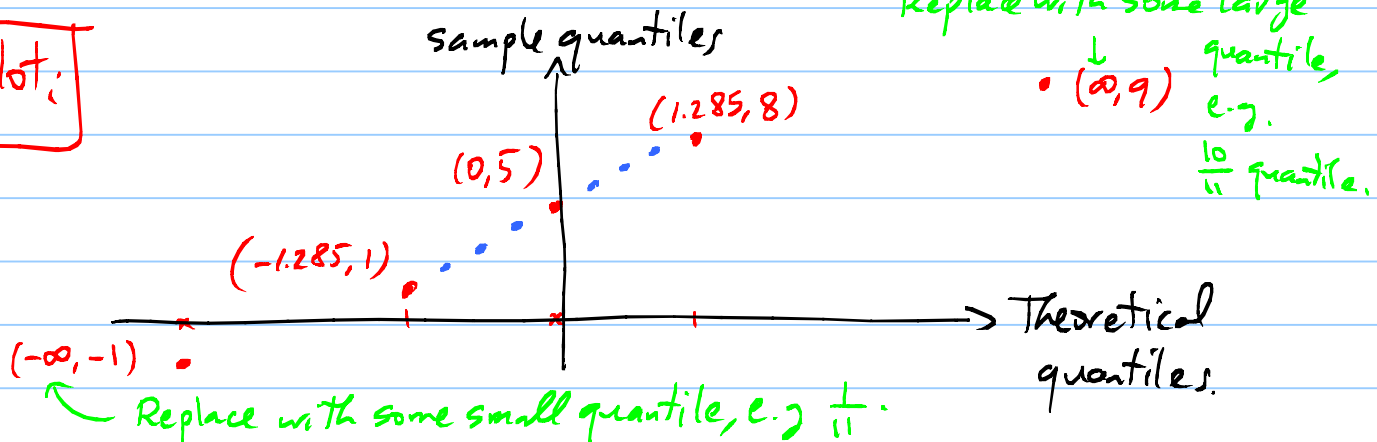Easier Q: how do we know if our data come from std. Normal?

$\frac{f}{2}$

90th percentile
(0.9 quantile)

median

10th percentile
(0.1 quantile)

$z = \frac{x - \mu}{\sigma}$

.1 quantile
$= -1.285$

.9 quantile
$= 1.285$

0.5 quantile $= 0$.

## Example: (Very Crude!)
Data, sorted:

$$-1, \quad 1, \quad 3, \quad 4, \quad 4.5, \quad 5, \quad 5.5, \quad 6, \quad 6.5, \quad 8, \quad 9$$

$0^{th}$ quantile    0.1 quantile    - - -    0.5 quantile    - - -    0.9 quantile    1.0 quantile

→ Theoretical quantiles:

$-\infty$    $-1.285$    - - -    $0$    - - -    $+1.285$    $\infty$

Replace with some large
quantile,
e.g.
$\cdot (\infty, 9)$
$\frac{10}{11}$ quantile.

## qq plot:

sample quantiles

$(1.285, 8)$

$(0, 5)$

$(-1.285, 1)$

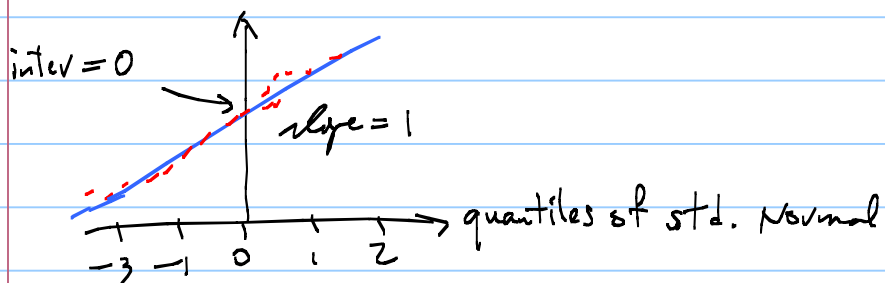$(-\infty, -1)$

→ Theoretical quantiles.
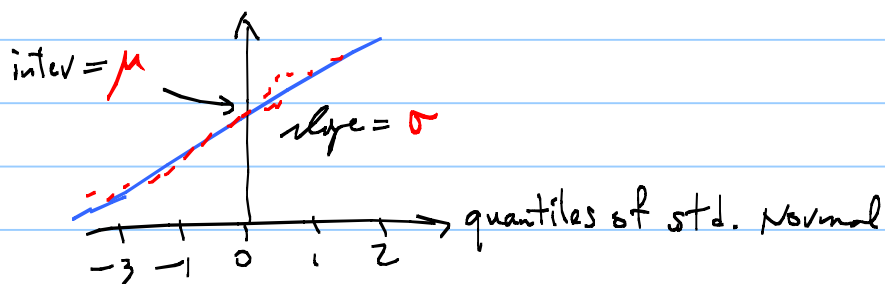
Replace with some small quantile, e.g. $\frac{1}{11}$.

If The histogram is consistent with a std. Normal, Then
the quantiles/percentiles of data should be equal/comparable
to those of The distr.. Then The qq plot should be a straight
diagonal line ( intercept=0, slope=1).

intev = 0

slope = 1

quantiles of std. Normal

-3  -1  0  1  2

If The data are |not| from std. normal, but from |$N(\mu, \sigma)$|,
The only thing That changes is That The slope becomes $\sigma$,
and The intercept becomes $\mu$. NoT too obvious, but pf. in book.

intev = $\mu$

slope = $\sigma$

quantiles of std. Normal
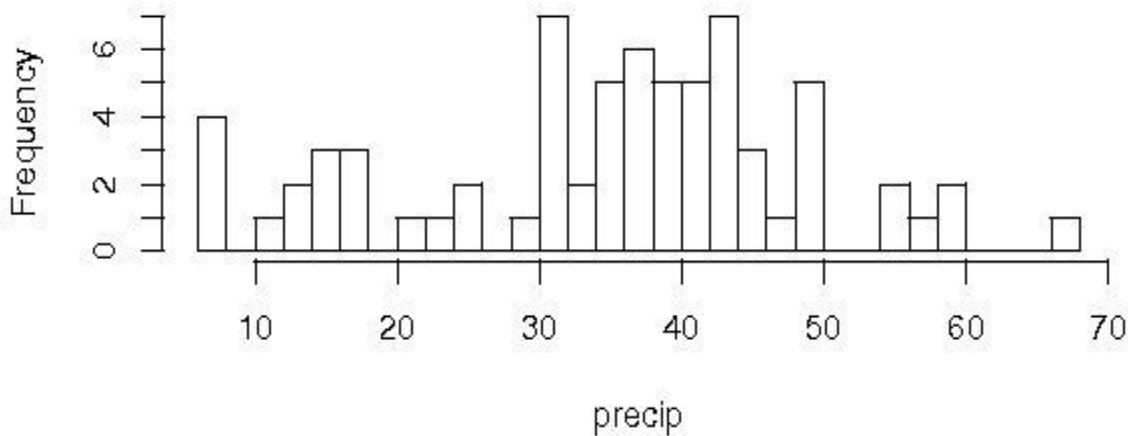
-3  -1  0  1  2

In R:  qqnorm(x)    where    x  is The vector of data.

## (Example)

From the histogram, it's hard to tell if the data come from a normal dist., especially because hists depend on bin size.
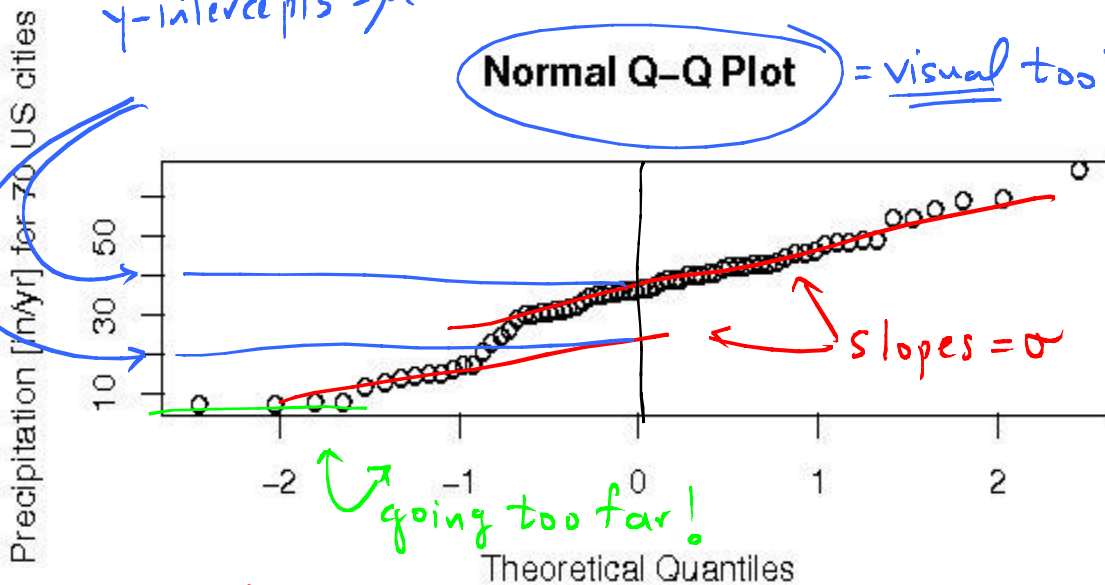
**Histogram of precip**



y-intercepts = $\mu$

**Normal Q-Q Plot** = visual tool.



slopes = $\sigma$

going too far!

The plot looks linear, mostly!

So, data are consistent with a Normal.

In fact, it looks like 2 different normals (Bimodal) with diff $\mu$'s, same $\sigma$ (slope).

(hw-lect11-1) For The uniform distr. between $a, b$, show $V(x) = \frac{1}{12}(b-a)^2$

(hw-lect11-2) Find The variance of The exp. dist. with param. $\lambda$.

Hint: $\int_0^\infty (y-1)^2 e^{-y} dy = 1$

(hw-lect11-3) In Example 1.22 (in text and in Lect), we found That on the average, out of 100 computers, 0.5 computers are defective.

a) What is The typical deviation we expect to see from This number (still out of 100)?

b) Suppose we do not know That The proportion of defective computers is 0.005. Then out of 100 computers, what is The maximum value we expect to see for typical deviation?

(hw-lect11-4) Find The area within $\mu_x \pm \sigma_x$ for

a) binomial $(n=20, \pi = \frac{1}{4})$
b) Poisson $(\lambda = 5)$
c) Normal $(\mu = 5, \sigma = 1)$

(hw-lect11-5) Do a qq plot of each of The 2 cont. vars. in The data from hw-lect1. (By R). Describe/Interpret The results.

Note: If you find out That There is not much you can say about The qqplot, it may be That your data is not appropriate.

This is another chance to correct The error, because later you will be doing more hw problems using your data.

So, see me, if you are not sure.