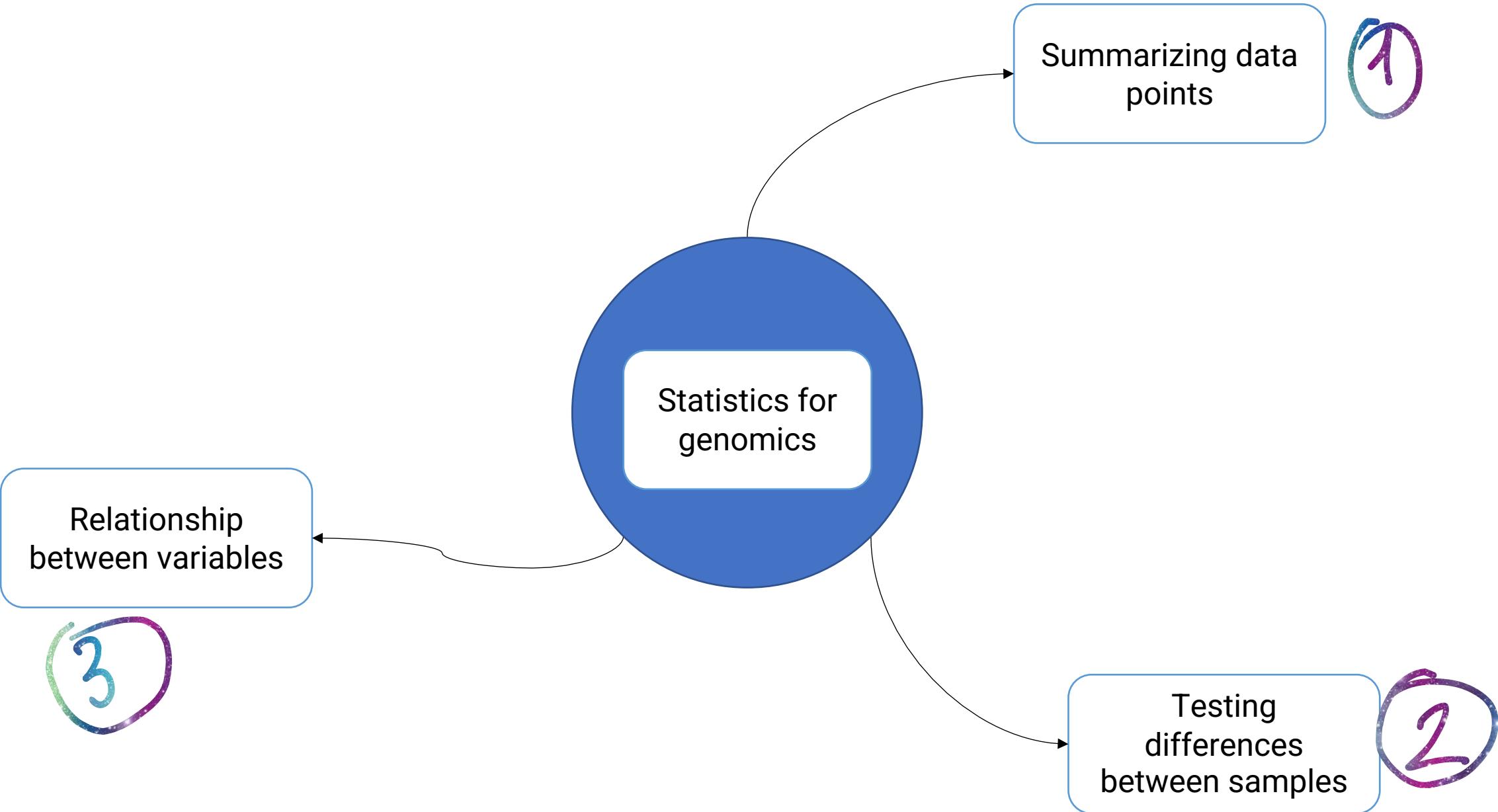


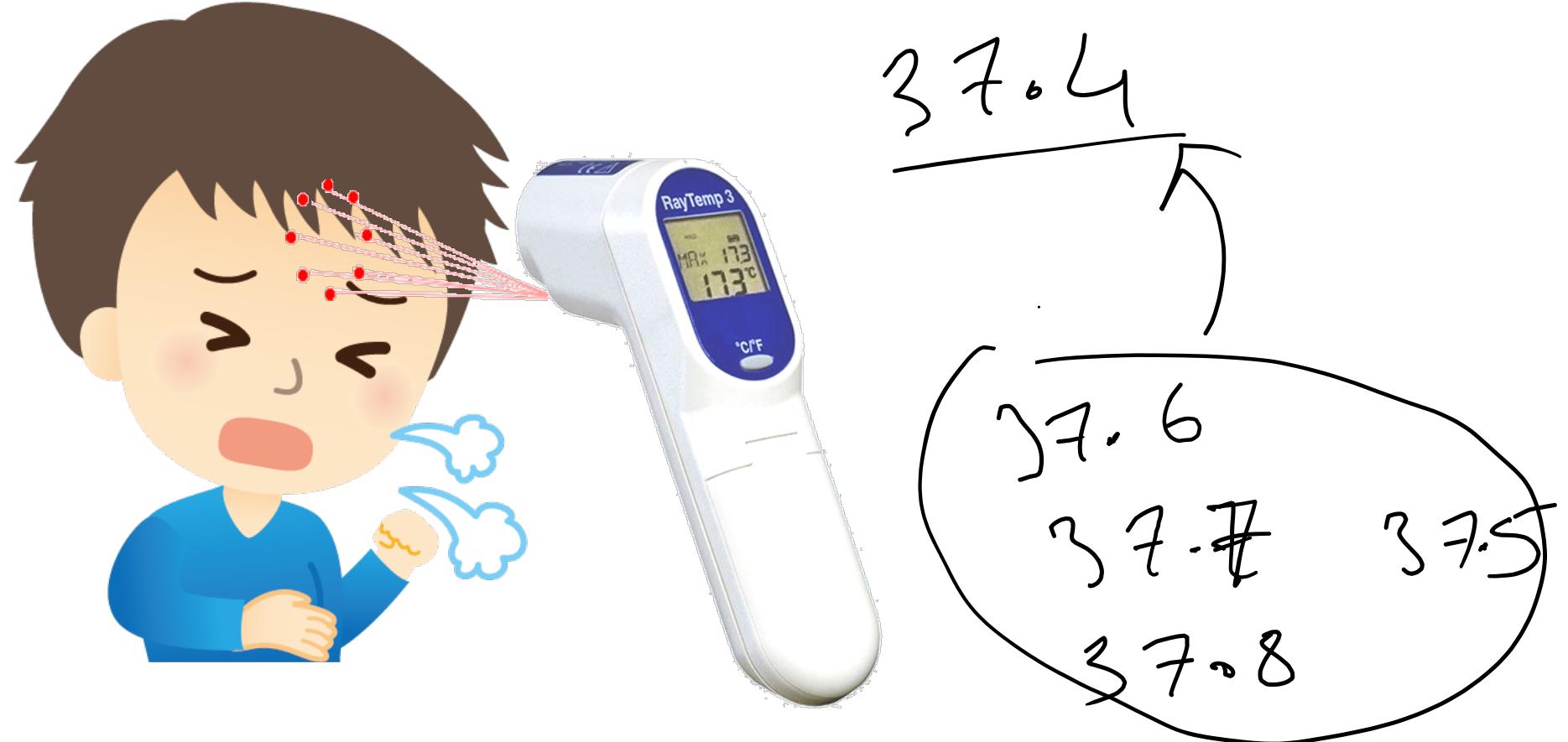
Computational genomics:  
hands on course

**Statistics for genomics**

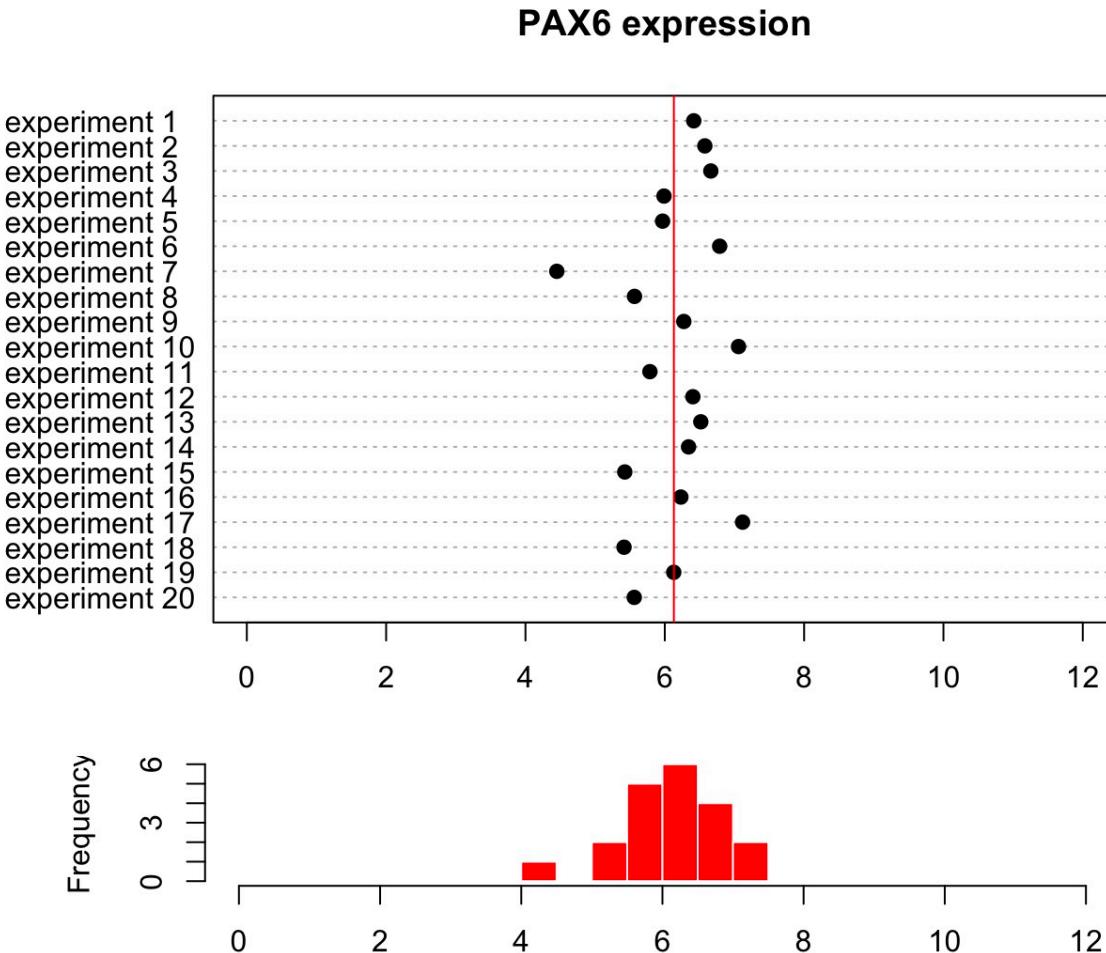


# How to summarize collection of data points:

## The idea behind statistical distributions



# How to summarize collection of data points: The idea behind statistical distributions



general  
points

# How to summarize collection of data points:

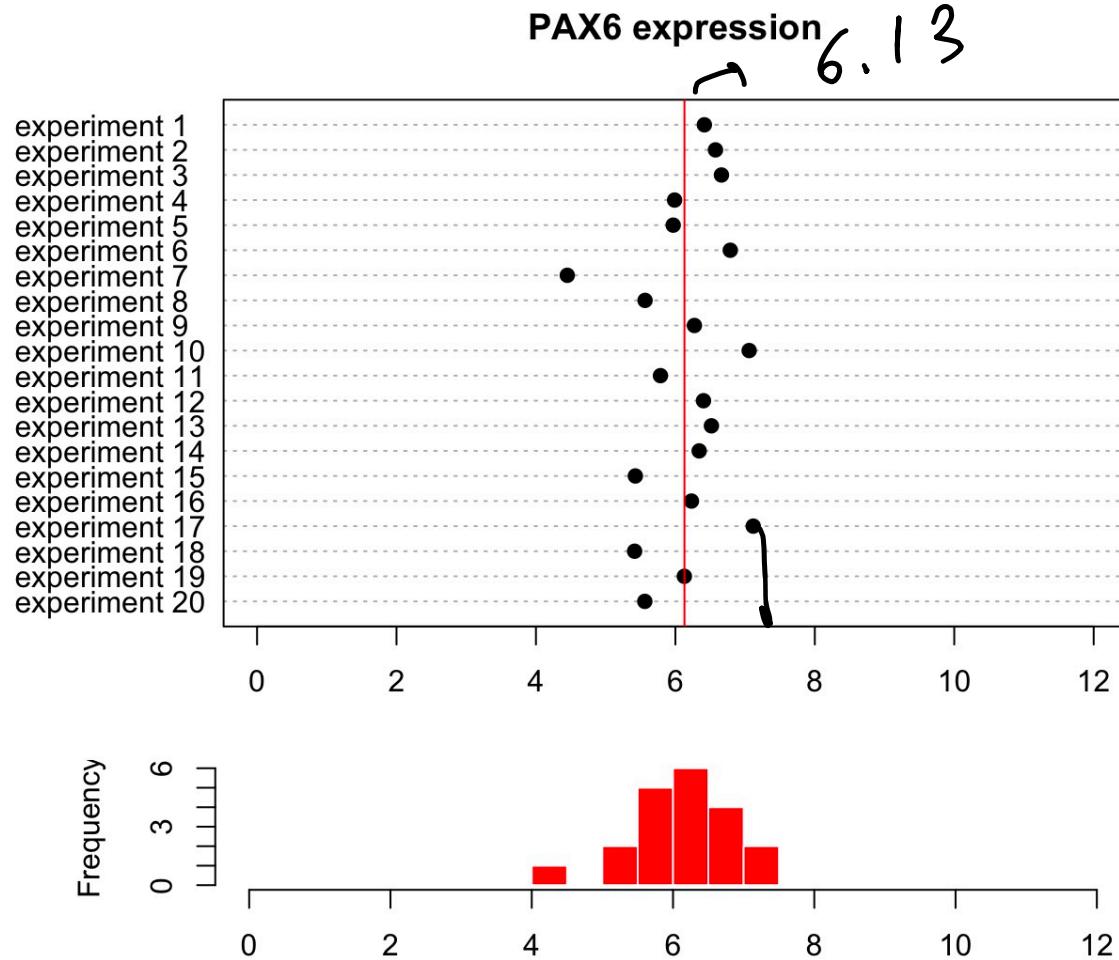
## The idea behind statistical distributions

- Experimentation -> data
- Each data point is different
  - Biological variation or measurement error is to blame

How to describe a collection of data points

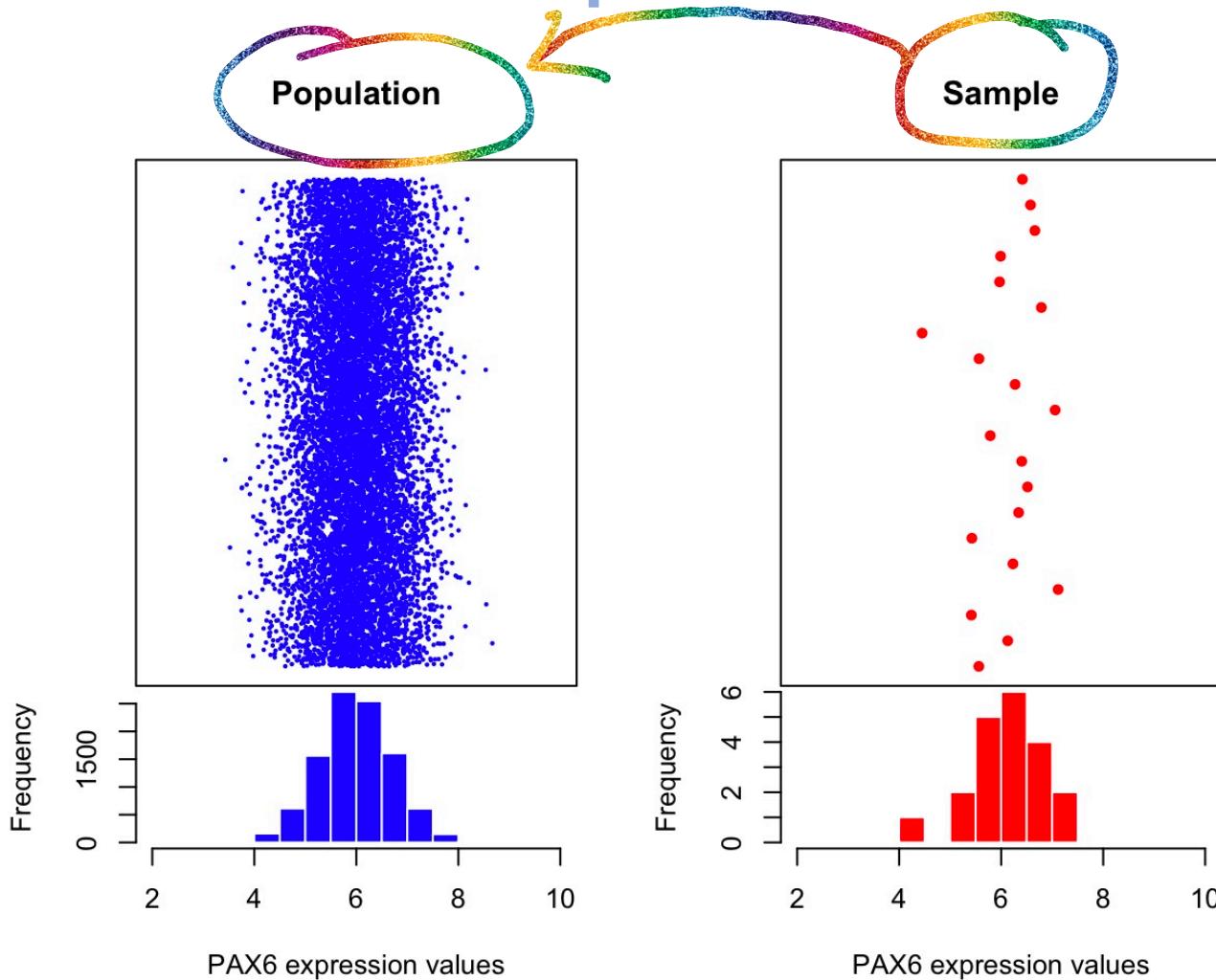
- 
- What is the most frequent value ?
  - How variable are your measurements ?

# Describing the central tendency: mean and median



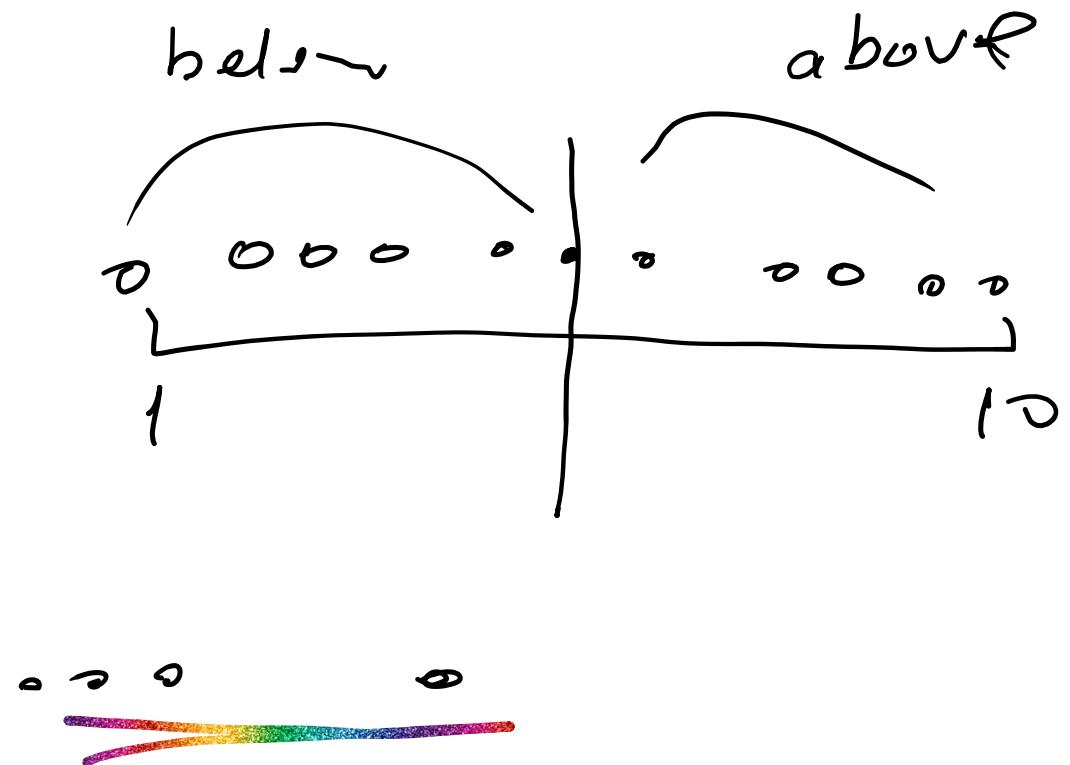
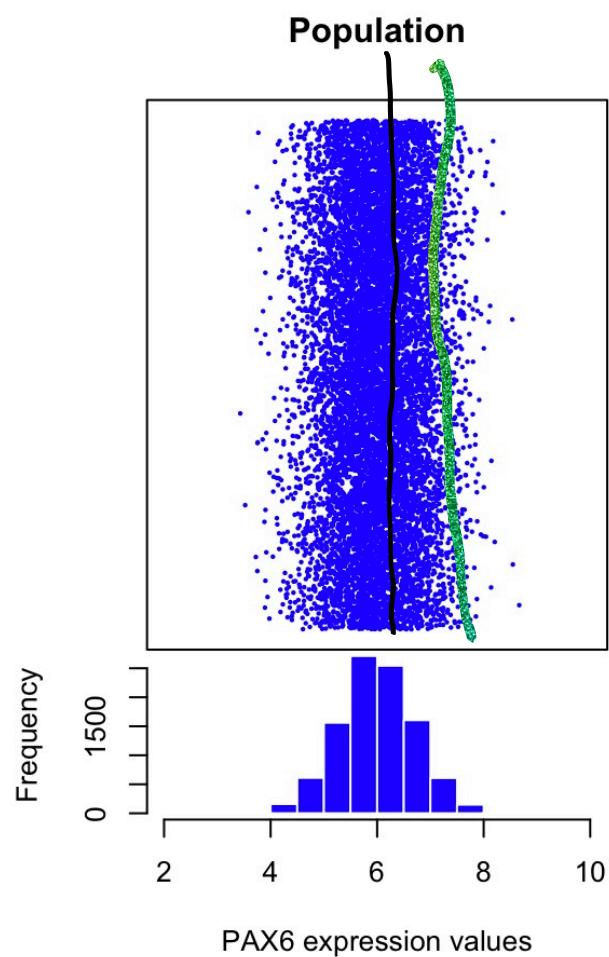
$$\bar{X} = \sum_{i=1}^n x_i / n$$
$$4 + 4.5 - 7.2 \\ 20 \\ = 6.13$$

# Describing the central tendency: population → sample

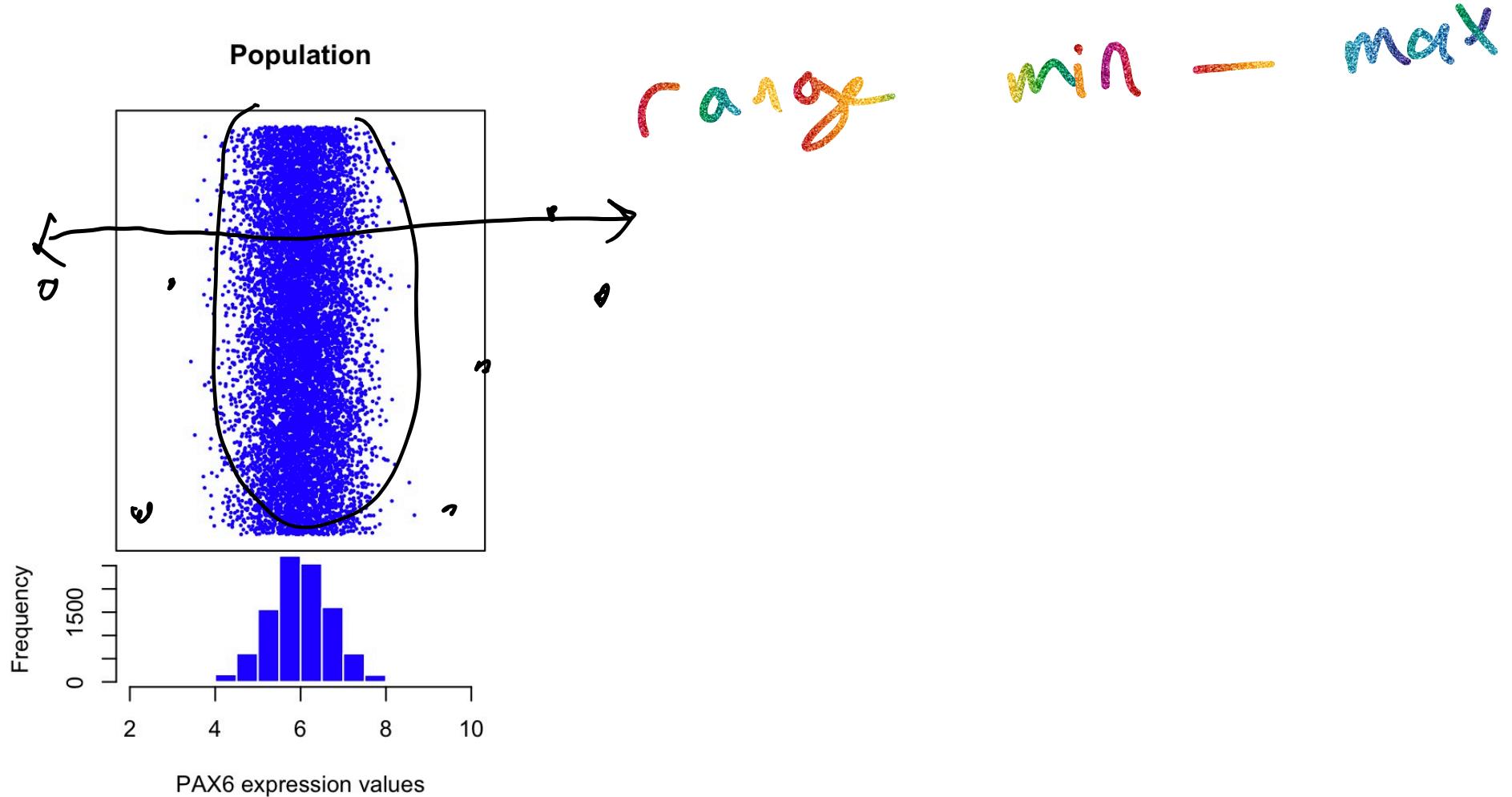


# Describing the central tendency:

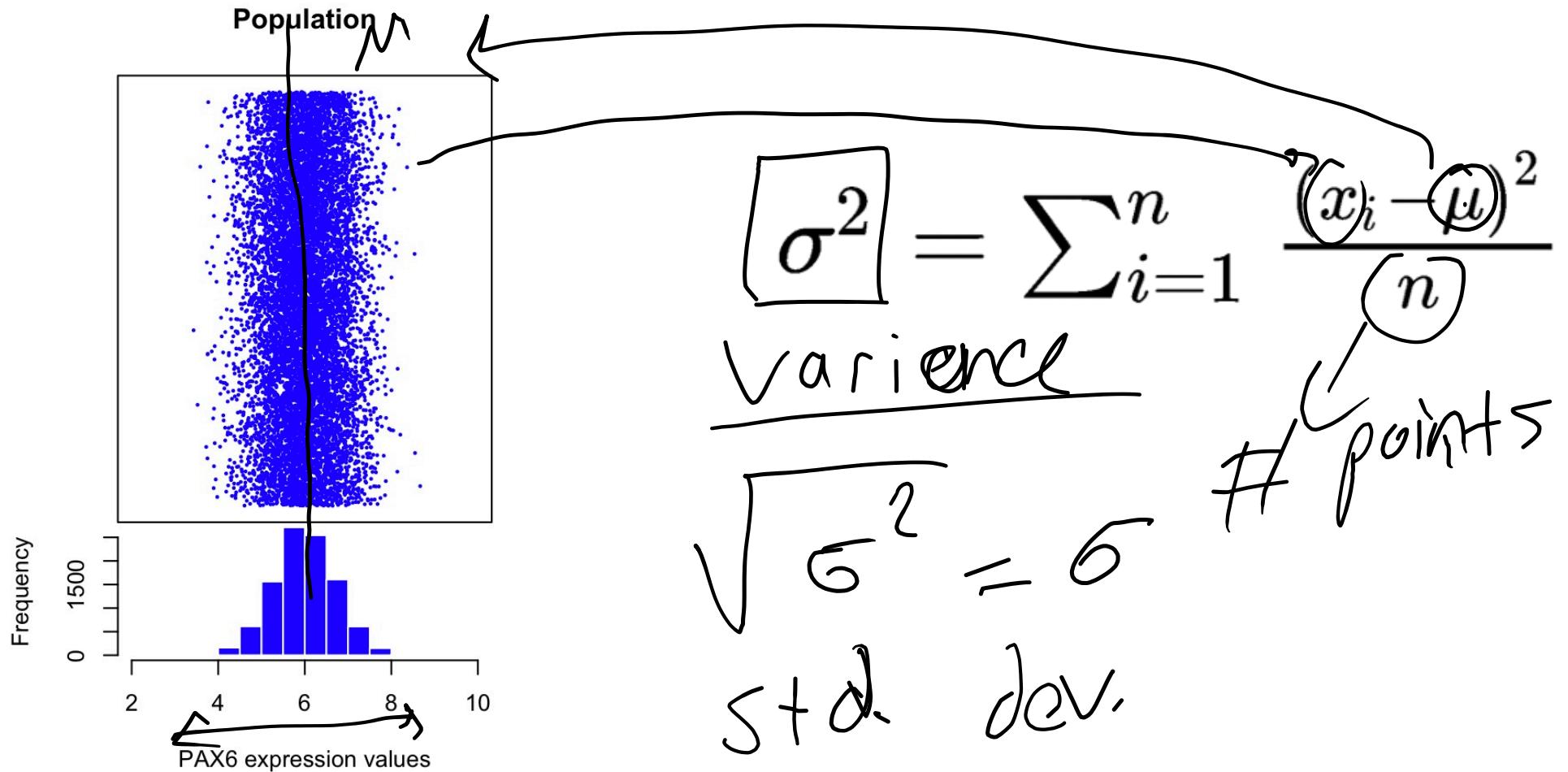
median



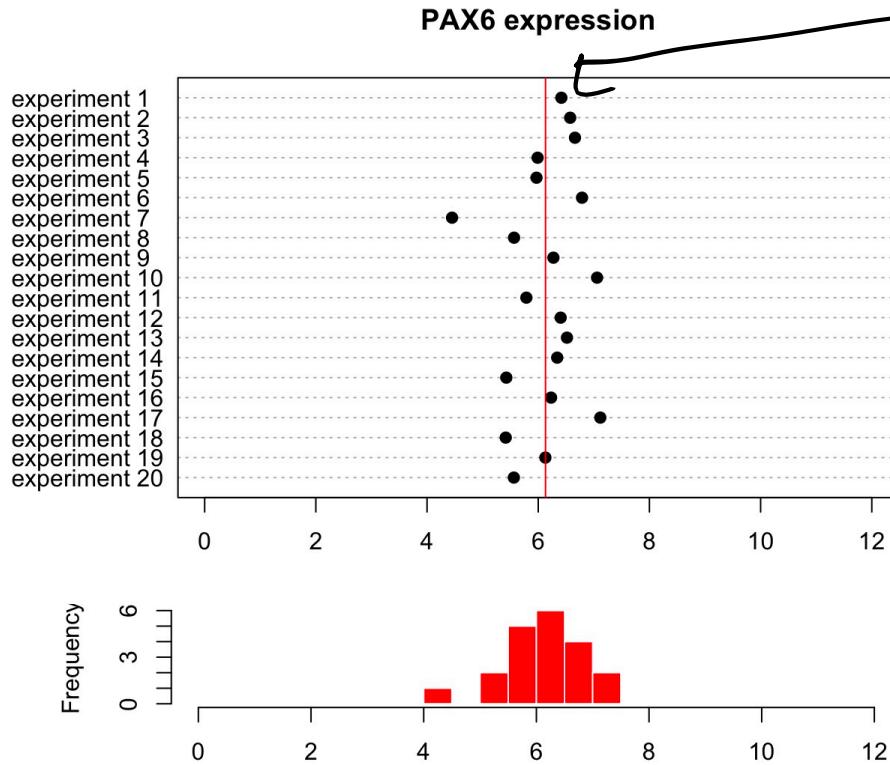
# Describing the spread: measurements of variation



# Describing the spread: measurements of variation



# Describing the spread: measurements of variation

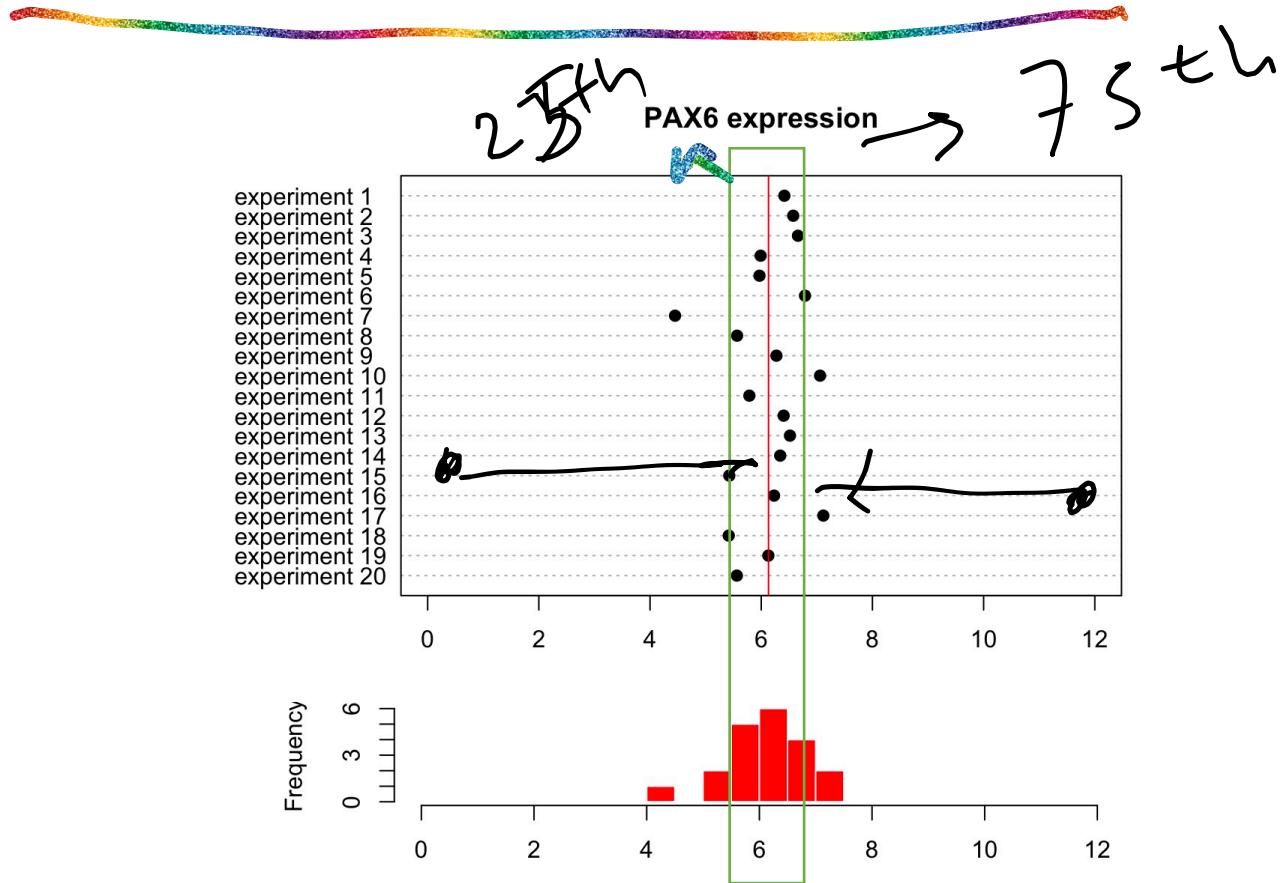


sample mean

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

# points - 1

# Describing the spread: interquartile range



# Describing the spread and central tendency: R code

`x=rnorm(20, mean=6, sd=0.7)`

`mean(x)`  
`median(x)`

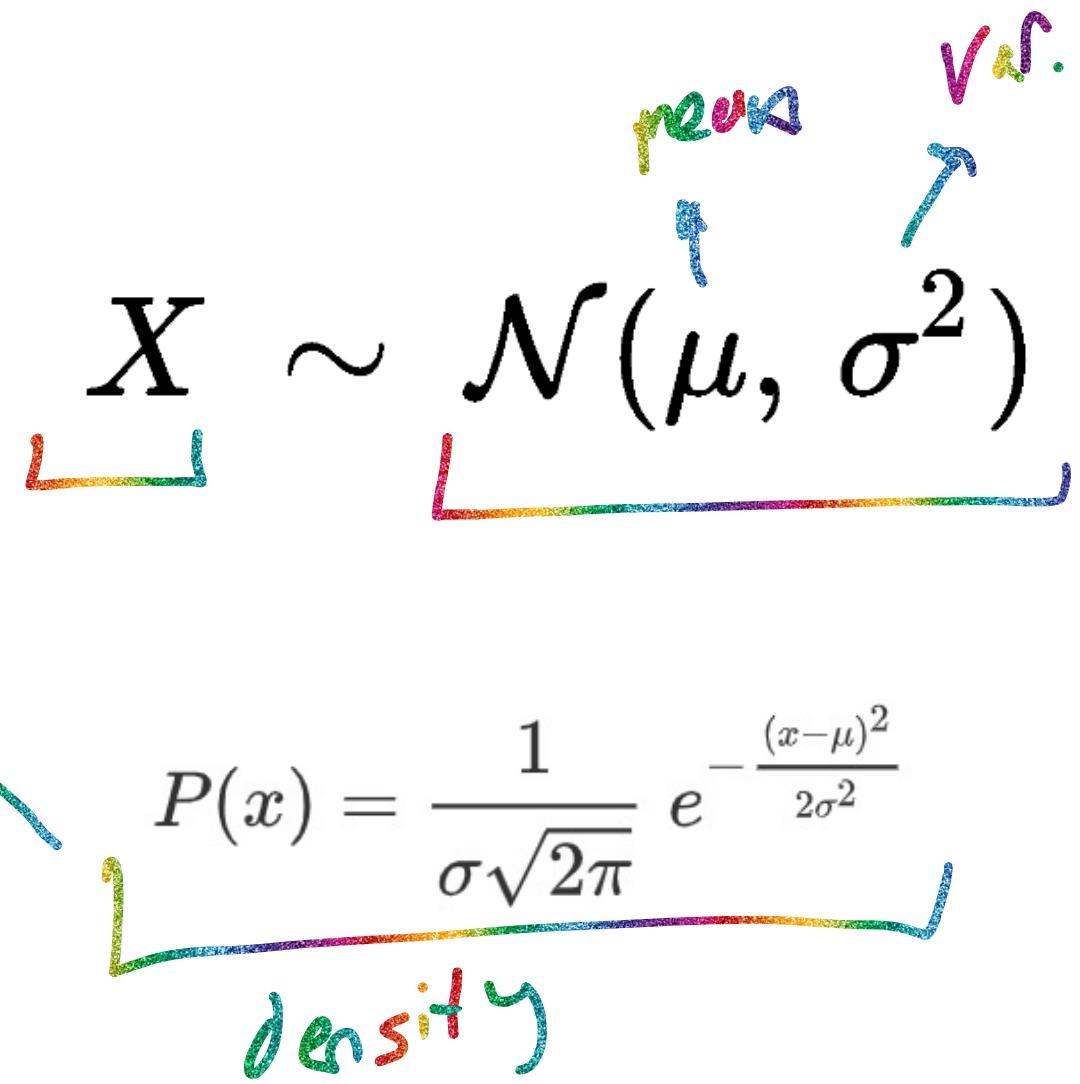
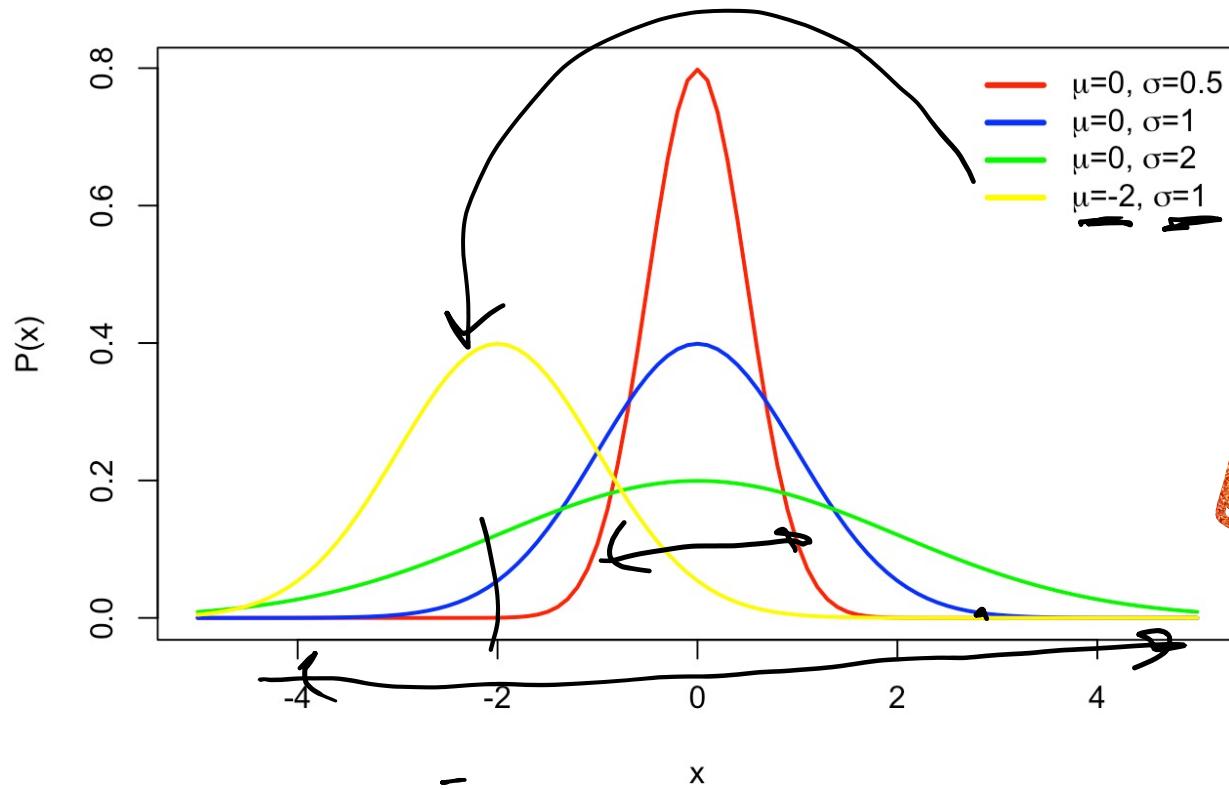
`var(x)`  
`sd(x)`  
`IQR(x)`

Sample  
cent. tend.  
spread

The diagram illustrates the use of R code to generate a sample and calculate descriptive statistics. It starts with the command `x=rnorm(20, mean=6, sd=0.7)`, where the value 20 is circled in red. An arrow points from this circled value to the word "sample". Below this, two functions are shown: `mean(x)` and `median(x)`, grouped by a bracket and labeled "cent. tend.". To the right of this bracket, another arrow points to the word "spread". At the bottom, three more functions are grouped by a bracket: `var(x)`, `sd(x)`, and `IQR(x)`. An arrow points from this bracket to the word "spread".

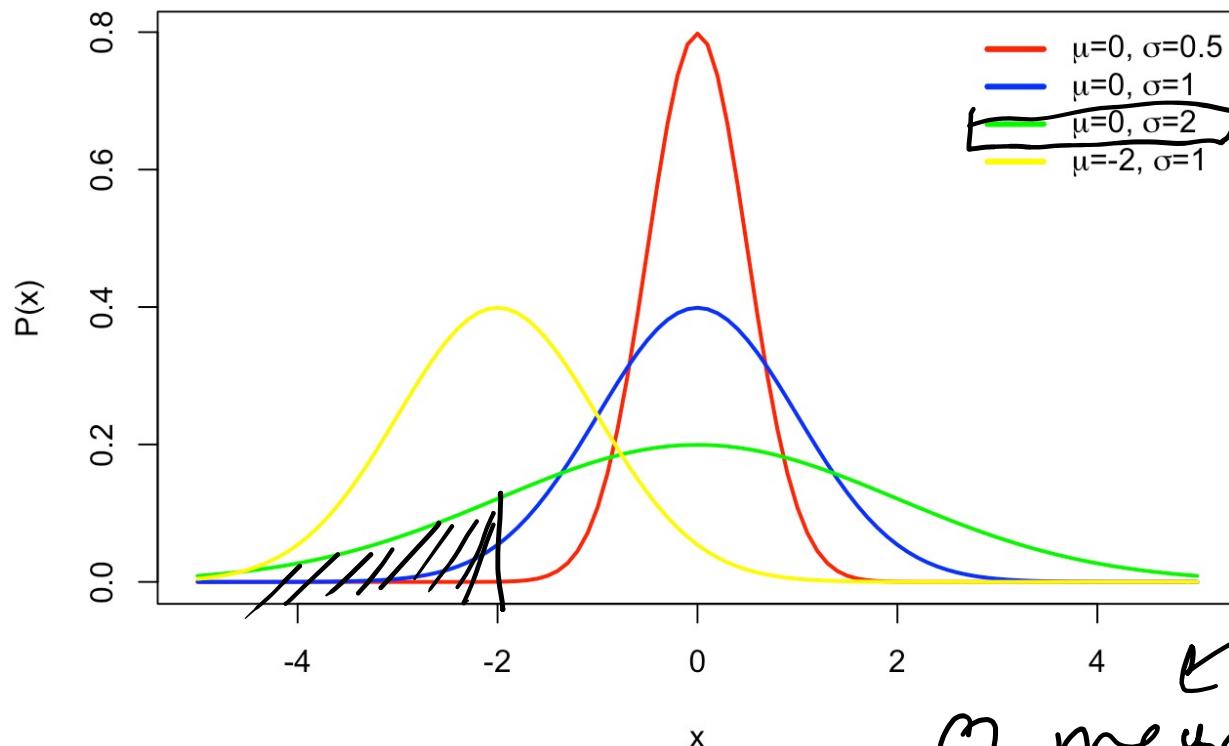
# Frequently used statistical distributions

## Normal distribution



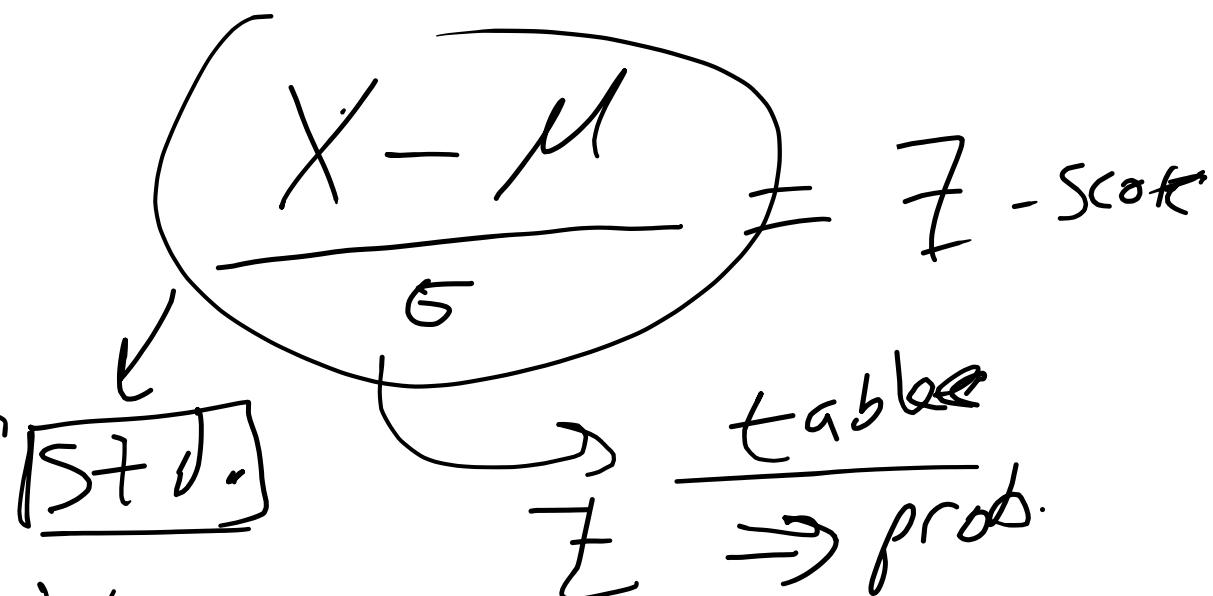
# Frequently used statistical distributions

## Normal distribution



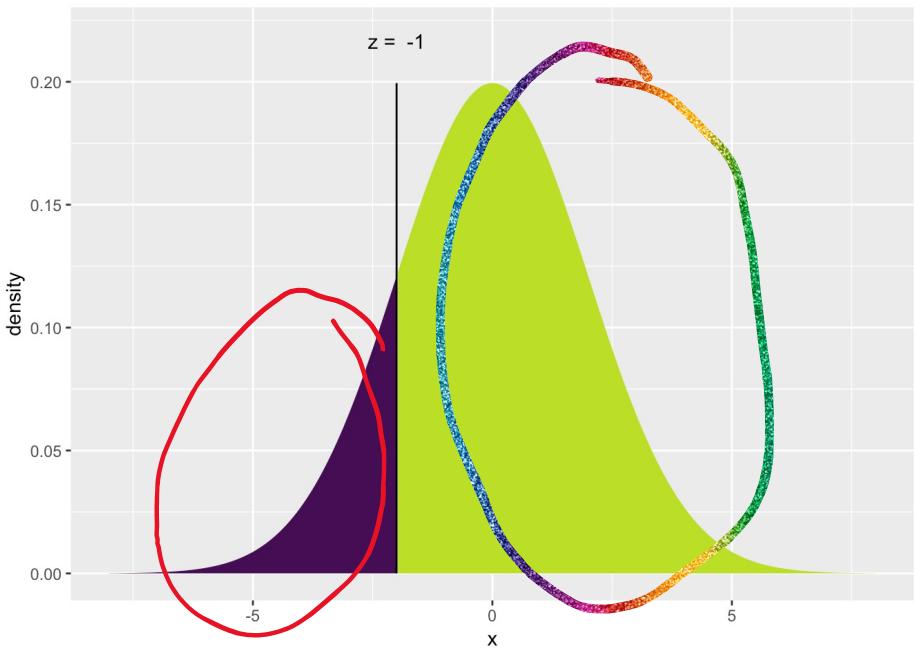
0 mean  
1 std. dev.

$$P(X \leq -2 | \mu = 0, \sigma = 2)$$



# Frequently used statistical distributions

## Normal distribution



Larea U. C  $\Rightarrow$  prob.

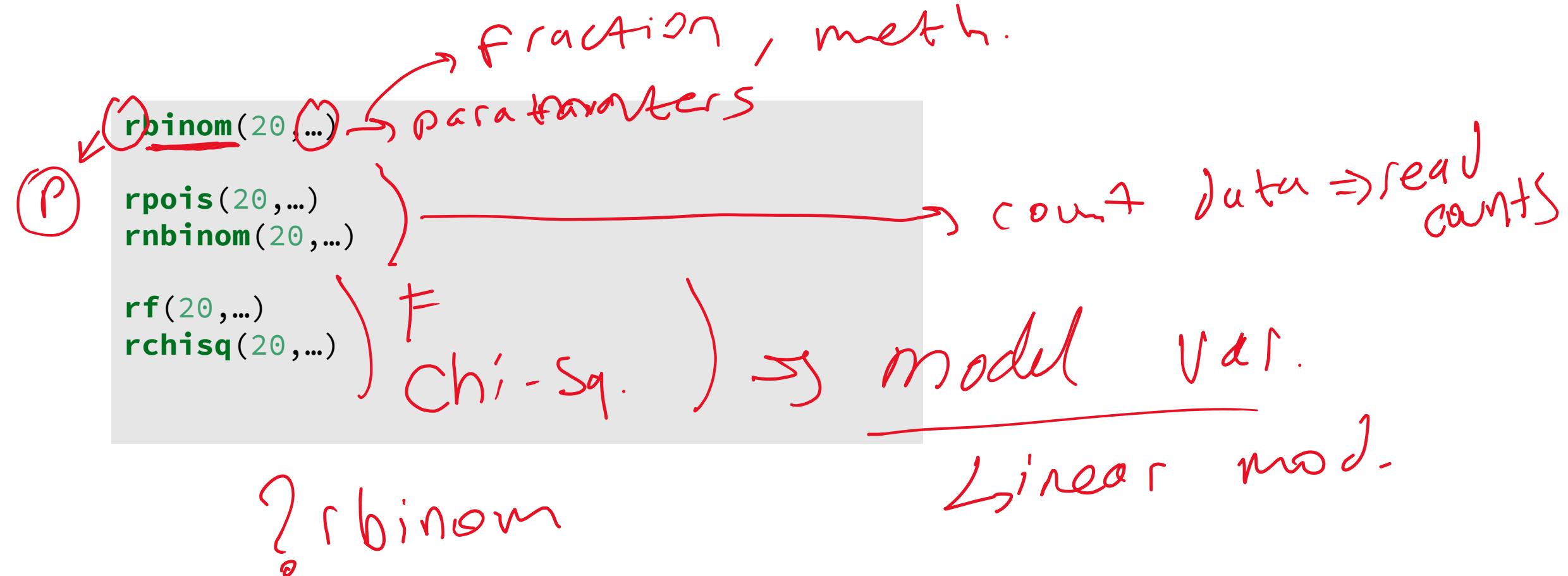
$$P(X \leq -2 \mid \mu = 0, \sigma = 2)$$

```
#get the probability of P(X <= -2) | mean=0 and sd=2  
pnorm(-2, mean=0, sd=2)  
  
#get the probability of P(X > -2) | mean=0 and sd=2  
pnorm(-2, mean=0, sd=2, lower.tail = FALSE)  
  
#get 5 random numbers from norm. dist. | mean=0 and sd=2  
rnorm(5, mean=0, sd=2)
```

↳ sample

# Frequently used statistical distributions

## Other distributions: R code



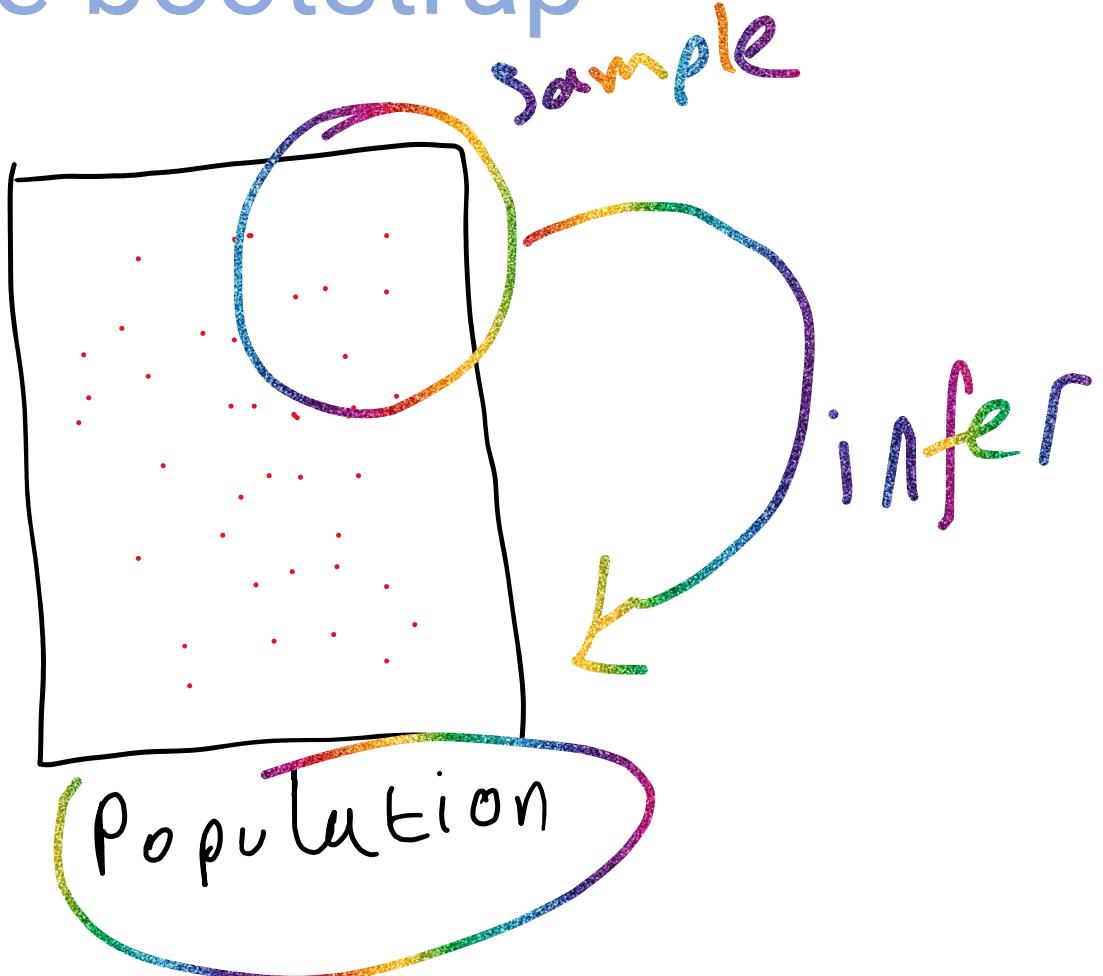
# Precision of estimates: confidence intervals

- How well a statistic estimated from a sample represents the population statistic ?

range for stat.

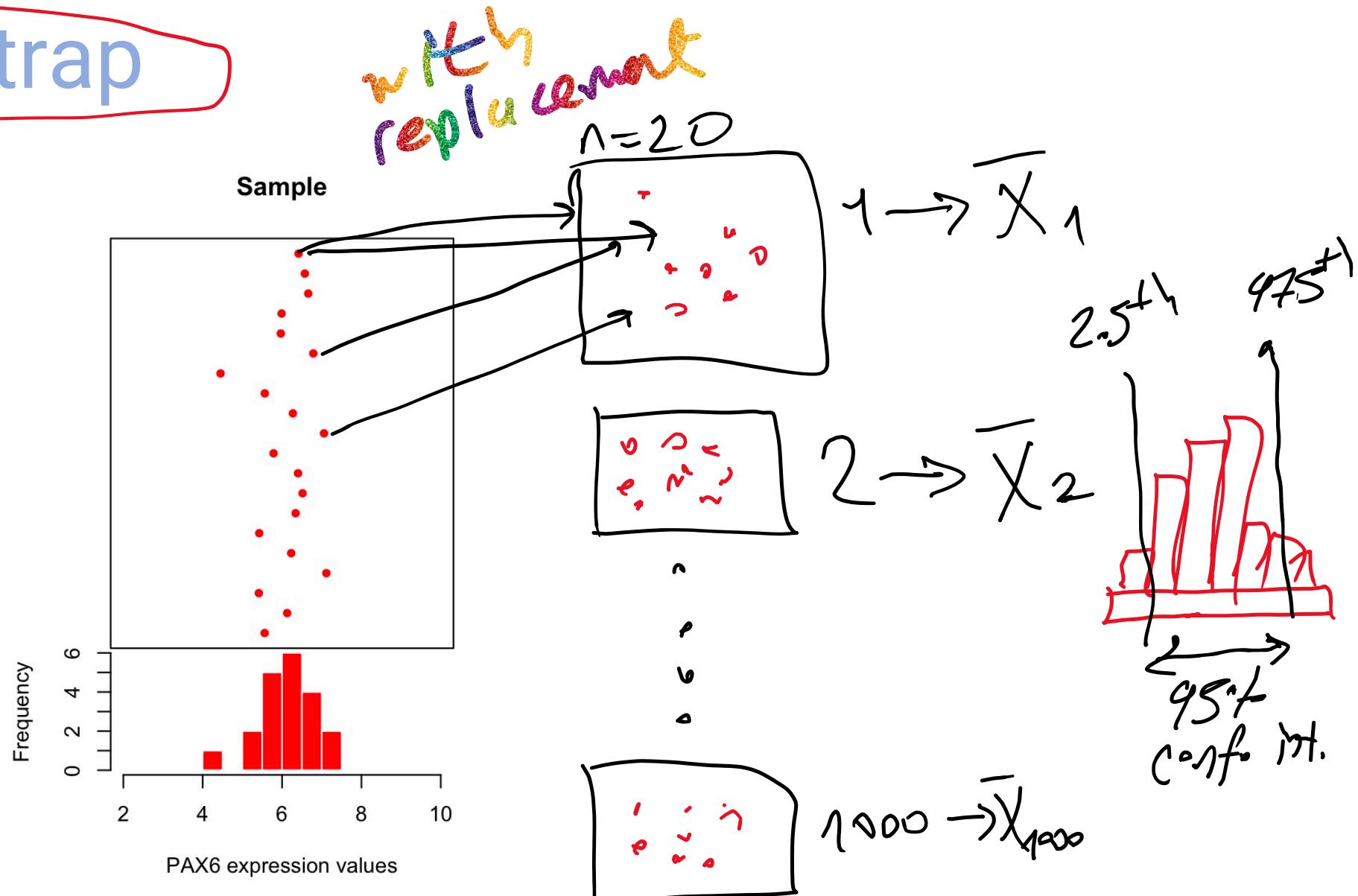
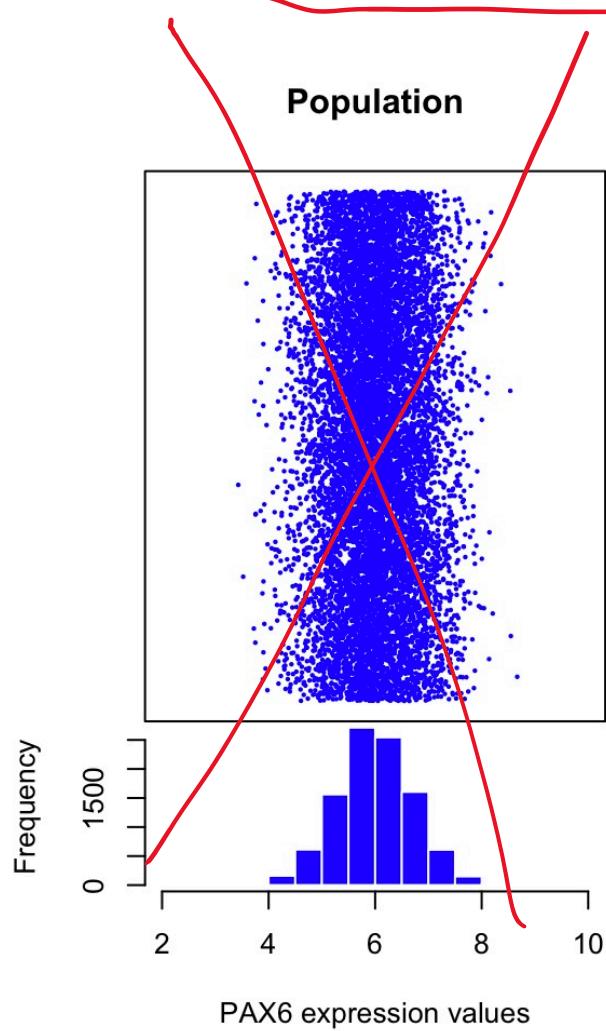
mean  
std. dev.

# Precision of estimates: The bootstrap

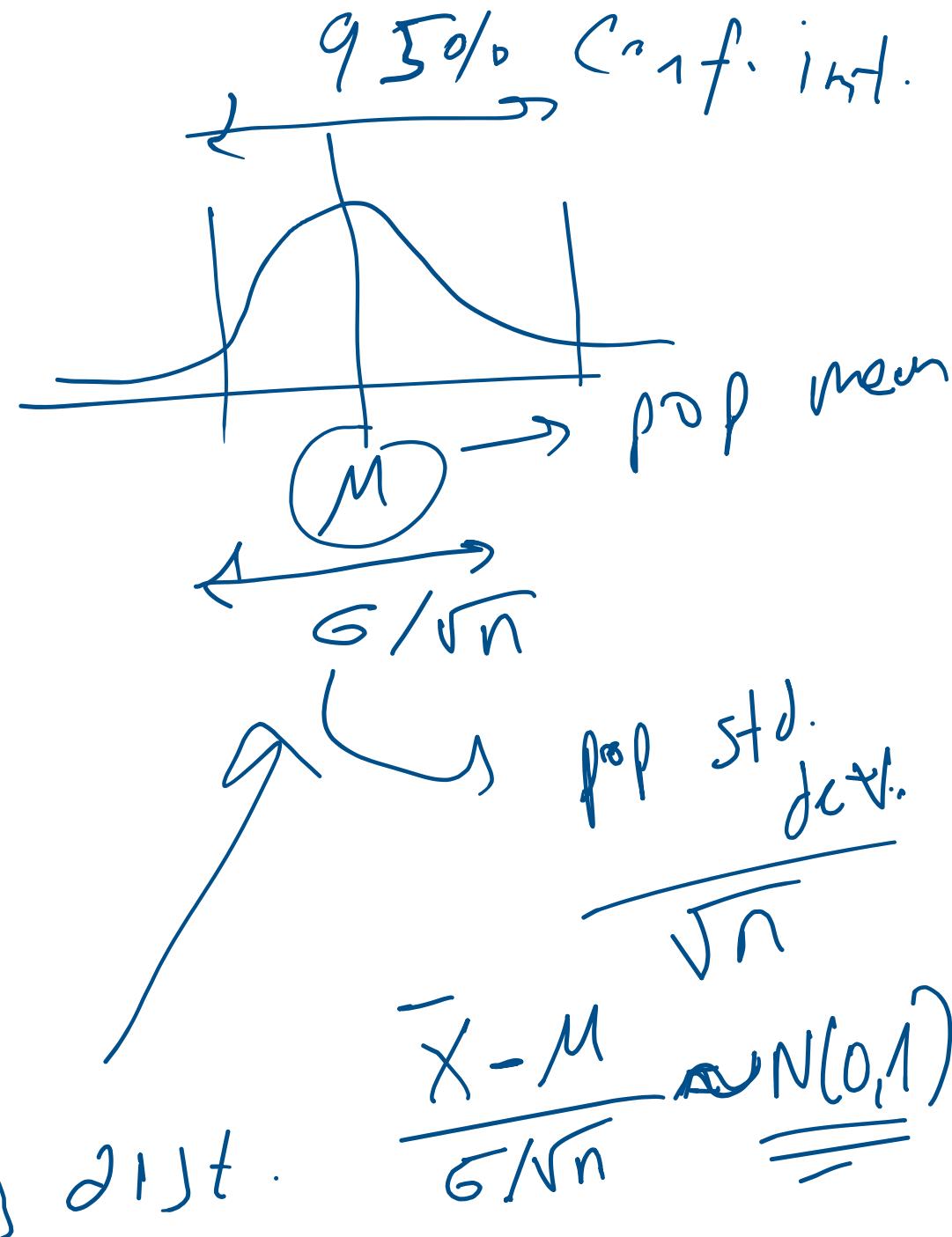
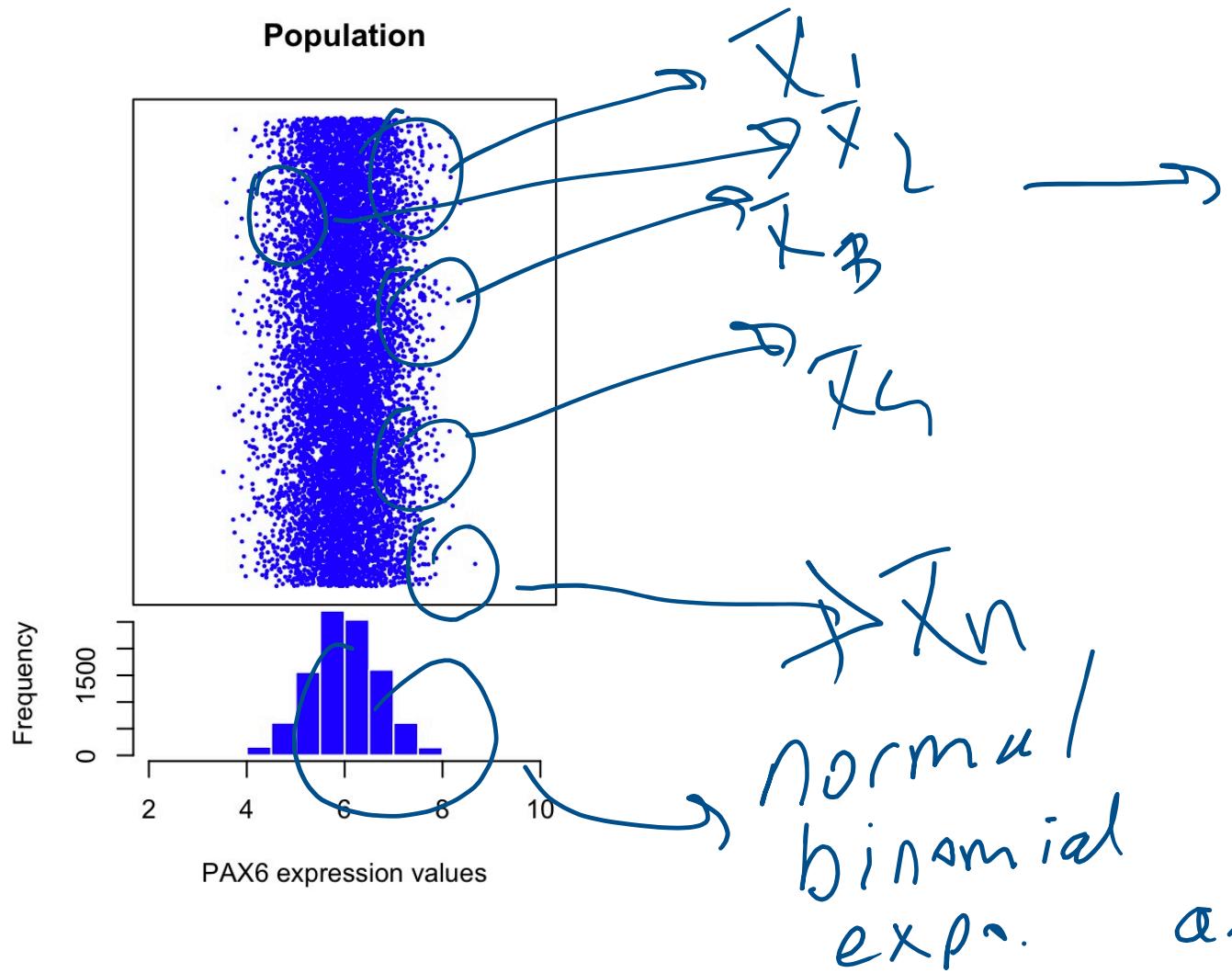


# Precision of estimates:

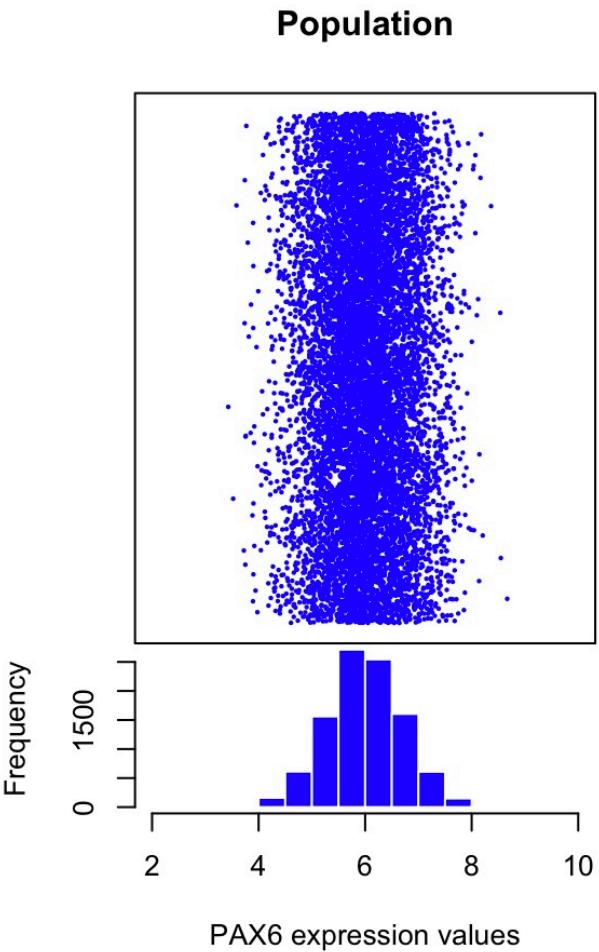
## The bootstrap



# Precision of estimates: Central Limit Theorem



# Precision of estimates: Central Limit Theorem



$$Z = \frac{\bar{X} - \mu}{\sigma / \sqrt{n}}$$

$P(-1.96 < Z < 1.96) = 0.95$

$P(-1.96 < \frac{\bar{X} - \mu}{\sigma / \sqrt{n}} < 1.96) = 0.95$

$P(\mu - 1.96\sigma / \sqrt{n} < \bar{X} < \mu + 1.96\sigma / \sqrt{n}) = 0.95$

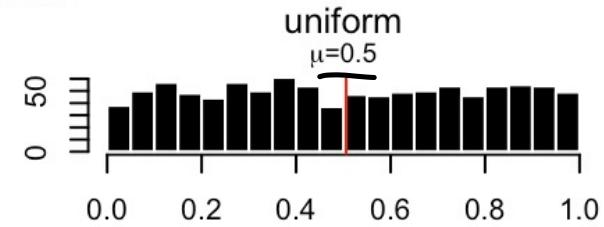
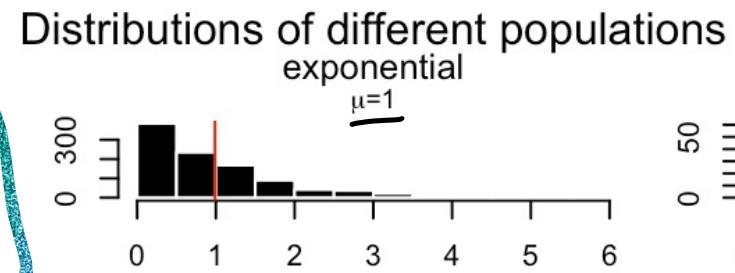
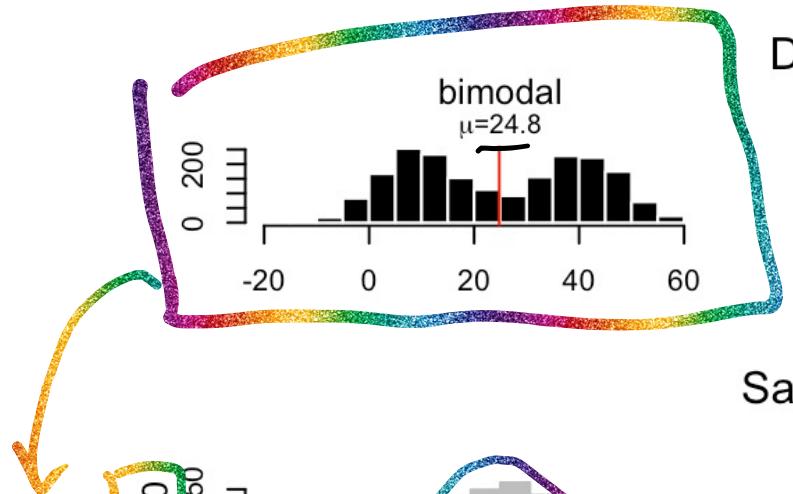
$\rightarrow P(\bar{X} - 1.96\sigma / \sqrt{n} < \bar{X} < \bar{X} + 1.96\sigma / \sqrt{n}) = 0.95$

$confint = [\bar{X} - 1.96\sigma / \sqrt{n}, \bar{X} + 1.96\sigma / \sqrt{n}]$

$90\%$

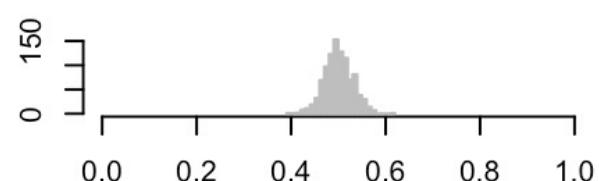
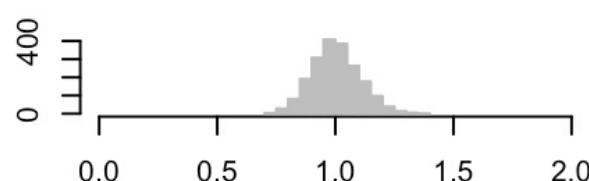
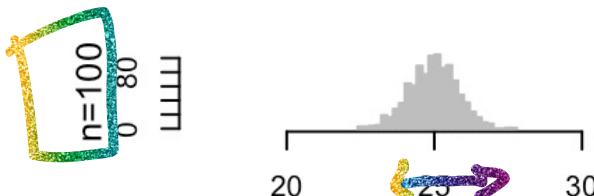
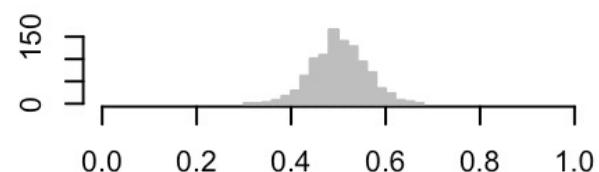
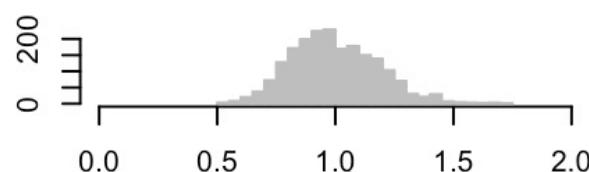
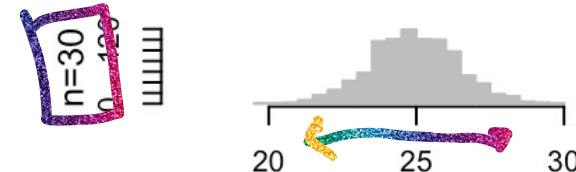
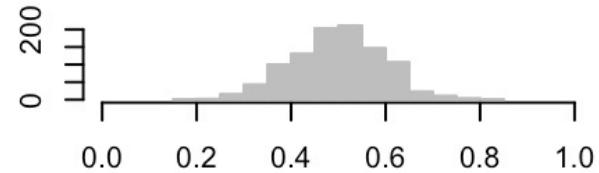
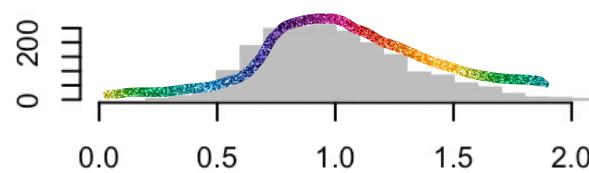
$\bar{X} \pm Z_{\alpha/2} \sigma / \sqrt{n}$

5th, 95th  $\Rightarrow N(0, 1)$



dist.  
me<sup>+</sup>(x)  
↓  
Normal

Sampling distribution of sample means



CLT problem: we need population std. dev.

$$\bar{X} \pm Z_{\alpha/2} \sigma / \sqrt{n}$$

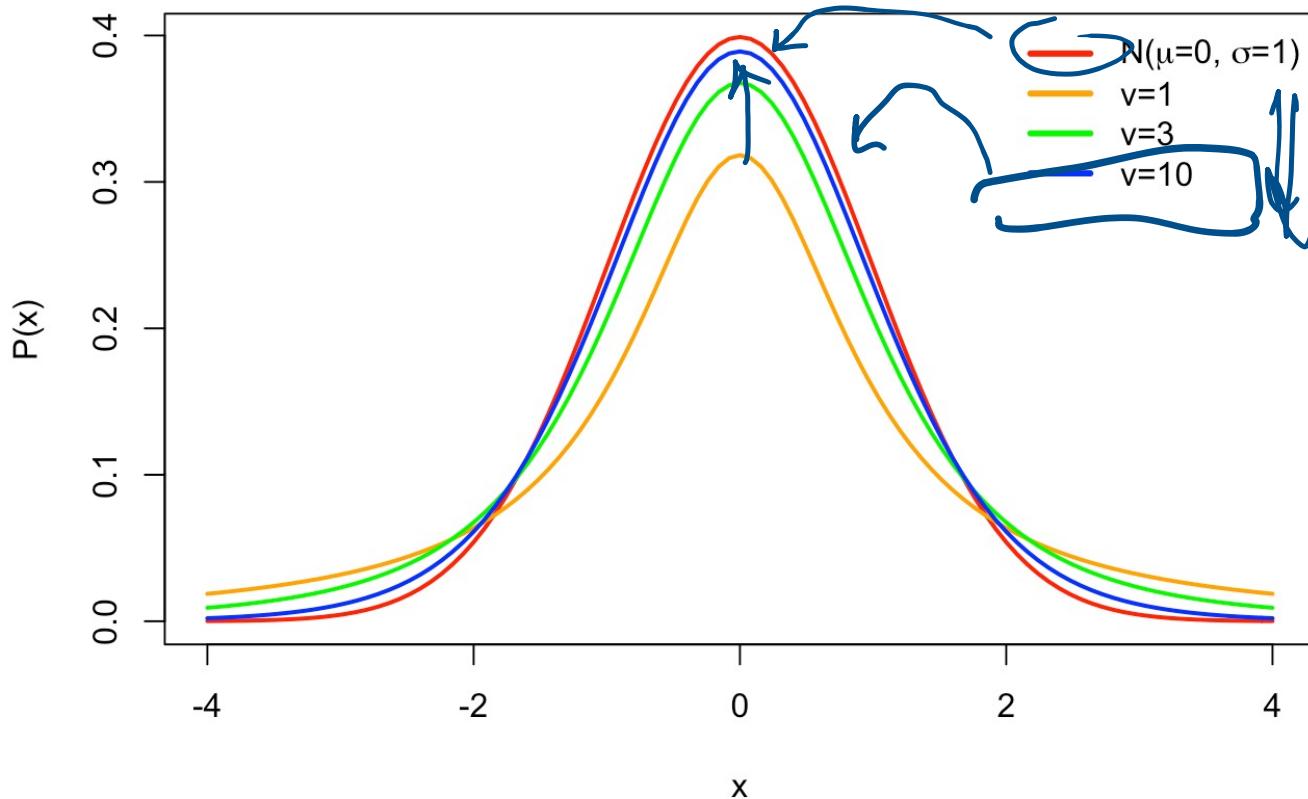
$$\frac{\bar{X} - \mu}{\sigma / \sqrt{n}} \sim N(0,1)$$

no idea

$$\bar{X} \pm t_{\alpha/2} s / \sqrt{n}, \text{ std dev}$$

$$\frac{\bar{X} - \mu}{s / \sqrt{n}} \sim t \text{ distribution}$$

# Why t distribution ?



$d.f = n-1$   
# points - 1

# How to test for differences between samples:

## Hypothesis testing

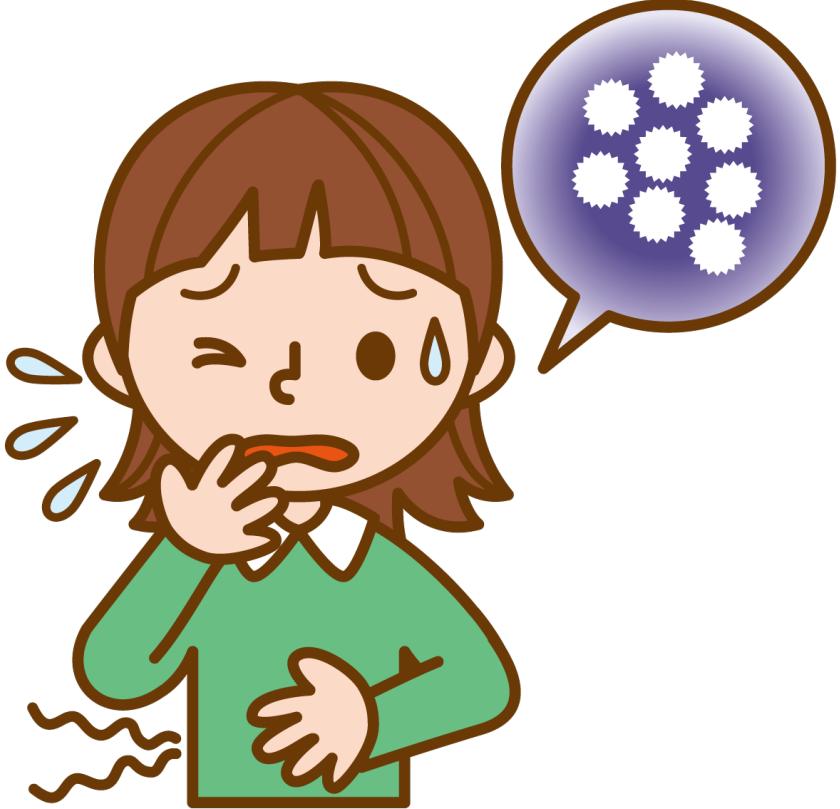
# How to test for differences between samples:

## Hypothesis testing



Healthy

VS.



Sick

# How to test for differences between samples:

## Hypothesis testing



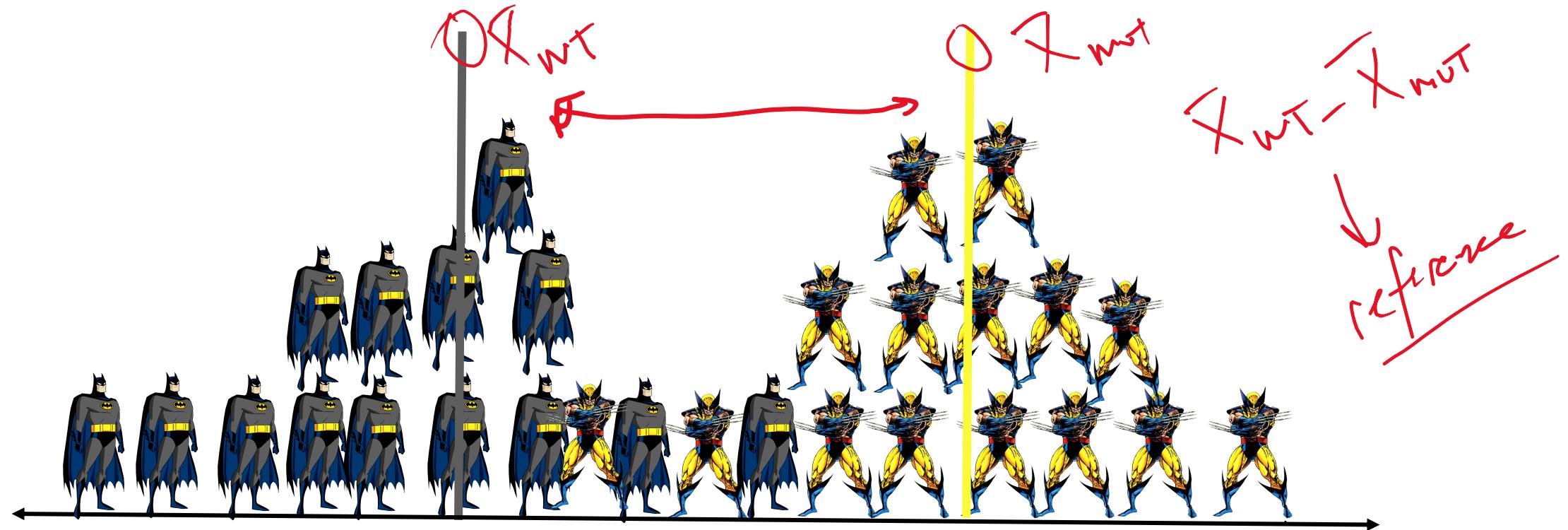
Wild-type (WT)



Mutant (Mut)

# How to test for differences between samples:

Hypothesis testing: Difference between the means

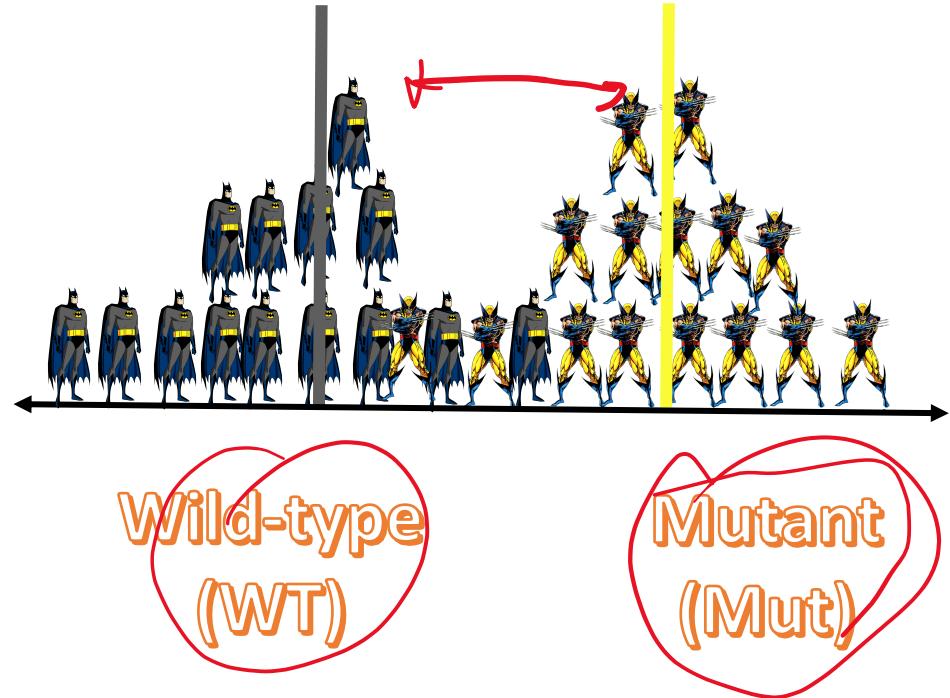


Wild-type (WT)

Mutant (Mut)

# How to test for differences between samples:

## Hypothesis testing: Difference between the means

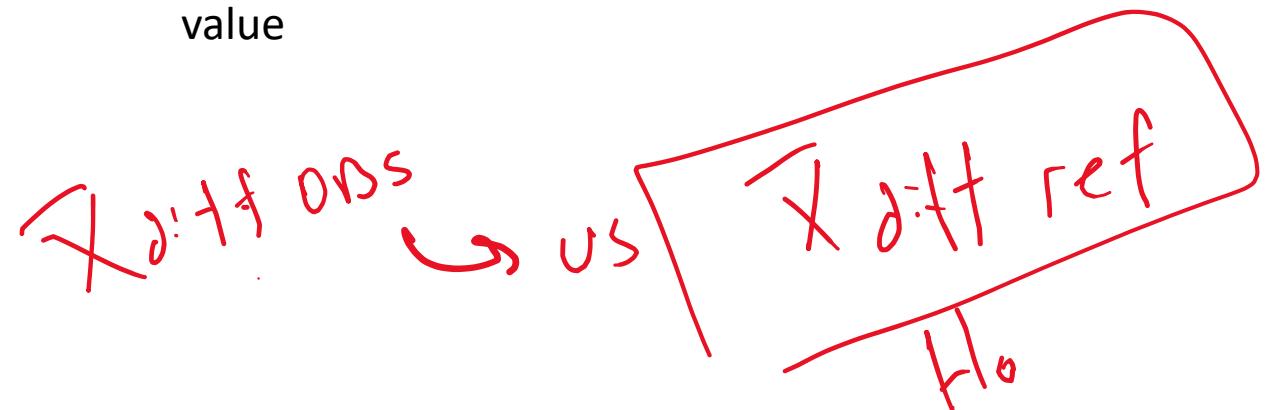


### Core ideas:

Decide on a hypothesis to test, often called “null hypothesis”  $H_0$        $H_1$

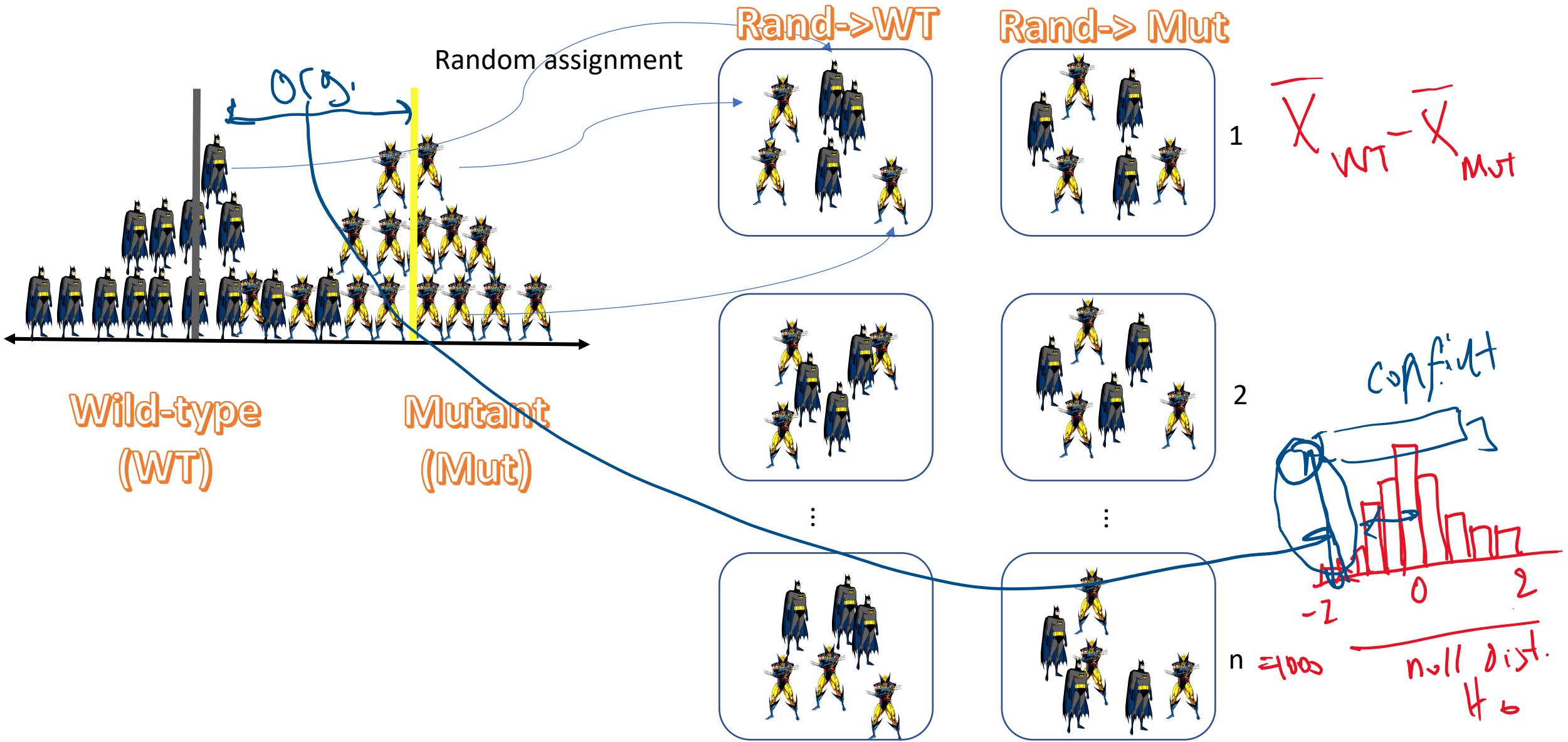
Decide on a statistic to test the truth of the null hypothesis.

Calculate the statistics and compare to a reference value



# How to test for differences between samples:

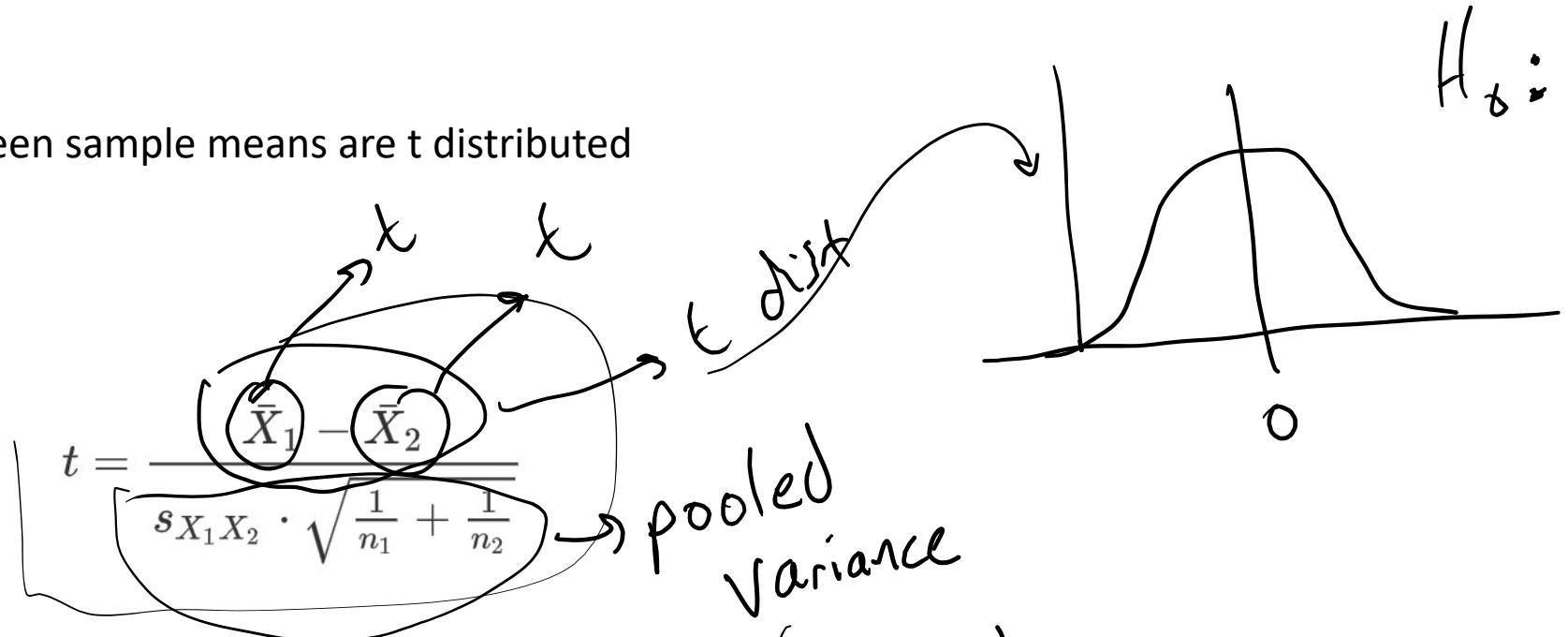
## Randomization based test



# How to test for differences between samples:

## Hypothesis testing: t-test

The differences between sample means are t distributed



$$s_{\bar{X}_1 \bar{X}_2} = \sqrt{\frac{(n_1 - 1)s_{\bar{X}_1}^2 + (n_2 - 1)s_{\bar{X}_2}^2}{n_1 + n_2 - 2}}$$

Annotations explain the components:

- "assume eq. variance" points to the term  $(n_1 - 1)s_{\bar{X}_1}^2 + (n_2 - 1)s_{\bar{X}_2}^2$  under the square root.

# How to test for differences between samples:

## Hypothesis testing: t-test

Unequal variance assumption: Welch's t-test

Default t-test

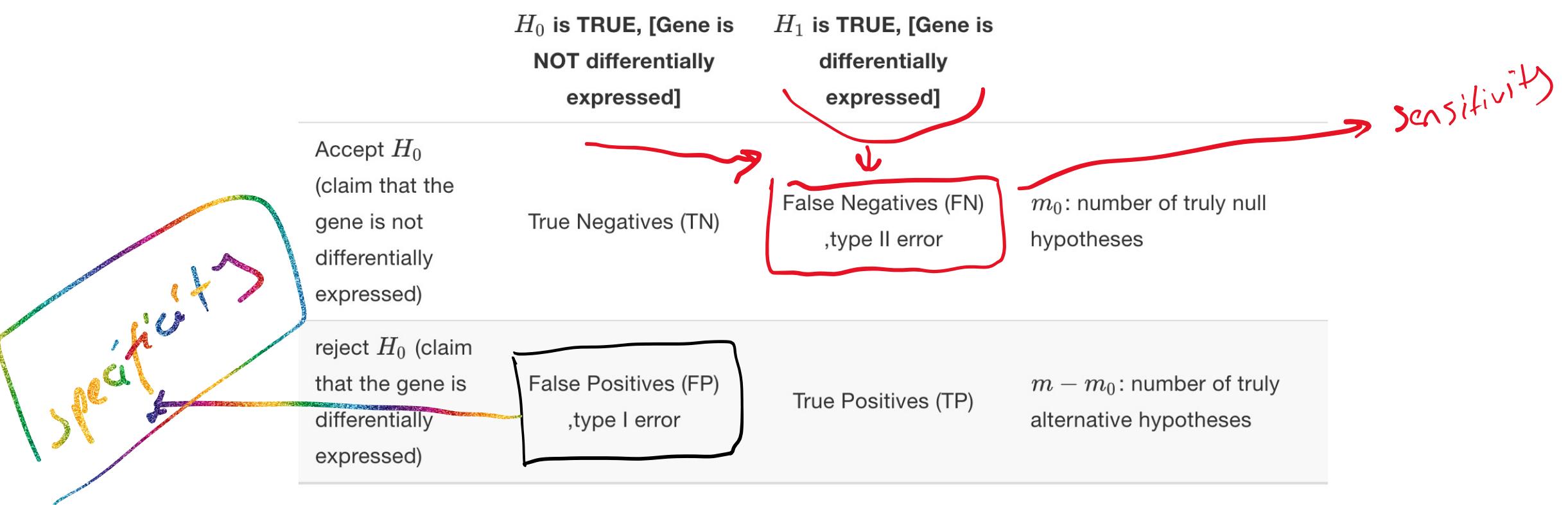
$$t = \frac{\bar{X}_1 - \bar{X}_2}{s_{\bar{X}_1 - \bar{X}_2}}$$

$$s_{\bar{X}_1 - \bar{X}_2} = \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}$$

$$\text{d.f.} = \frac{(s_1^2/n_1 + s_2^2/n_2)^2}{(s_1^2/n_1)^2/(n_1 - 1) + (s_2^2/n_2)^2/(n_2 - 1)}$$

# How to test for differences between samples:

## Hypothesis testing: multiple testing correction



# How to test for differences between samples:

## Hypothesis testing: multiple testing correction

General idea:

- more FPs as we test more ↑
- Push P-values to higher values to avoid making these mistakes

Bonferroni

$$\frac{m \cdot p_i}{\# \text{ tests}}$$

False Discovery Rate

$$\frac{m \cdot p_i}{i}$$

Q values

Depends on estimated proportion of True null hypothesis

# How to test for differences between samples:

## Hypothesis testing: moderated t-tests

General idea: Use information from the many tests you are doing

Doing many tests is typical procedure for gene expression analysis

Borrow information for variation estimates

$$\hat{V}_g = aV_0 + bV_g$$

Annotations in red:

- Above the equation:  $\frac{1}{2} = a - b$
- To the left of  $\hat{V}_g$ : Shrunken var.
- To the right of  $V_0$ : gene var.
- To the right of  $bV_g$ : background
- To the right of the entire equation: other var. est.

# How to test for differences between samples:

## Hypothesis testing: moderated t-tests

Moderated tests in action: a simple example

- Simulate 1000 genes, 3 test and 3 control groups
- No real difference in any of the genes
- Observe the effect of moderated and unmoderated tests

$$\hat{V}_g = (V_0 + V_g)/2$$

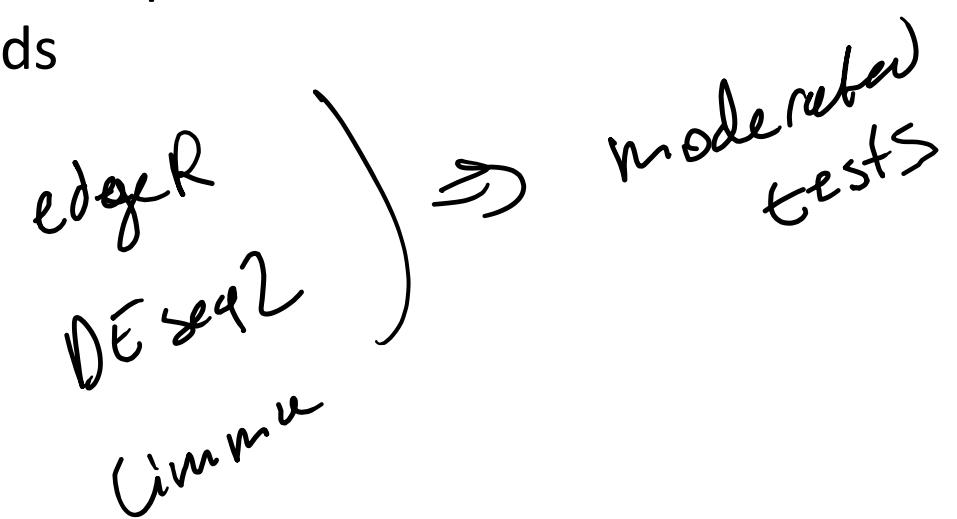
# How to test for differences between samples:

## Hypothesis testing: moderated t-tests

Why is this important? moderation of parameter estimates are used by modern diff. exp. methods

edgeR  
DEseq2  
Cimme

moderated tests



# Relationship between variables:

## Linear regression

We often need to model the relationship between variables

$$Y = \beta_0 + \beta_1 X + \epsilon$$

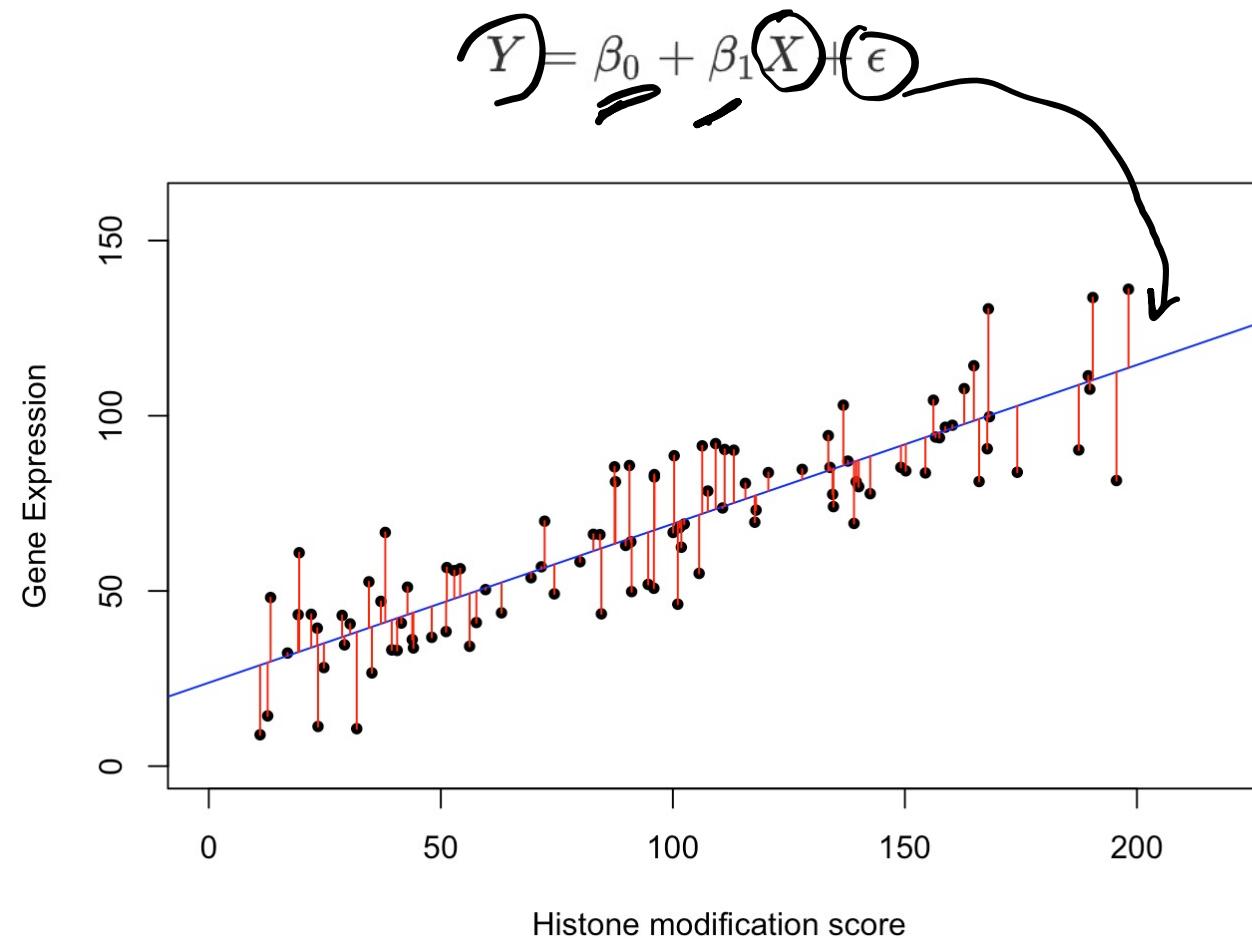
Dependent variable  
Response variable

Independent variable  
Explanatory variables  
Predictor variable

The diagram illustrates the components of a linear regression equation. The dependent variable (Y) is shown as the outcome being predicted. The independent variable (X) is shown as the factor used to predict the outcome. The error term ( $\epsilon$ ) represents the difference between the observed value and the predicted value. The coefficients ( $\beta_0$  and  $\beta_1$ ) represent the intercept and slope of the regression line, respectively.

# Relationship between variables:

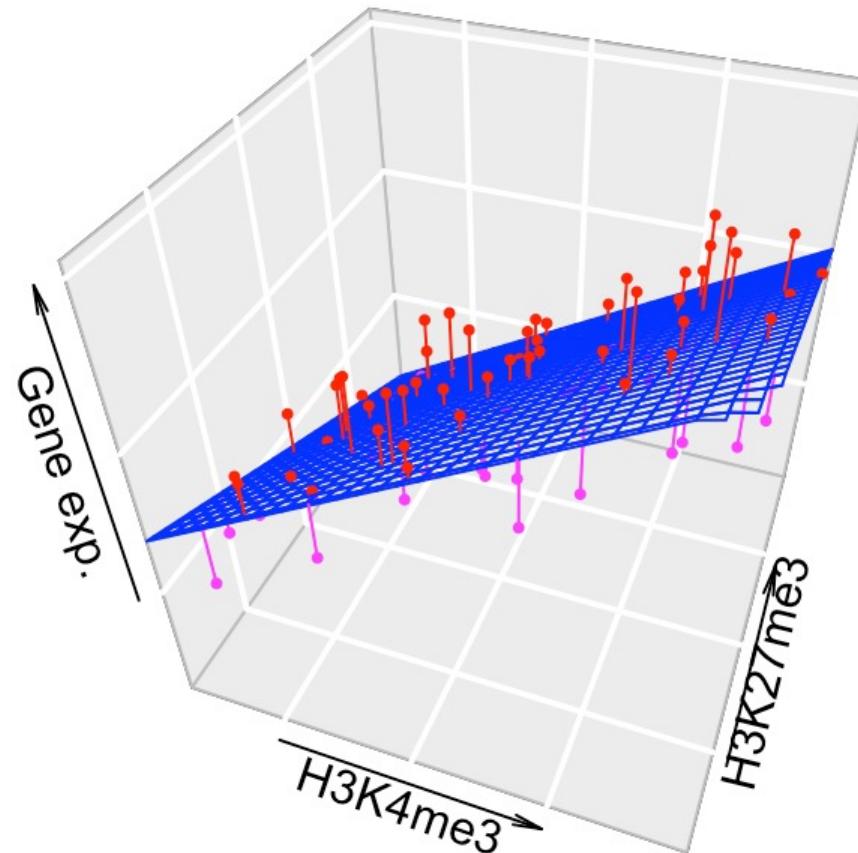
## Linear regression



# Relationship between variables:

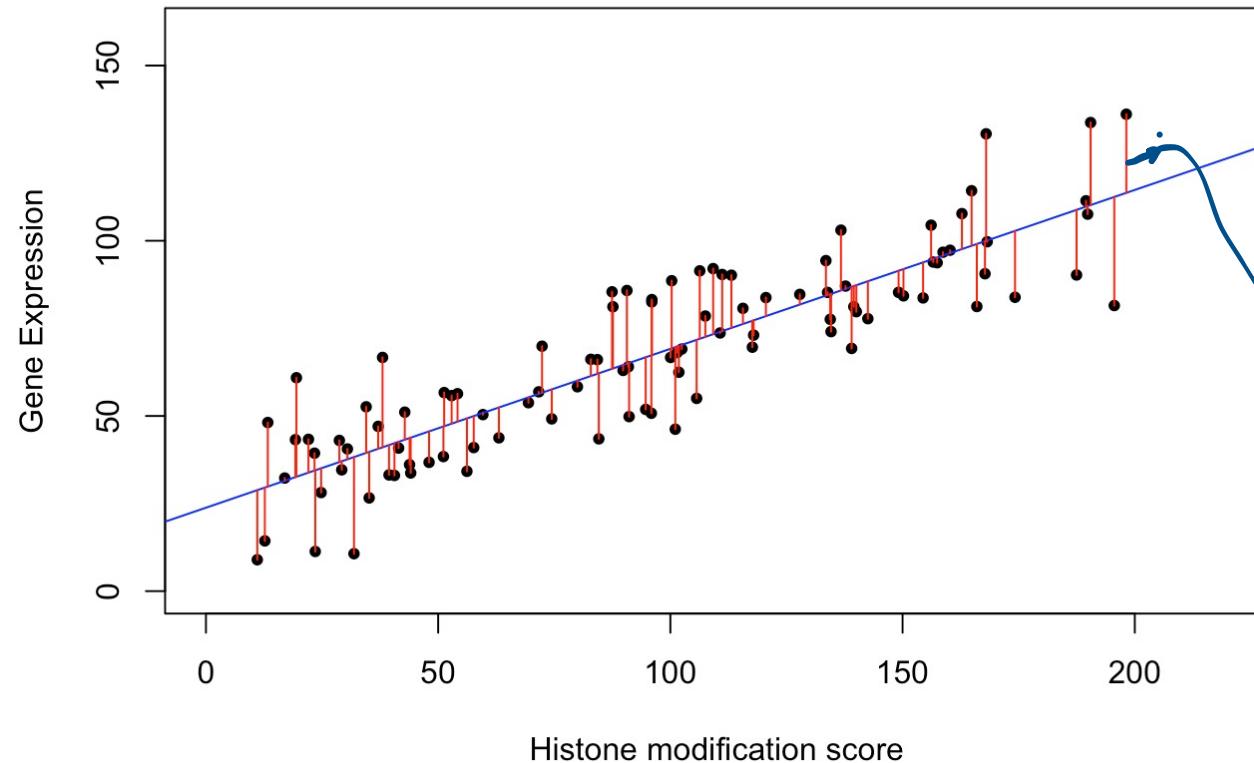
## Linear regression

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \epsilon$$



# Relationship between variables: How to fit a line ?

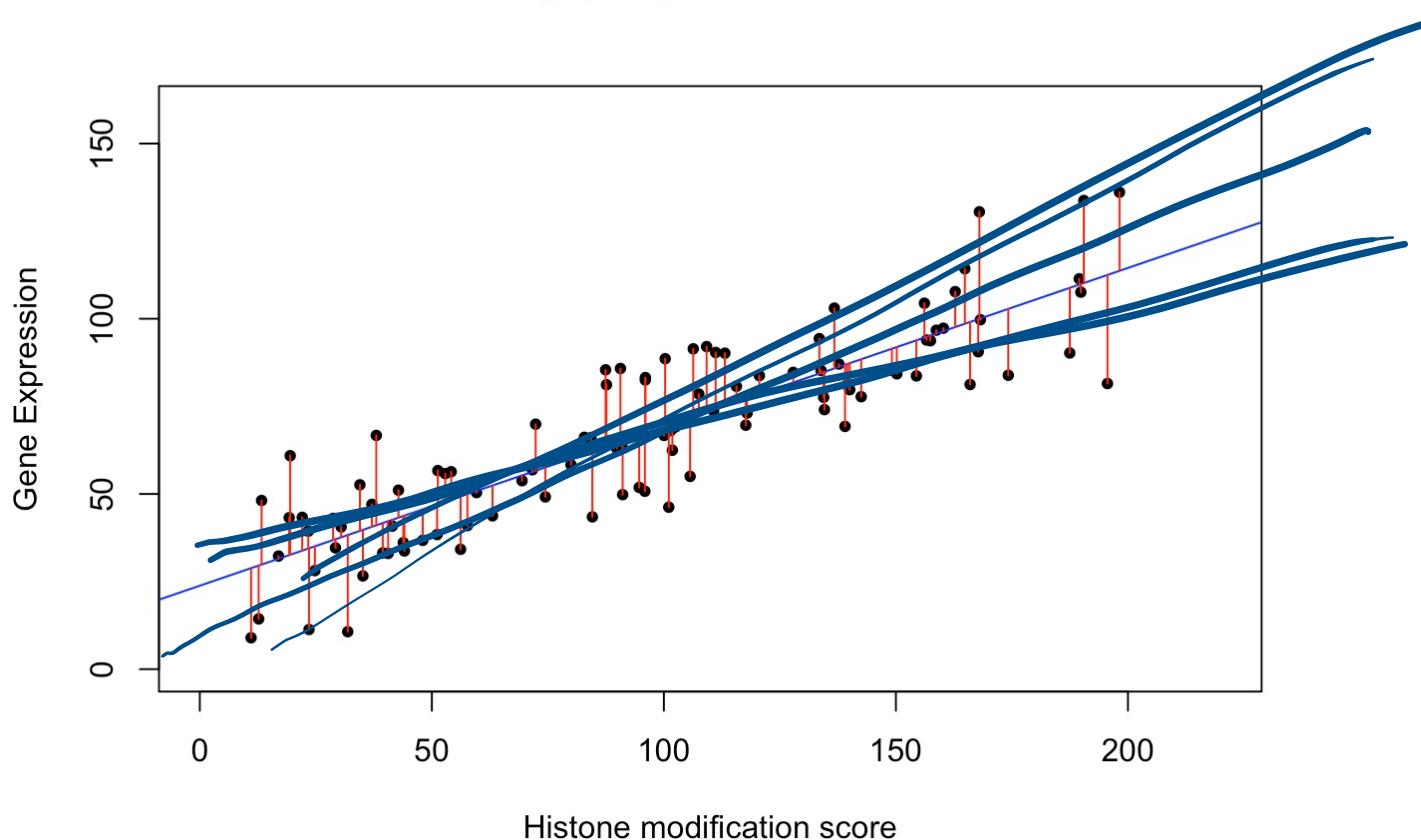
$$Y = \beta_0 + \beta_1 X + \epsilon$$



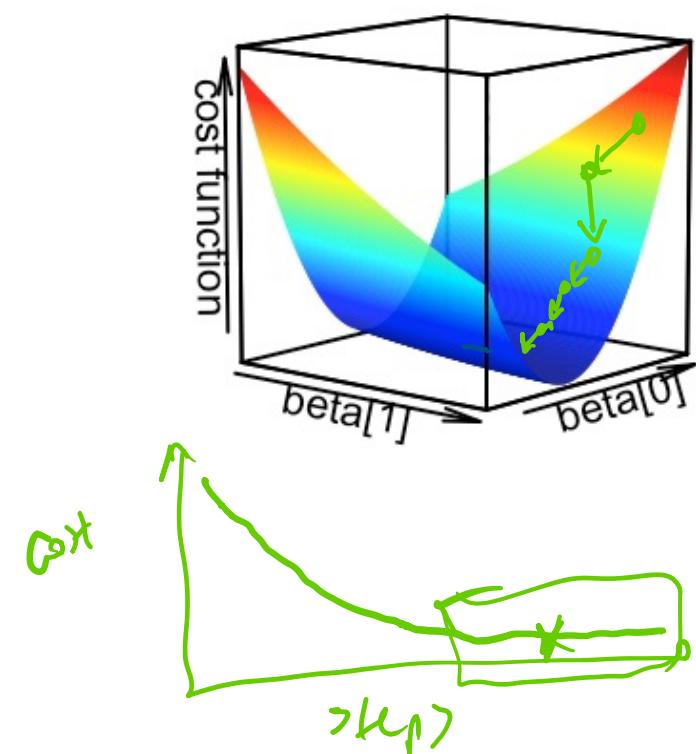
$$\min \sum (y_i - (\beta_0 + \beta_1 x_i))^2$$

# Relationship between variables: How to fit a line ? -> cost or loss function

$$Y = \beta_0 + \beta_1 X + \epsilon$$

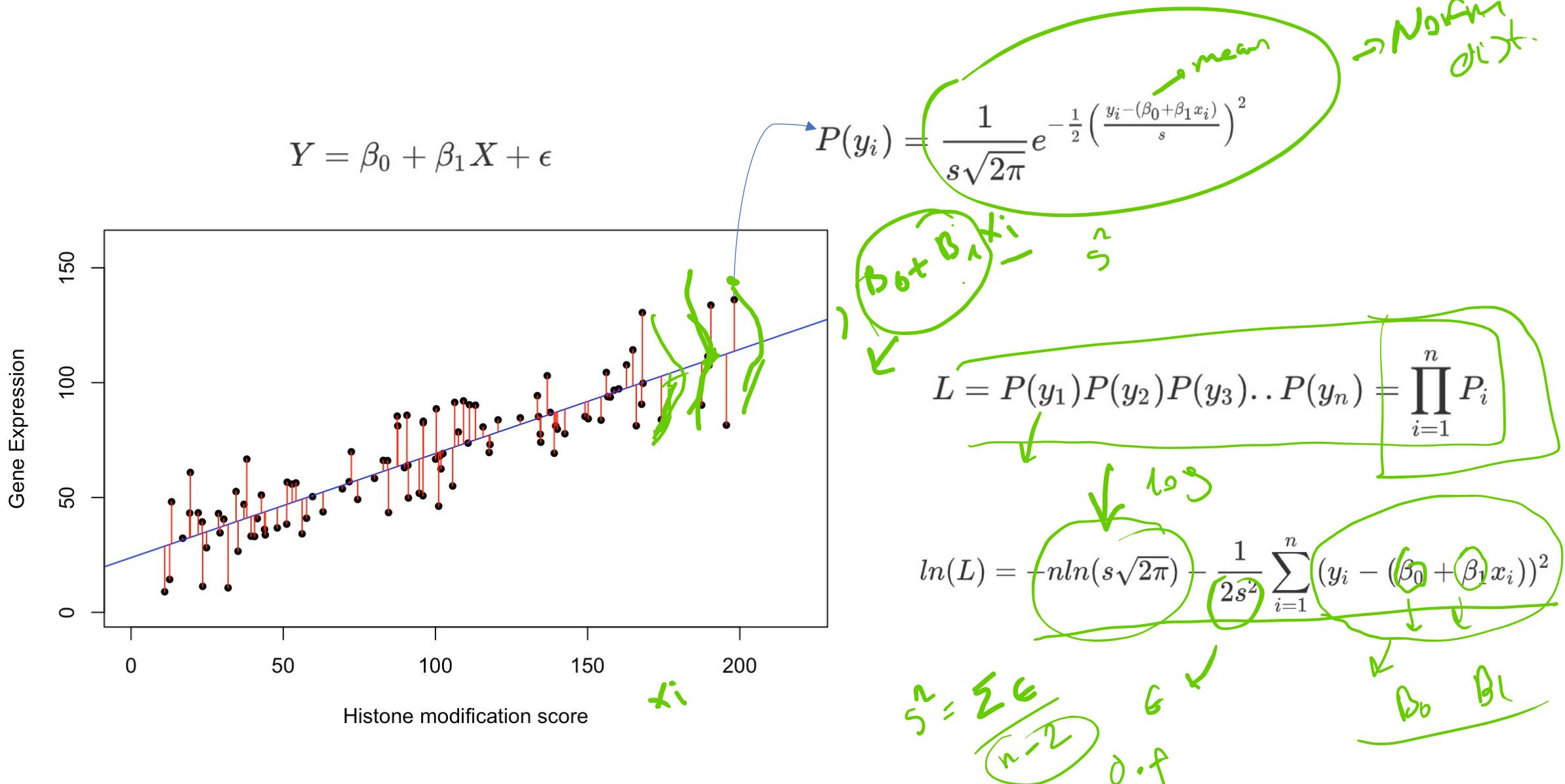


$$\min \sum (y_i - (\beta_0 + \beta_1 x_i))^2$$



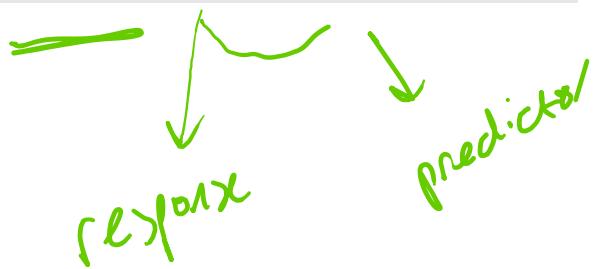
# Relationship between variables:

## How to fit a line? -> maximum likelihood function

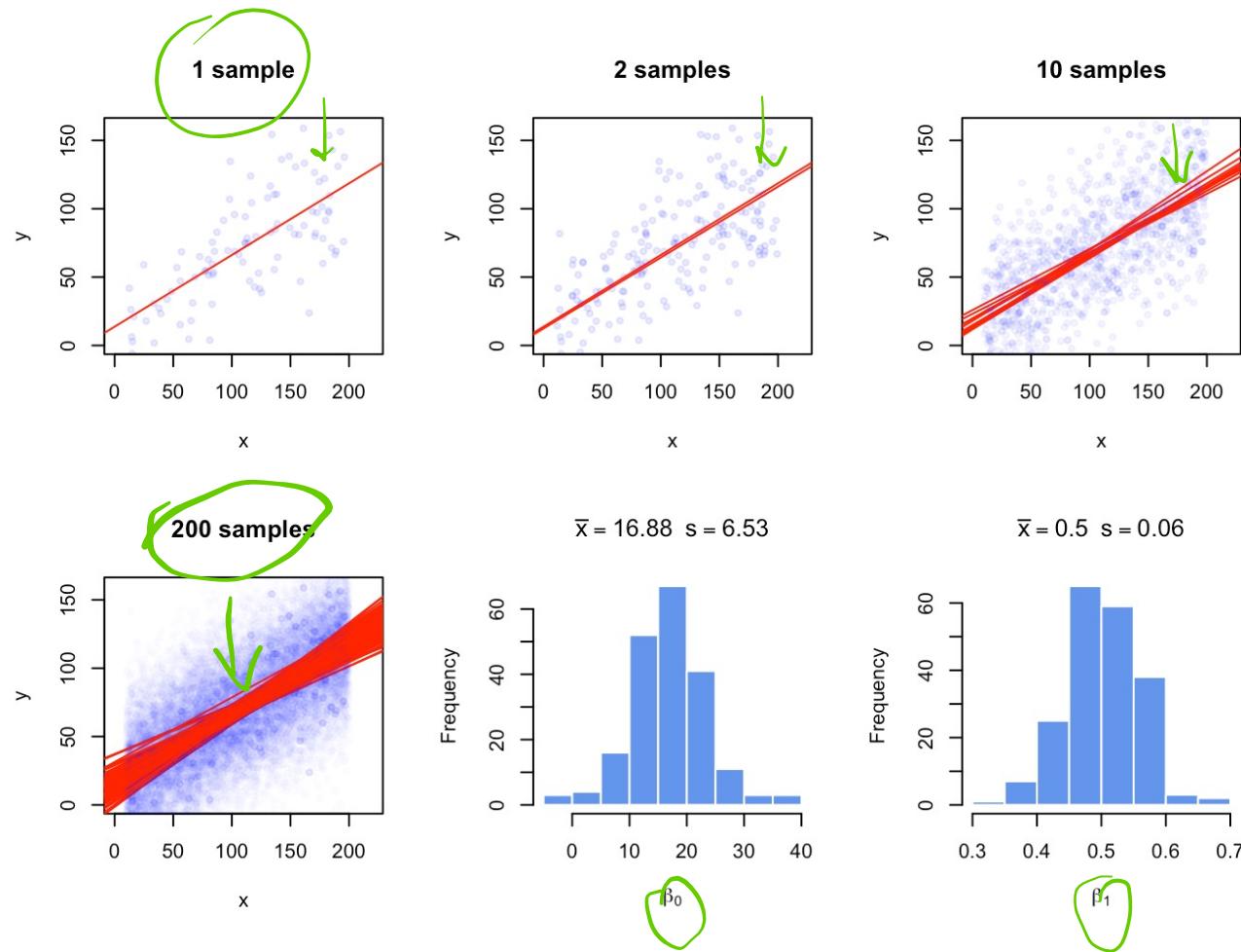


# Relationship between variables: Fitting lines in R

```
mod1=lm(y~x)
```



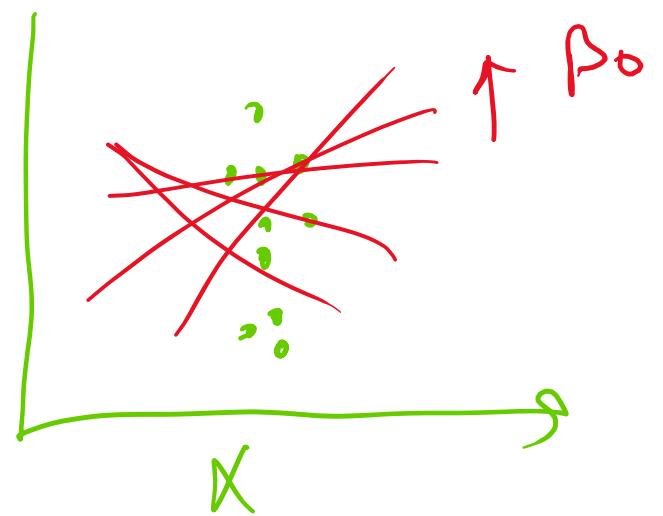
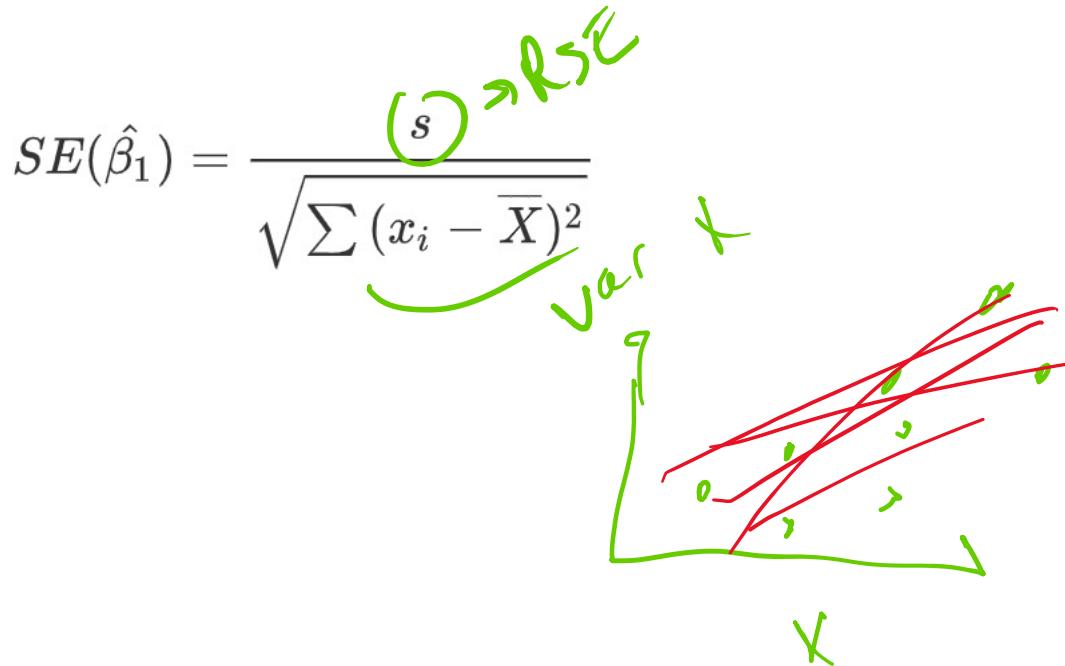
# Relationship between variables: Error of the coefficients



# Relationship between variables:

## Error of the coefficients

$$s = RSE = \sqrt{\frac{\sum (y_i - (\beta_0 + \beta_1 x_i))^2}{n-2}} = \sqrt{\frac{\sum \epsilon^2}{n-2}} \Rightarrow RSE$$



# Relationship between variables: Accuracy of the model

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	90.40951	6.43208	14.06	<2e-16 ***
x	2.08742	0.09775	21.36	<2e-16 ***
---				
Signif. codes:	0 ‘***’ 0.001 ‘**’ 0.01 ‘*’ 0.05 ‘.’ 0.1 ‘ ’ 1			

Residual standard error: 16.96 on 48 degrees of freedom

Multiple R-squared: 0.9048, Adjusted R-squared: 0.9028

F-statistic: 456 on 1 and 48 DF, p-value: < 2.2e-16

$$s = RSE = \sqrt{\frac{\sum (y_i - \hat{Y}_i)^2}{n - p}} = \sqrt{\frac{RSS}{n - p}}$$

# Relationship between variables:

## Accuracy of the model

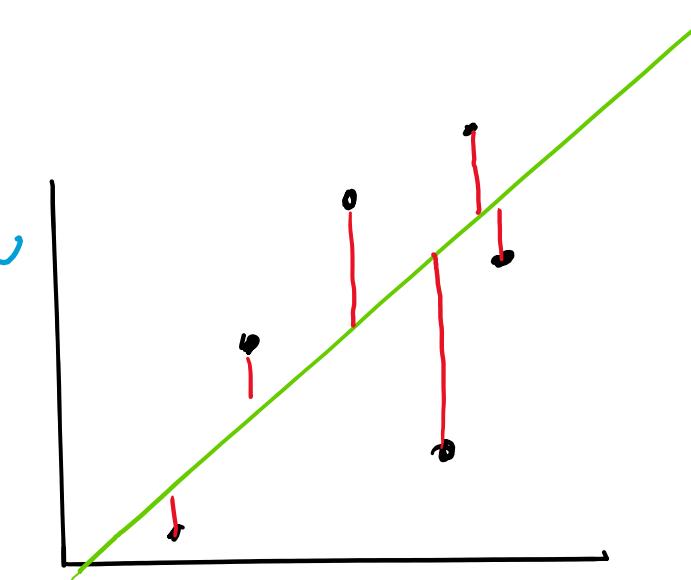
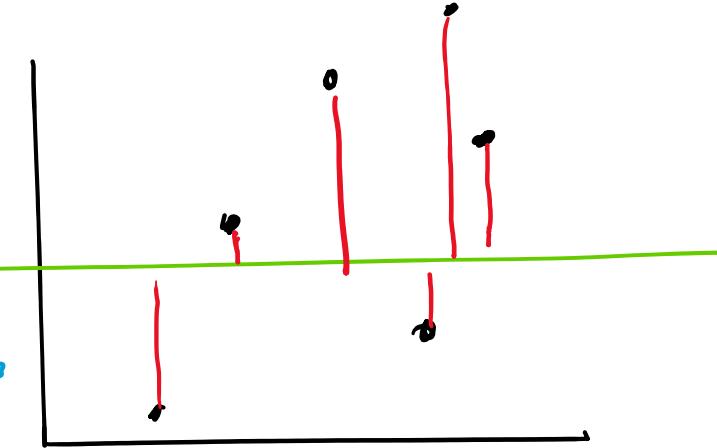
$$R^2 = 1 - \frac{RSS}{TSS} = \frac{TSS - RSS}{TSS} = 1 - \frac{RSS}{TSS}$$

$\sum (y_i - \hat{Y}_i)^2$

(0, 1)

Explained variability : R-squared

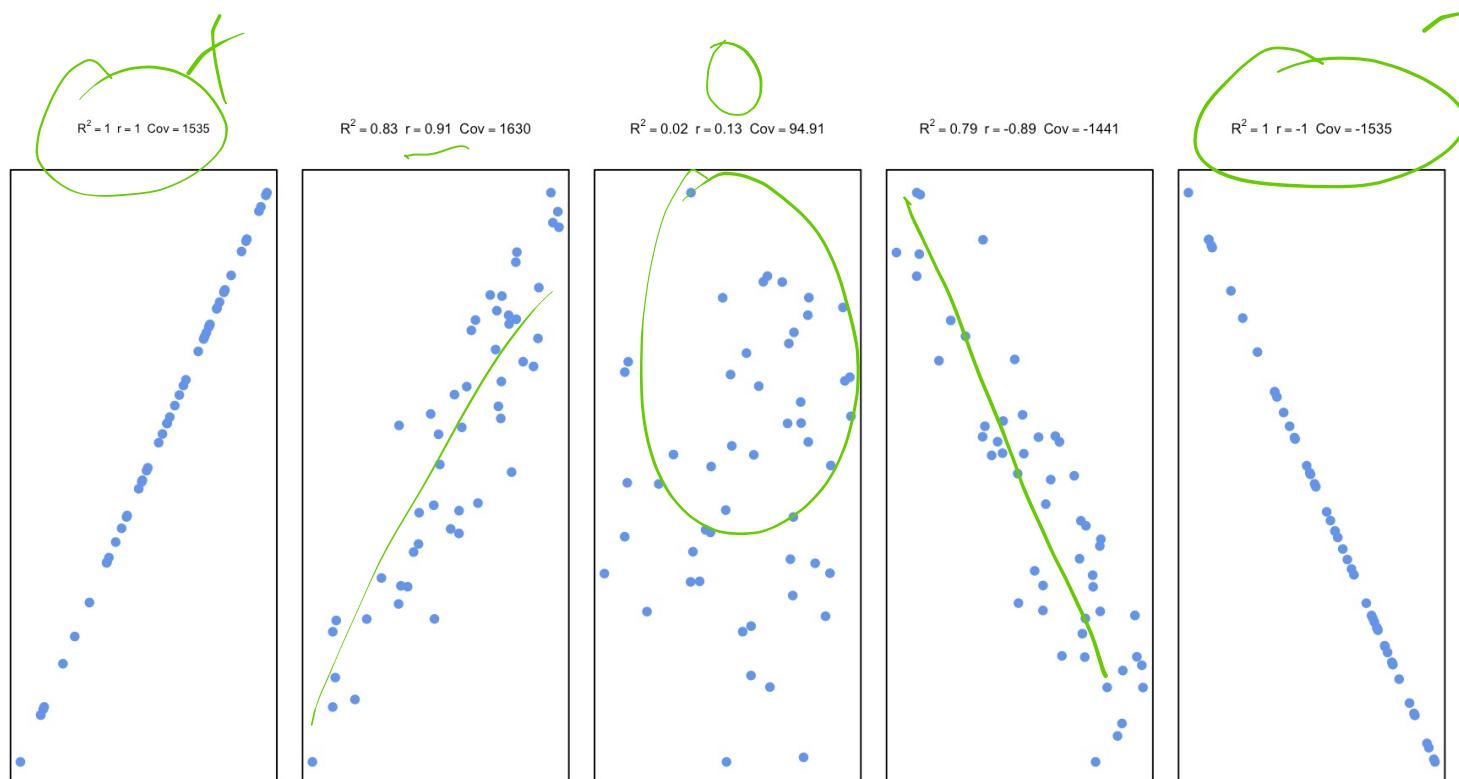
F-statistic  $\rightsquigarrow p\text{-value}$



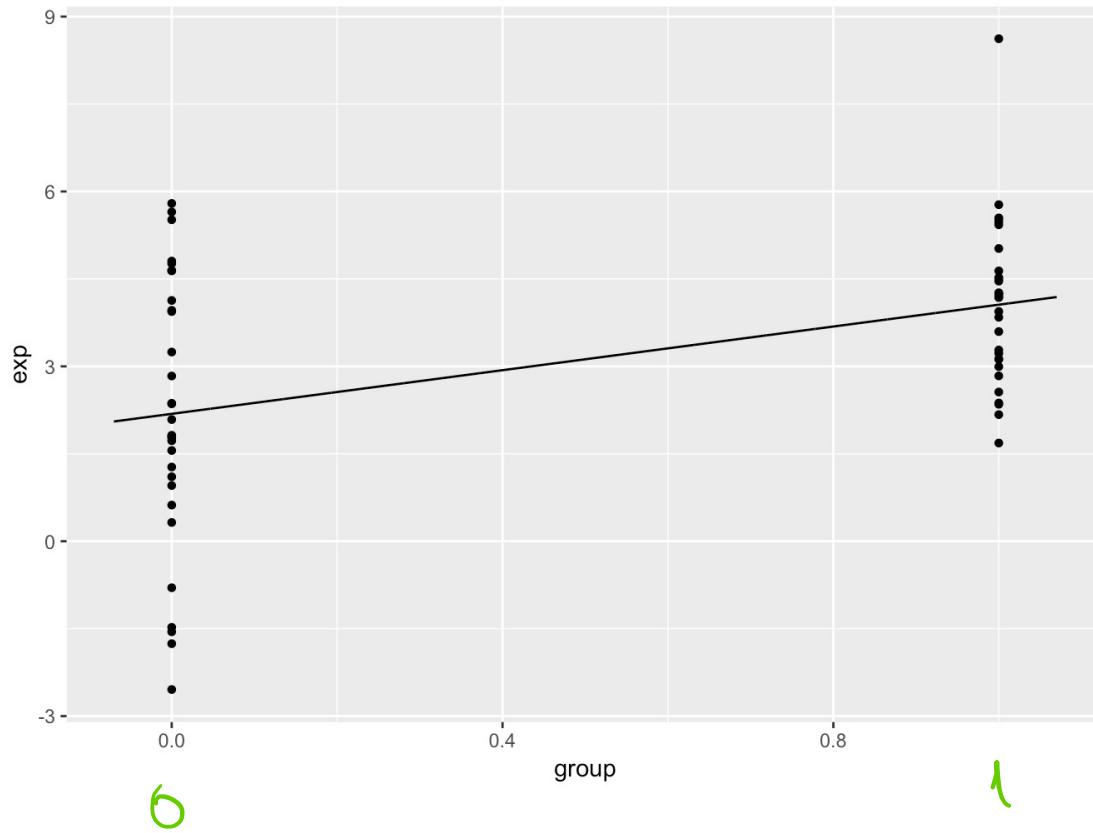
# Relationship between variables:

## Correlation

$$r_{xy} = \frac{cov(X, Y)}{\sigma_x \sigma_y} = \frac{\sum_{i=1}^n (\tilde{x}_i - \bar{x})(\tilde{y}_i - \bar{y})}{\sqrt{\sum_{i=1}^n (\tilde{x}_i - \bar{x})^2 \sum_{i=1}^n (\tilde{y}_i - \bar{y})^2}}$$



# Relationship between variables: Regression with categorical variables



[0,1] ~ Gene exp

Why regression framework for GE?

- Control covariates → Age, Sex
- Have more groups ↗ mech.

# Statistics for Genomics:

## Recap

