

## Lecture 9 (Ch. 2)

We are in the process of defining measures of width and spread for hists and dists. And this is where we are:

### histogram location

Sample mean:

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

~ typical  $x$  (based on sample)

### histogram spread

Sample variance:

$$s^2 = \frac{1}{n-1} \left[ \sum_{i=1}^n (x_i - \bar{x})^2 \right] \quad \text{defn. } S_{xx}$$

= computational formula

Sample std. dev. =  $s$ .

~ typical deviation/spread  
(based on sample).

Other/related measures of spread:

$s^2$  itself

$$S_{xx} = \sum_i (x_i - \bar{x})^2$$

= numerator of  $s^2$ .

### distribution/population location

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

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

~ typical  $x$  (based on pop./dist)

Last time, we found the mean of  $N(\mu, \sigma)$  and  $\text{Pois}(\lambda)$ .

$$\mu_x = \mu$$

$$\mu_x = \lambda$$

Today, we find the  $\mu_x$  of  $\text{Binomial}(n, \pi)$ , and then

dist. spread

?

(below)

E.g. # of coins tossed, or sample size, ...

Mean of Binomial ( $n, \pi$ )?

$$E[x] = \sum_{x=0}^n \frac{n!}{x!(n-x)!} \pi^x (1-\pi)^{n-x} \cdot x$$

or  
 $\mu_x$

$x=0$  contributes zero to the sum

relabel  $\sum_x$   
and  
note that  
 $\frac{x}{x!} = \frac{1}{(x-1)!}$

$$= \sum_{x=1}^n \frac{n!}{(x-1)!(n-x)!} \pi^x (1-\pi)^{n-x}$$

$y = x-1$

$$= \sum_{y=0}^{n-1} \frac{(n-1)!}{y!(n-y-1)!} \pi^{y+1} (1-\pi)^{n-y-1}$$

$$= n\pi \sum_{y=0}^{n-1} \frac{(n-1)!}{y!(n-y-1)!} \pi^y (1-\pi)^{n-y-1}$$

$$= (n+1)\pi \sum_{y=0}^m \frac{m!}{y!(m-y)!} \pi^y (1-\pi)^{m-y}$$

$m = n-1$

$$= 1 = \sum_{y=0}^m p(y)$$

$$= \underbrace{(n+1)}_n \pi$$

$$E[x] = n \cdot \pi$$

2 params of binomial Note  
Note  $E[x]$  is not a function of  $\pi$ .

E.g. 1.23:

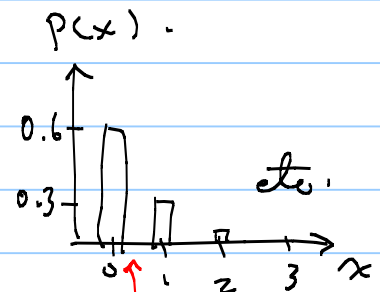
# of Bads out of 100	$x$	0	1	2	3	4
	$p(x)$	.6058	.3044	.0757	.0124	...

$$E[x] = \sum_{x=0}^{100} x p(x) = 0(.6058) + 1(.3044) + \dots$$

$$= n\pi = 100(.005) = 0.5$$

Easy way.

Hard way.



On avg. 0.5 out of 100  
(i.e. 1 out of 200)  
computers are defective.

Now, let's finish The 2x2 table:

### histogram location

Sample mean:

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

~ typical  $x$  (based on sample)

### distribution/population location

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

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

~ typical  $x$  (based on pop./dist.)

### histogram spread

Sample variance:

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

= computational formula

Sample std. dev. =  $s$ .

~ typical deviation/spread (based on sample).

Other/related measures of spread:

$s^2$  itself

$$S_{xx} = \sum_i (x_i - \bar{x})^2$$

= numerator of  $s^2$ .

### dist./pop. spread

NEW

dist./pop. Variance

$$\sigma_x^2 \equiv V[x]$$

$$= \begin{cases} \sum_x (x - E[x])^2 p(x) \\ \int_{-\infty}^{\infty} (x - E[x])^2 f(x) dx \end{cases}$$

Don't drop This  $x$ , like The book does.

dist./pop. standard dev. =  $\sigma_x$

~ typical dev./spread (based on pop./dist.)

Now, let's find The var. of some of our named distributions

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

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

Binomial  $(n, \pi)$ :  $\sigma_x^2 = V[X] = \sum_x (x - n\pi)^2 p(x) \stackrel{\substack{\mu_x \text{ of Binomial} \\ p(x) \text{ of Bin.}}}{=} \dots = n\pi(1-\pi)$ .  
 See hw and/or prelab for better understanding This formula!

Poisson  $(\lambda)$ :  $\sigma_x^2 = V[X] = \sum (x - \lambda)^2 \hat{p}(x) \stackrel{\text{Poisson}}{=} \dots = \lambda$   
 Recall  $E[X] = \lambda \leftarrow \text{same!}$

Summary (Some of These are done in hw)

	binomial $(n, \pi)$	poisson $(\lambda)$	$N(\mu, \sigma)$	Unif $(a, b)$	Exp $(\lambda)$
$E[X] = \mu_x$	$n\pi$	$\lambda$	$\mu$	$(a+b)/2$	$1/\lambda$
$V[X] = \sigma_x^2$	$n\pi(1-\pi)$	$\lambda$	$\sigma^2$	$(b-a)^2/12$	$1/\lambda^2$

Again note The diff. between  $\mu_x (\equiv E[X])$  and  $\mu$ .

$\uparrow$   
 mean of any dist.

$\uparrow$   
 $\mu$  param. of  $N(\mu, \sigma)$ .

READ!

Jargon

We now have measures of location & spread for hists & dists.

histograms	vs.	distributions / pop.
Sample mean $\bar{x}$	vs.	distv. mean $E[x] \equiv \mu_x$
" Variance $s^2$	vs.	" Variance $V[x] \equiv \sigma_x^2$
" Std. dev. $s$	vs.	" Std. dev. $\sigma_x$
sample "Statistic"		pop. "parameters"

One says that  $\bar{x}$  is a point estimate of  $\mu_x$ , Etc.  
This statement may seem obvious, because of the similarity of the names (e.g. sample mean and pop. mean); but it's not!

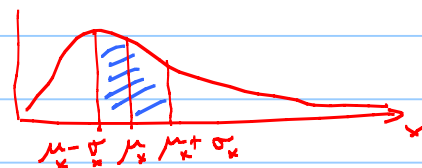
Let's not forget "areas". Before Ch.2 we used to say things like

For  $N(\mu, \sigma)$ : 68% of the area is within 1  $\sigma$  of  $\mu$ .

But now, that translates to

68% of the area is within 1 std dev. of the mean.

And now we can say things like That  
for any distr. e.g. Poisson, ---



Computing "areas" like this will (eventually) allow us to provide some measure of confidence as to what  $\mu_x$  is, based on observed data.

Recall,

For hists: "area" = prob. of times  $x$  is observed to be within ...

For dists: "area" = prob. of times  $x$  is expected to be within ...

because we don't know the pop./distr. But if we assume the distr. that describes the pop., then we would expect ...

# Switching gears

(99 plots)

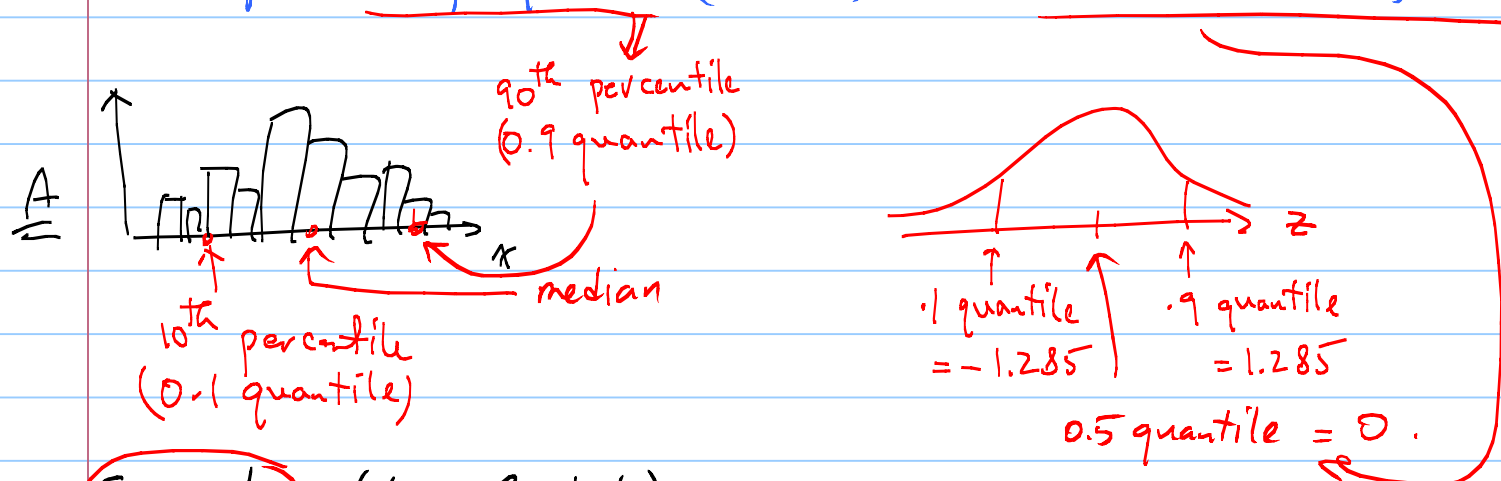
## A very powerful tool!

The business of estimating pop. params from sample stats refers to any distr. E.g., one says that  $\bar{x}$  and  $s$  provide point estimates of  $\mu$  and  $\sigma$  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 distr?

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

A: compare sample quantiles (of data) with distr. (or Theoretical) quantiles.



Example: (Very Crude!) Here is (sorted) data:

Data = -1, +1, 3, 4, 4.5, 5, 5.5, 6, 6.5, 8, 9

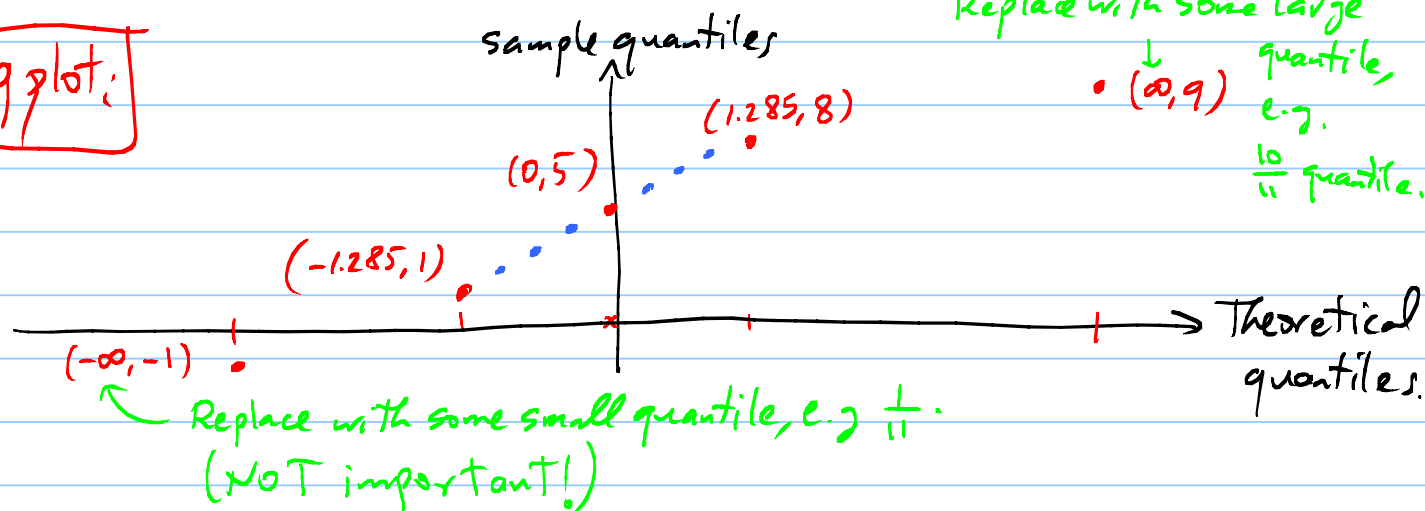
quantile prob. = 0 0.1 0.2 0.3 0.4 0.5 0.6 0.7 0.8 0.9 1.0

theoretical quantile =  $-\infty$  -1.285 --- 0 --- +1.285  $\infty$

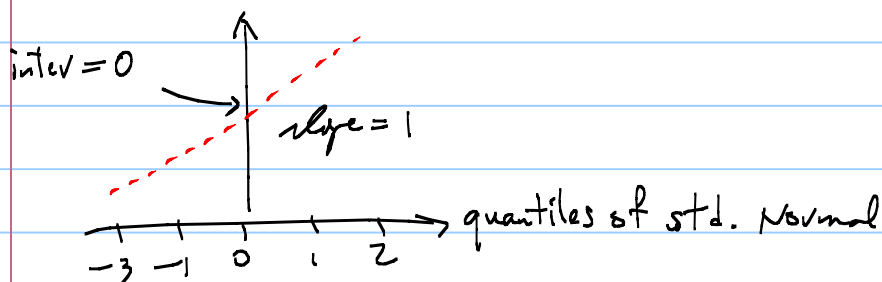
Now, we want to see if The sample quantiles "line-up" with The distribution quantiles. In Ch. 3 we will learn That The best way of seeing if 2 columns of numbers "line-up" is to plot one vs. The other, i.e.

	0	0.1	0.2	---	0.9	1.0
y = sample quantile	-1	+1	3		8	9
x = dist. quantile	$-\infty$	-1.285		---	1.285	$+\infty$

qq plot:

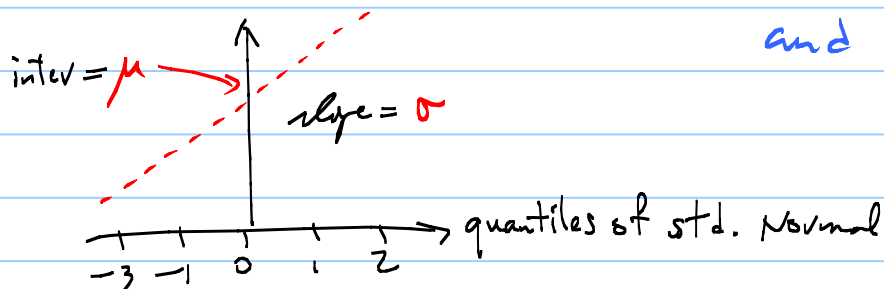


If the histogram is consistent with a std. Normal, then the quantiles/percentiles of data should be equal/comparable to those of the distr.. Then the qq plot should be a straight diagonal line (intercept=0, slope=1).



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  $\mu$ .

The proof is easy, but later. For now, focus on the concept and the use of qqplots

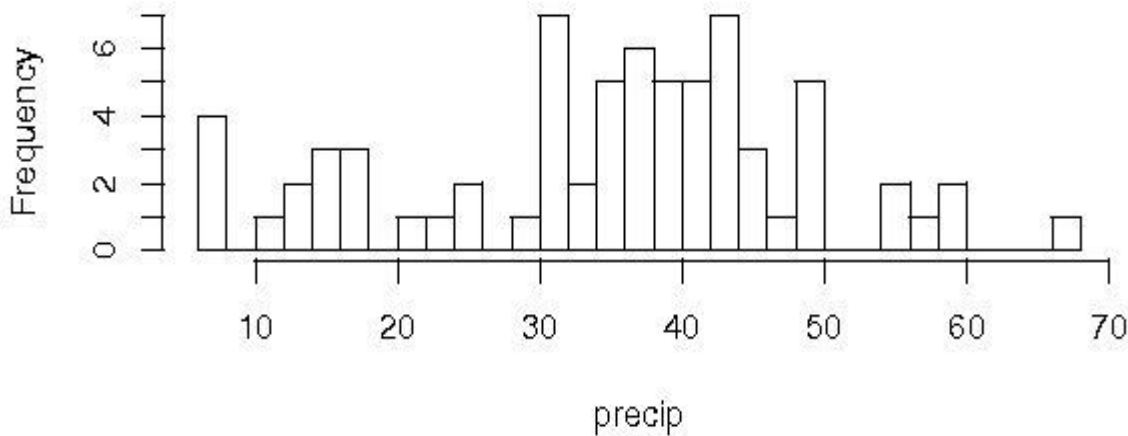


In R: `qqnorm(x)` where  $x$  is the vector of data.

## Example

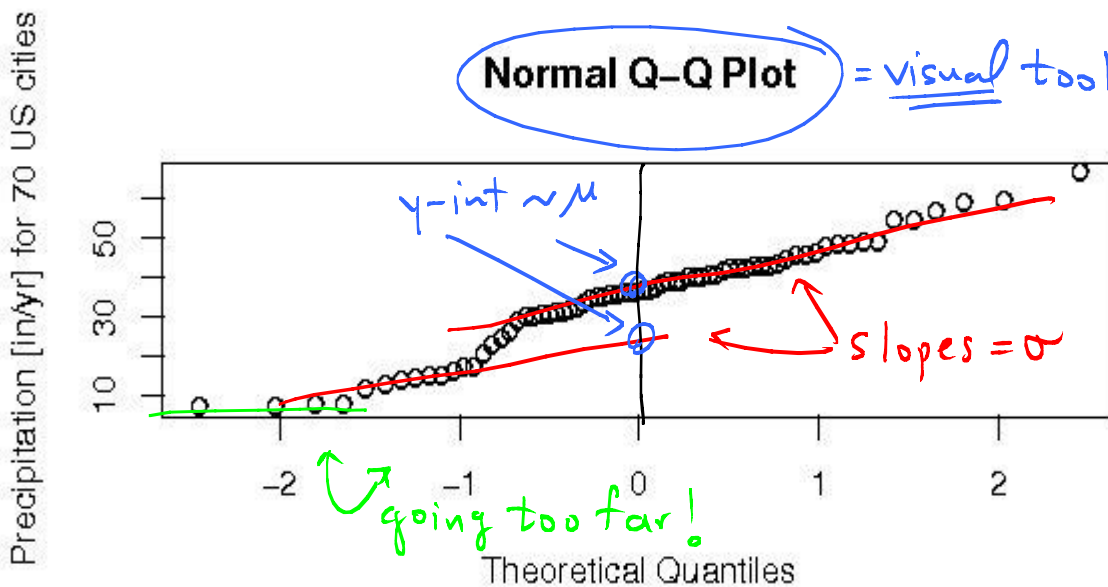
From The histogram, it's hard to tell if The data come from a normal dist., especially because hists depend on bin size.

Histogram of precip



Normal Q-Q Plot

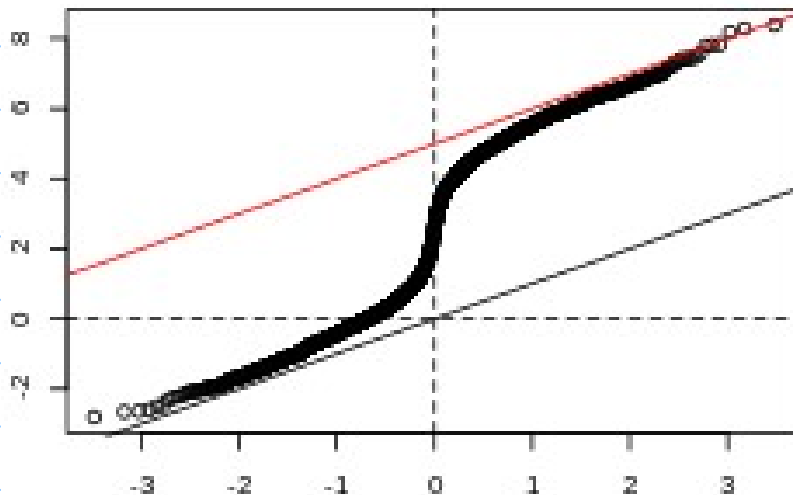
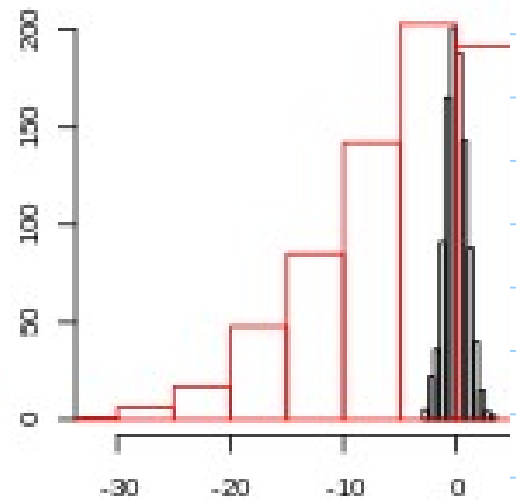
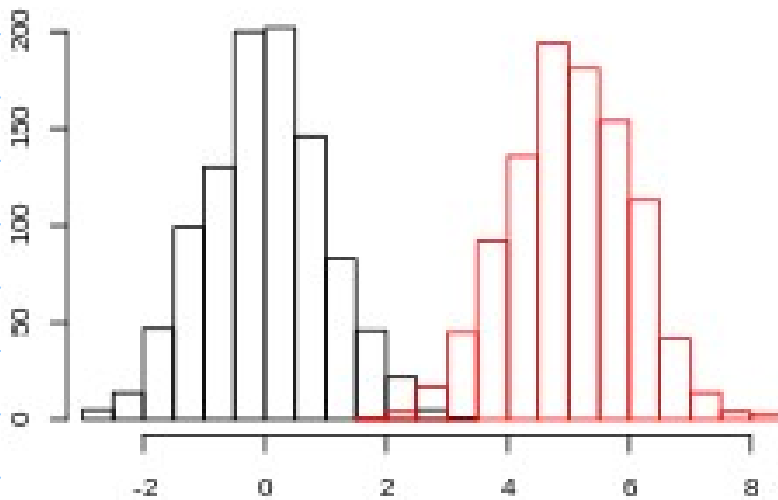
= visual tool.



- ⇒ The qq-plot looks mostly linear ⇒ Data are consistent with Normal dist.
- ⇒ Looking deeper, in fact it looks like data may have come from a bi-modal distr. (ie. 2 normals with same  $\sigma$  but different  $\mu$ 's) (However, FKI, see next page)

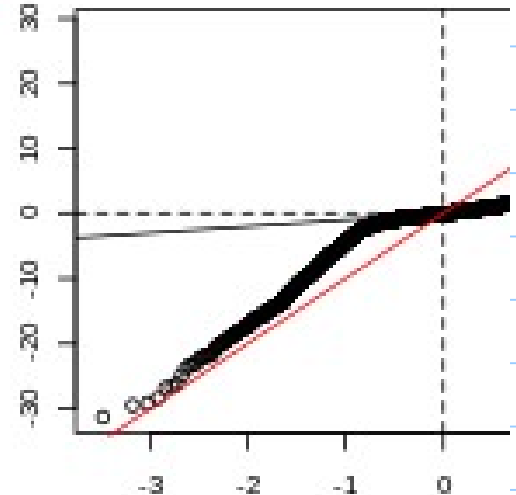


FYI



n = 1000

```
x = rnorm(n, 0, 1)
y = rnorm(n, 5, 1)
hist(x, xlim=range(x,y))
hist(y, add=T, border=2)
```

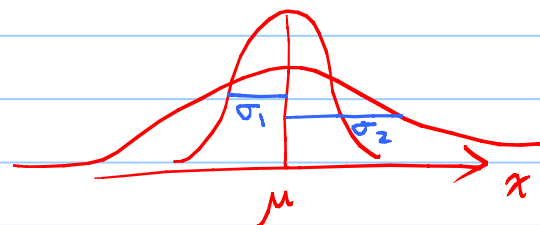
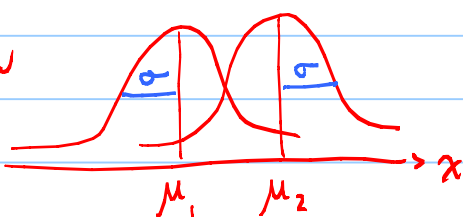


n = 100

```
x = rnorm(n)
y = rnorm(n, 5, 1)
hist(x, xlim=range(x,y))
hist(y, add=T, border=2)
```

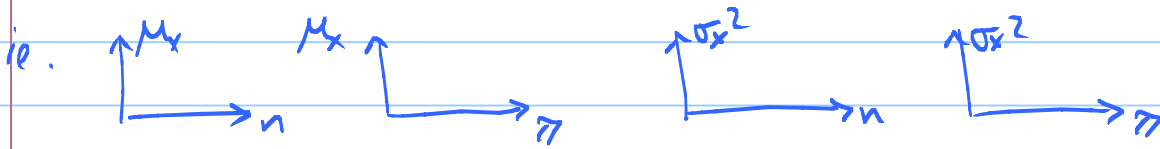
A "bend" in the qq-plot can be interpreted as

- 1) two normals with the different  $\mu$ 's, and the same sigma (if the tail ends of qqplot are parallel)
- 2) two normals with the same  $\mu$ , and different sigma's.



### hw\_lect9-1

a) Consider the binomial dist. with params  $n, \pi$ . Draw four figures that show qualitatively how its mean ( $\mu_x$ ) and variance ( $\sigma_x^2$ ) vary with  $n$  and  $\pi$ .



Suppose we toss  $n=100$  unfair coins, with an unknown  $\pi$ .

- What is the expected number of heads out of  $n$ ? (The answer depends on  $\pi$ ),
  - What is the typical deviation in the number of heads out of  $n$ ? (The answer depends on  $\pi$ ).
  - What is the largest typical deviation of the number of heads out of  $n$ ? (The answer is a number!)
- Hint: Consult your graph of variance vs.  $\pi$ , in part a).

### hw-lect9-2

For the exponential distribution with parameter  $\lambda$ , find the variance. Hint: You may use this integral:

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

hw\_lect9\_3: Find the prob that  $x$  is within 1 standard deviation of the mean, for

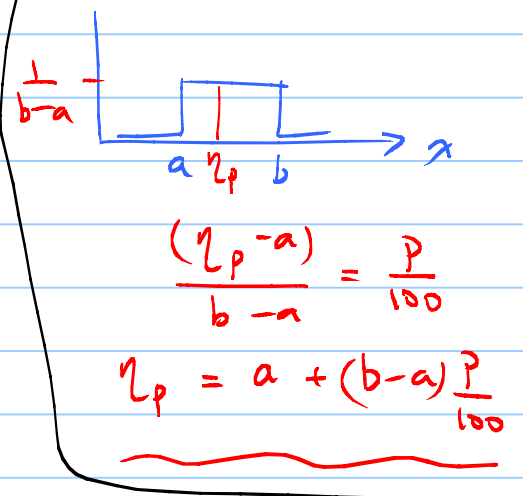
- binomial ( $n=20, \pi = \frac{1}{4}$ )
- poisson ( $\lambda=5$ )
- Normal ( $\mu=5, \sigma=1$ )

### hw-lect9-4

It can be shown that the  $p$ th percentile of  $\text{Unif}(a,b)$  is given by  $\eta_p(a,b) = a + (b-a) \frac{p}{100}$

- What's the  $p$ th percentile of  $\text{Unif}(0,1)$ , i.e.  $\eta_p(0,1)$ ?
- Write  $\eta_p(a,b)$  in terms of  $\eta_p(0,1)$ .
- What will the plot of  $\eta_p(a,b)$  vs.  $\eta_p(0,1)$  look like? What are the slope &  $y$ -intercept?

Later, check the soln to see the moral



### hw-lect9\_5

Do a qq-plot of each of the 2 cont. vars. in the data you collected. By R. Describe/Interpret the result.

Note: If you find out that there is not much you can say about the qq-plot, 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.