*Lecture 14 (Ch 3-5)*
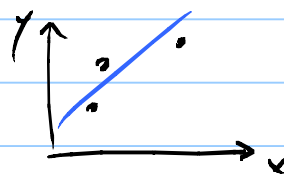
We know that one can overfit data on $x, y$, if one uses a high-order polynomial in poly. regression. Recall that the main reason this can happen is because such a regression model will have a lot of parameters.

In <u>multiple regression</u> there is yet another way that overfitting can happen even w/o including high-order terms in the model.

Consider 3 cases on $y$ and $x_1$:

A model like $y = \alpha + \beta x_1$ (a line)

Cannot overfit that data

But a model like $y = \alpha + \beta_1 x_1 + \beta_2 x_2$ ← (a plain) overfits completely.
The reason is because in that case the 3 cases are in 3D (not 2D), and there is always one plain that goes thru 3 points exactly.

Note that the additional variable $x_2$ can even be completely unrelated to $y$! It can even be just random values!

In other words, by arbitrarily making the space big, we opened up the possibility of overfitting.

So one can <u>overfit</u> even a multiple regression model without any non-linear (eg. quadratic, cubic, ...) terms.

You may think this is happening only because I have 3 cases here. But even with more cases, one can still overfit by simply including more (even random) predictors in the model, if there are many more params in regression than cases.

This overfitting problem is not specific to regression. <u>ALL</u> models can overfit when they are too large. CS students: WATCH OUT!

One last thing before we leave regression (until Ch. 11)

Here is an explanation of $df = n-1, n-2, \cdots, n-(k+1)$:

Q $\quad \bar{Y} = \frac{1}{n} \sum_{i}^{n} Y_i \qquad$ why $n$ ?

A $\quad \{Y_1, Y_2, \cdots, Y_n\}$ are all indep. $\implies df = n$.

of numerator

Q $\quad s^2 = \frac{1}{n-1} \sum_{i}^{n} (Y_i - \bar{Y})^2 \quad$ why $n-1$ ?

A $\quad \{Y_1 - \bar{Y}, Y_2 - \bar{Y}, \cdots, Y_n - \bar{Y}\}$ are not all indep.

There is 1 constraint on them : $\sum_{i}^{n} (Y_i - \bar{Y}) = 0$

$\therefore df = n - 1$.

of numerator

This is one reason why $\sum (Y_i - \bar{Y})^2$ is divided by $n-1$.

Similar reasoning implies that the df for SSE is $n - (k+1)$, which is why we define $s_e^2$ as $\frac{SSE}{n-(k+1)}$. Note for $k=1$ (i.e. simple linear regression), $s_e^2 = \frac{SSE}{n-2}$.

In simple linear regression : $y = \alpha + \beta x$

$$SSE = \sum_{i=1}^{n} (y_i - \hat{y}_i)^2 \quad \text{has} \quad df = n-2 \quad \leftarrow 2 \text{ constraints}$$

**First constraint** $\quad \sum_{i=1}^{n} (y_i - \hat{y}_i) = 0$

Pf: $\frac{1}{n} \sum (y_i - \hat{y}_i) = \frac{1}{n} \sum (y_i - \hat{\alpha} - \hat{\beta} x_i)$

$$= \frac{1}{n} \sum y_i - \frac{1}{n} \sum (\hat{\alpha} + \hat{\beta} x_i)$$

$$= \frac{1}{n} \sum y_i - \left(\hat{\alpha} + \hat{\beta} \frac{1}{n} \sum x_i\right)$$

$$= \left[\bar{y} - (\hat{\alpha} + \hat{\beta} \bar{x})\right] = 0$$

$$\underbrace{\qquad\qquad} = \bar{y} \quad (2^{nd} \text{ Normal equ.})$$

**2$^{nd}$ constraint** : $\sum (y_i - \hat{y}_i) x_i = 0$,

Pf: $\frac{1}{n} \sum (y_i x_i - \hat{y} x_i) = \frac{1}{n} \sum [x_i y_i - x_i (\hat{\alpha} + \hat{\beta} x_i)]$

$$= \overline{xy} - \hat{\alpha} \bar{x} - \hat{\beta} \overline{x^2}$$

$$= \overline{xy} - \bar{x} \overline{(\bar{y} - \hat{\beta} \bar{x})} - \hat{\beta} \overline{x^2}$$

$$= (\overline{xy} - \bar{x}\bar{y}) - \hat{\beta} (s_x^2) = 0.$$

$$\underset{\Downarrow}{} \quad \frac{\overline{xy} - \bar{x}\bar{y}}{s_x^2} \quad (1^{st} \text{ Normal equ.})$$

This page is only FYI, for now

All of Ch.3 has been about understanding the relationship between several continuous variables. What about categ. vars?

For categorical data the relationship is best captured through the contingency table: C-table
↖ aka confusion matrix.

**Data**

| x | Y |
|---|---|
| Yes | High |
| Yes | Low |
| Yes | High |
| No | High |
| Yes | High |
| No | Low |
| No | Low |
| perhaps | medium |
| perhaps | Low |

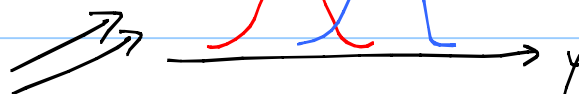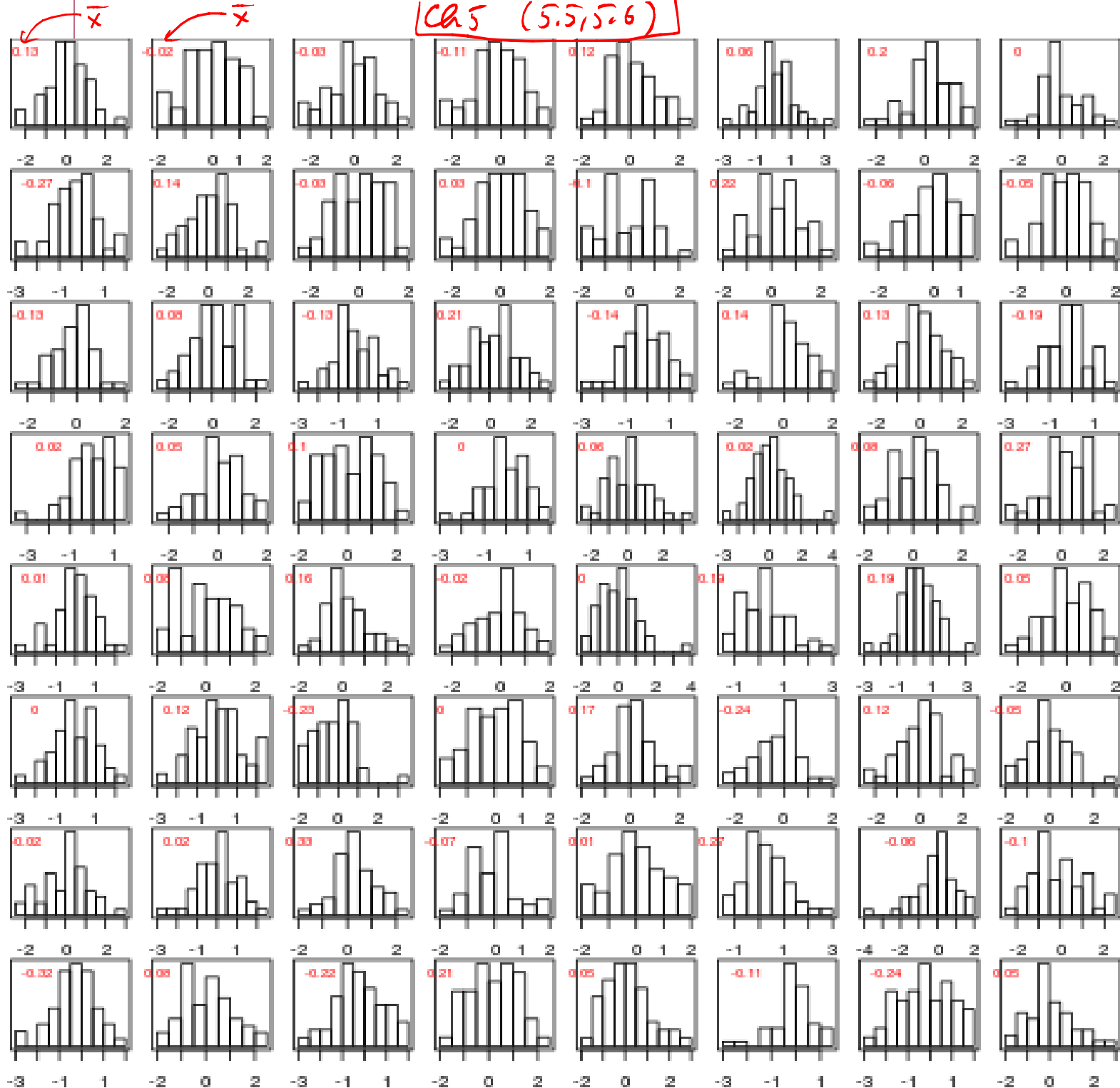| | | Y | | |
|---|---|---|---|---|
| x | | High | Low | Medium |
| | Yes | 3 | 1 | 0 |
| | No | 1 | 2 | 0 |
| | perhaps | 0 | 1 | 1 |

∃ Relationship between x and y.
↑
Maybe "positive" or "negative".

3 variables X, Y, Z ⟹ Cube = Set of Contingency Tables.

**Q**: what about mixed (discrete and cont)?

E.g. { X = 0, 1
      { y = Continuous

X=0    X=1

**A**: conditional histograms → y

Top annotations: $\bar{x}$ ... $\bar{x}$
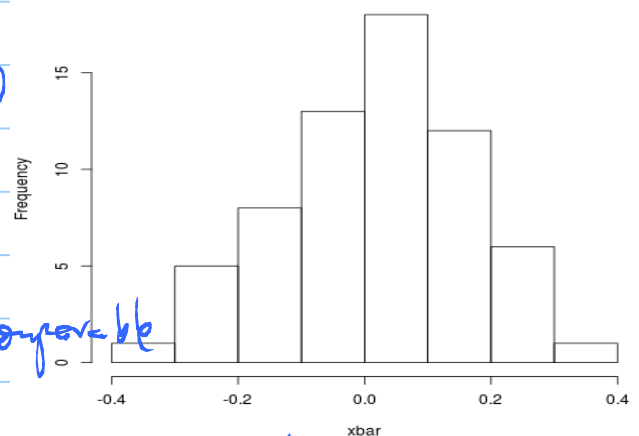
Ch 5 (5.5, 5.6)

```
ntrial = 64
xbar = numeric(ntrial)
par(mfrow=c(8,8))
  for( trial in 1:ntrial ){
  x = rnorm(50, 0, 1)        ← Try rexp(50,1)
  hist(x, breaks=10)
  xbar[trial] = mean(x)
  }
hist(xbar, main="")
```

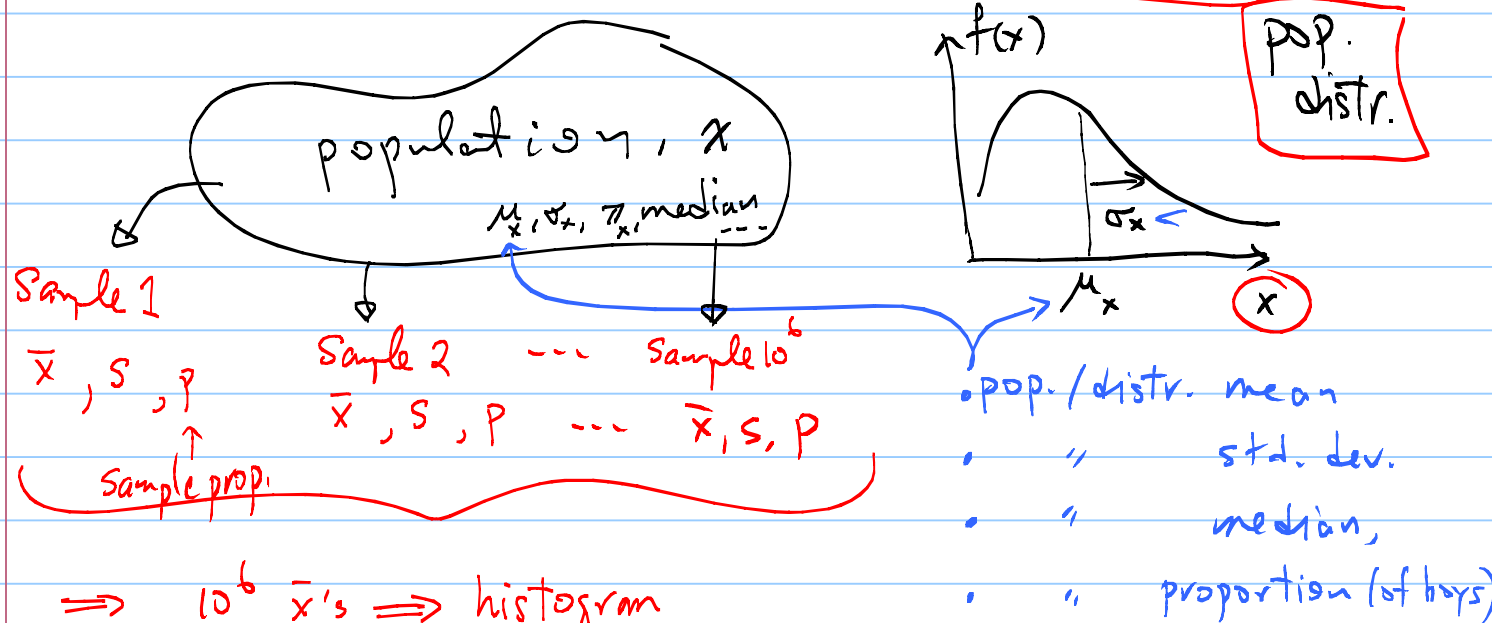Q: What's $\bar{x}$ in each hist above?
   What's the mean of the $\bar{x}$'s?  } comparable

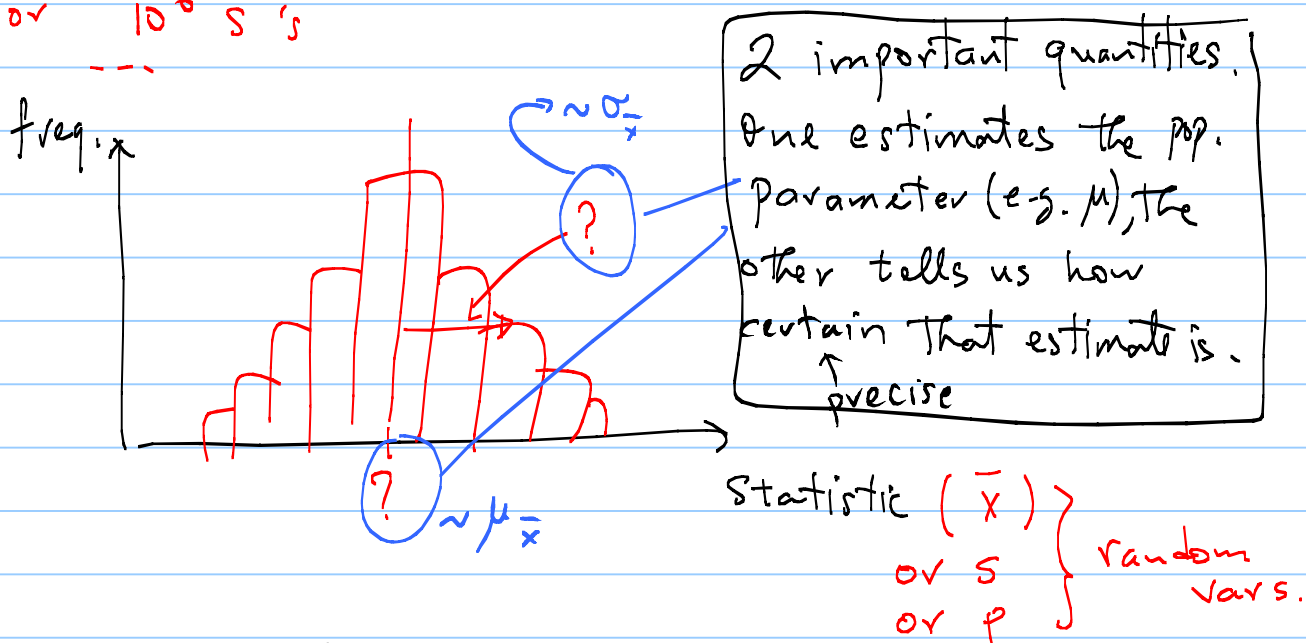Q: What's s in each hist above?
   What's s of the $\bar{x}$'s?  } very different!

# Sampling Distribution : **Extremely Important !!**

population, $x$

$\mu_x, \sigma_x, \pi_x$, median ---



POP. distr.

$f(x)$

$\sigma_x$

$\mu_x$

$\bar{x}$

**Sample 1**
$\bar{x}, S, \hat{p}$
↑
Sample prop.

**Sample 2**
$\bar{x}, S, P$

--- **Sample $10^6$**
$\bar{x}, S, P$

- pop./distr. mean
- " std. dev.
- " median,
- " proportion (of boys)

$\Rightarrow 10^6 \; \bar{x}$'s $\Rightarrow$ histogram
or $10^6 \; S$'s
---

freq. $_x$

$\sim \sigma_{\bar{x}}$

?

?

$\sim \mu_{\bar{x}}$

**2 important quantities.**
One estimates the pop. parameter (e.g. $\mu$), the other tells us how certain that estimate is.
↑
precise

Statistic ($\bar{x}$)
or $S$    } random vars.
or $P$

The sampling dist. (of the sample mean) is a **distribution**, ie. a $p(x)$ or an $f(x)$ that can be derived mathematically, or simply assumed as a description of the population of all $\bar{x}$'s. The only reason I talk about a histogram is to make the concept of the sampling dist. more intuitive. The histogram is sometimes called the "empirical sampling dist."

Note that the sampling distr. is the distribution of a sample statistic.
For example, the sampl. distr. of the sample mean,
tells us how the sample means are distributed.

Similarly, the sampl. distr. of the sample proportion,
tells us how the sample proportions are distributed.   Etc.

---

**Q** What is the sampling distr. of $\bar{x}$ ? Normal, Poisson,...?

**A** Later!

But even without knowing the dist., we can still find its
mean ( $E[\bar{x}]$ or $\mu_{\bar{x}}$ ) and Variance ( $V[\bar{x}]$ or $\sigma_{\bar{x}}^2$ ):

If the population (ie. distribution) has mean $\mu_x$ and std. dev. $\sigma_x$, then

Mean of the Sampling distr. of sample mean $(\mu_{\bar{x}})$ :

Std. dev. " " " " " " " " $(\sigma_{\bar{x}})$ :

$$\mu_{\bar{x}} = E[\bar{x}] = \mu_x \quad \longleftarrow \text{pop. mean}$$

$$\sigma_{\bar{x}} = \sqrt{V[\bar{x}]} = \sigma_x/\sqrt{n}$$

pop. std. dev. $\longrightarrow$

$\longleftarrow$ sample size

proof, below.

"sometimes called "standard error of mean."

Derivation: Suppose we do not know the distr. of the population ($p(x)$, $f(x)$), but we do know its $\mu_x$ and $\sigma_x$

Of course, if you do know the pop. distr., then you can compute $\mu_x$, $\sigma_x$ as before:

$$E[x] \equiv \mu_x = \sum_x x\, p(x) \qquad \left(\text{or } \int x\, f(x)\, dx\right)$$

$$V[x] = \sigma_x^2 = \sum_x (x - \mu_x)^2\, p(x) \qquad \left(\text{or } \int \text{---} dx\right)$$

Recall, $E[ax] = a\,E[x]$, $V[ax] = a^2\,V[x]$, $a = $ constant. Then

$$\mu_{\bar{x}} = E[\bar{X}] = E\left[\frac{1}{n}\sum_i^n x_i\right] = \frac{1}{n}\sum_i E[x_i] = \frac{1}{n}\,n\,\mu_x\left(\sum_i 1\right) = \mu_x.$$

$\underbrace{\qquad}_{\mu_x \ \forall i}$

$$\boxed{E[\bar{X}] \equiv \mu_{\bar{x}} = \mu_x}$$

The $i^{th}$ obs. is a random value, there is nothing special about the $i^{th}$ obs. So, just drop the "$i$". Then $E[x_i] = E[x] = \sum_x x\, p(x) = \mu_x$.

Alternatively, work out $E[x_i]$ for each $i$, e.g. $i = 1$

$$E[x_1] = \sum_{x_1} x_1\, p(x_1) = \mu_x, \quad E[x_2] = \mu_x, \quad \text{etc.}$$

$$\sigma_{\bar{x}}^2 = V[\bar{X}] = V\left[\frac{1}{n}\sum_i^n x_i\right] = \left(\frac{1}{n}\right)^2 \sum_i V[x_i] \qquad \begin{array}{l}\text{The var. of each} \\ \text{element in the pop.} \\ \text{is the var. of the pop.}\end{array}$$

$\underbrace{\qquad}_{\sigma_x^2}$

$$= \left(\frac{1}{n}\right)^2 \sigma_x^2 \left(\sum_{i=1}^n 1\right)^n = \frac{\sigma_x^2}{n} \implies \boxed{\sigma_{\bar{x}} = \sqrt{V[\bar{X}]} = \frac{\sigma_x}{\sqrt{n}}}$$

$$\bar{x} = \frac{1}{n}\sum^n x_i \longrightarrow \mu_x \qquad s_{\bar{x}} =$$

In Summary:

$\mu_{\bar{x}} \equiv E[\bar{x}] = \mu_x$   Tells us that we can use the sample mean (from the one sample of size $n$) to estimate the pop. mean $\mu_x$ with <u>accuracy</u>. ← see bottom of page.

$\sigma_{\bar{x}} \equiv \sqrt{V[\bar{x}]} = \frac{\sigma_x}{\sqrt{n}}$   Tells us that the typical deviation in $\bar{x}$ is $\frac{\sigma_x}{\sqrt{n}}$, and so it tells us how <u>precise</u> ← certain. is our estimate of $\mu_x$.
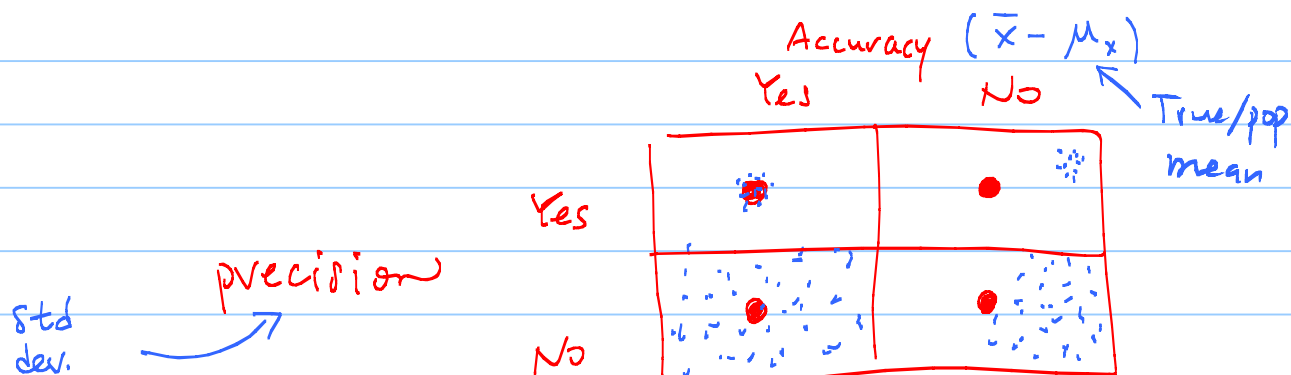
Note that $\mu_x, \sigma_x, \mu_{\bar{x}}, \sigma_{\bar{x}}$ are means and std. dev. of distributions, NOT of data. We are dealing with <u>distributions</u>, even though the thought exp. involved a <u>hist</u>.

$$\mu_x = \sum_x x \cdot p(x), \int x f(x) dx \quad ; \quad \sigma_x^2 = \sum_x (x - \mu_x)^2 p(x), \int (x - \mu_x)^2 f(x) dx$$

$$\boxed{FYI}$$

$\bar{x}$ and $s_x$ are measures of Accuracy & Precision:

and so $\mu_{\bar{x}}$ and $\sigma_{\bar{x}}$,



Accuracy $(\bar{x} - \mu_x)$

True/pop mean

precision

std dev.

Now, what is the sampling distr. of sample means?

**Thm** If the pop. is Normal $(\mu, \sigma)$, then the sampling dist. of $\bar{x}$ is Normal with

params: $N(\mu_{\bar{x}} = \mu_x = \mu, \ \sigma_{\bar{x}} = \frac{\sigma_x}{\sqrt{n}} = \frac{\sigma}{\sqrt{n}})$    **Central Limit Theorem (CLT)**

even if the pop. is NOT normal, as long as $n =$ large (say $> 30$)

─────────────────────────────

I'll go over this ↓ again tomorrow.

Now that we know the distr. of $\bar{x}$, we can compute probs. pertaining to a random (future) $\bar{x}$.   e.g. prob$(a < \bar{x} < b)$:

1a) **If** pop. distr. $(p(x), f(x))$ is given, use it to compute $\mu_x, \sigma_x$:

Eg. $\mu_x \equiv E[x] = \sum x \, p(x)$,     $\sigma_x^2 \equiv V[x] = \sum (x - \mu_x)^2 p(x)$.
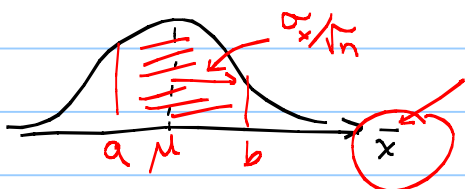
1b) **If** pop. distr. is not known, assume its $\mu_x, \sigma_x$ (Ch. 7, 8)

2) CLT $\Longrightarrow$ $\bar{x}$ is distributed as $N(\mu_x, \frac{\sigma_x}{\sqrt{n}})$

3) Standardize: $z = \frac{\bar{x} - \mu_{\bar{x}}}{\sigma_{\bar{x}}} = \frac{\bar{x} - \mu_x}{\sigma_x/\sqrt{n}} \sim N(0,1)$

4) prob$(a < \bar{x} < b)$   ← Sample mean.   Think about the meaning of this prob.

$= \text{prob}\left( \frac{a - \mu_x}{\sigma_x/\sqrt{n}} < \frac{\bar{x} - \mu_x}{\sigma_x/\sqrt{n}} < \frac{b - \mu_x}{\sigma_x/\sqrt{n}} \right)$

$\sigma_x/\sqrt{n}$     $z \sim N(0,1)$



$a \ \mu \ b$   $\bar{x}$     **Table I** .

Overfitting occurs in multiple regression even without higher powers of the predictors.
Let's see it. Consider the data on y and x1 made here:

```
set.seed(123)
n = 10
x1 = runif(n,-1,1)
y = 1 + 2*x1 + rnorm(n,0,1)
```

a) Make the scatterplot of y vs x1.
b) Perform simple linear regression, and report the R^2.
c) Generate another n cases from runif(-1,1), and call that data x2. Then repeat this
step three more times to generate x3, x4, and x5. In other words, in this step, generate
data on x2,x3,x4,x5, where they are all independent of each other, and none of them
are related to y.
d) Perform multiple linear regression on y,x1,x2,x3,x4,x5, and report R^2.


hw-lect 14-2  a) Write R code to produce The sampling distribution
of The sample __maximum__, for samples of size 50 taken from
a standard Normal. Use 5000 trials.
b) Then, repeat but for sample __minimum__.

Turn-in The code, and The resulting 2 histograms.

FYI, These distributions arise naturally when one tries to model
__extreme events__, e.g. The biggest storms, The strongest earthquakes,
The brightest stars, The smallest forms of life, etc.


hw-lect 14-3  Write R code to take 5000 samples of size
n=100 from an exponential distr. with parameter $\lambda = 2$,
and plot a qqplot of The 5000 means. Recall That if
The qqplot is a straight line, Then The histogram of The
sample means is Normal. This will show That The sampl.
dist. of sample means is Normal, even when the pop. is not!

A sampling distribution (e.g. of the sample mean) is a distribution, not a
histogram of observed sample means; the histogram of sample means discussed in class
is just an intuitive way of thinking about the sampling distribution; technically,
it's called the *empirical* sampling distribution. Of course, if the number of trials
is infinite, then the empirical sampling distribution (i.e., the histogram) approaches
the distribution. Anyway, to show that the sampling distribution is truly a
distribution (not a histogram), let's derive one mathematically - no data at all.

Consider a population described by a Bernoulli random variable, i.e., x = 0,1, following the
Bernoulli distribution, i.e., $p(x) = pi^x (1-pi)^{(1-x)}$ . Suppose we take samples of size 2.
a) Write down all the possible samples. Hint: there are only 4.
b) For each of the possible samples, compute the sample mean.
c) For each of the possible samples, compute the probability. Hint: Use Bernoulli.
d) Based on your answers to parts a-c, find the probability of each of the possible sample
means.

Note: your answer to part d *is* the sampling distribution of the sample mean! Note that it's
not a histogram, but a real distribution.