

Experimental Techniques

Last time:

> covariance

Today:

> Ch.9

> correlation and covariance

> Review

Main Results on Covariance

$$q_i = q(x_i, y_i) \leftarrow$$

$$\sigma_{xy} = \frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})(y_i - \bar{y}).$$

covariance — non zeros means they are not independent!

$$\sigma_q^2 = \left(\frac{\partial q}{\partial x}\right)^2 \sigma_x^2 + \left(\frac{\partial q}{\partial y}\right)^2 \sigma_y^2 + 2 \frac{\partial q}{\partial x} \frac{\partial q}{\partial y} \sigma_{xy}$$

std of q given non independent samples!

$$\sigma_q \leq \left| \frac{\partial q}{\partial x} \right| \sigma_x + \left| \frac{\partial q}{\partial y} \right| \sigma_y$$

naive estimate is still upper bound!
Schwarz inequality

never greater even when x_i, y_i are not independent.

Coefficient of Linear Correlation

Q: Given a set of measurements $(x_1, y_1) \dots (x_N, y_N)$, how well do they support the hypothesis that x and y are linearly related?

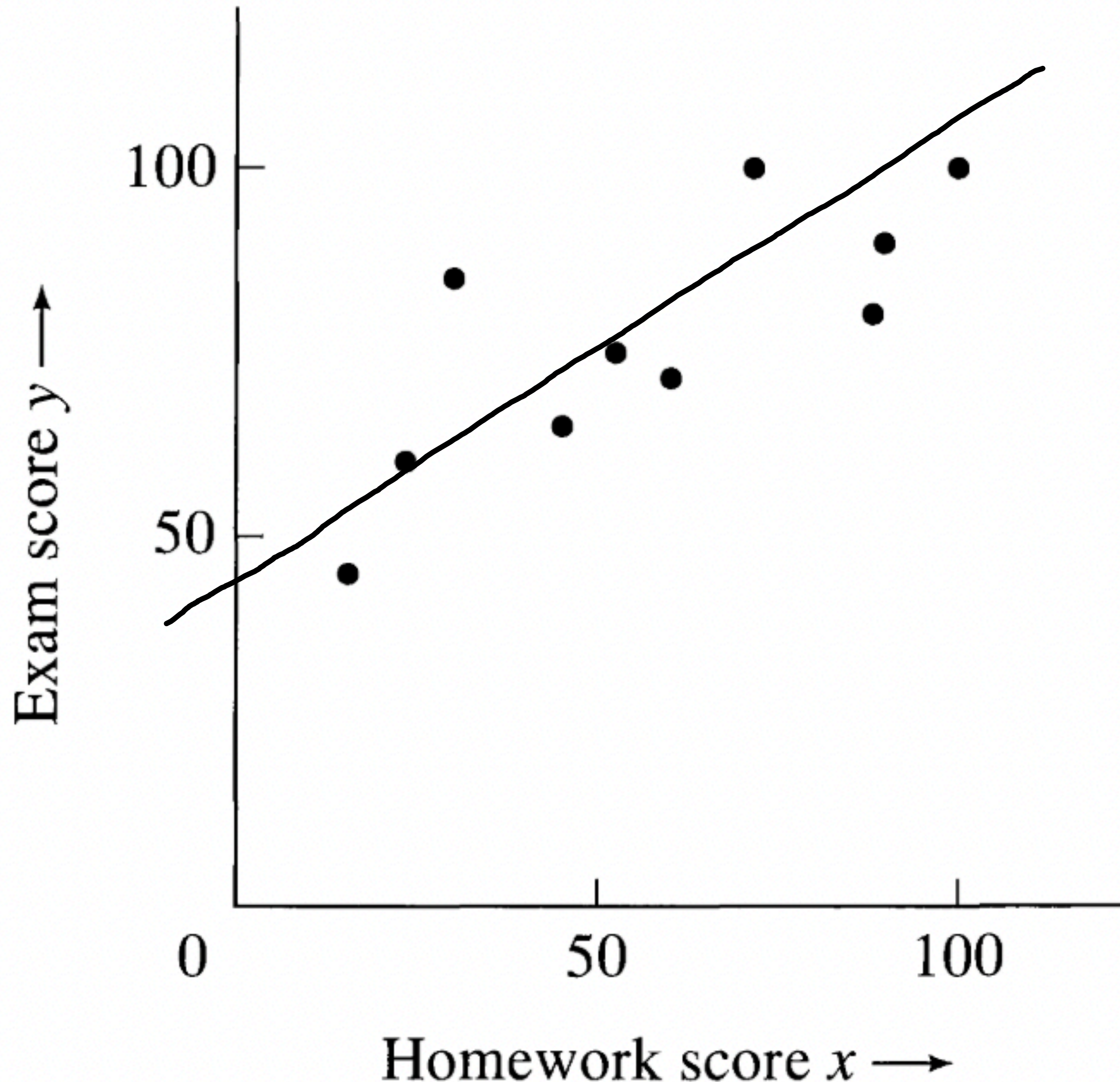
$$y = \underline{A} + \underline{B}x.$$

Using the method of least squares, we can find the values of A and B for the line that best fits the points

If we already have a reliable estimate of the uncertainties in the measurements, we can see whether the measured points do lie reasonably close to the line

We don't have uncertainty measures (each measurement is different range), how can we determine how well our data fits

Example



Professor plots HW vs Exam score
Hypothesis: better HW \rightarrow better score

The professor hopes to show that high exam scores tend to be correlated with high homework scores, and vice versa

This kind of experiment has no uncertainties in the points; each student's two scores are known exactly.

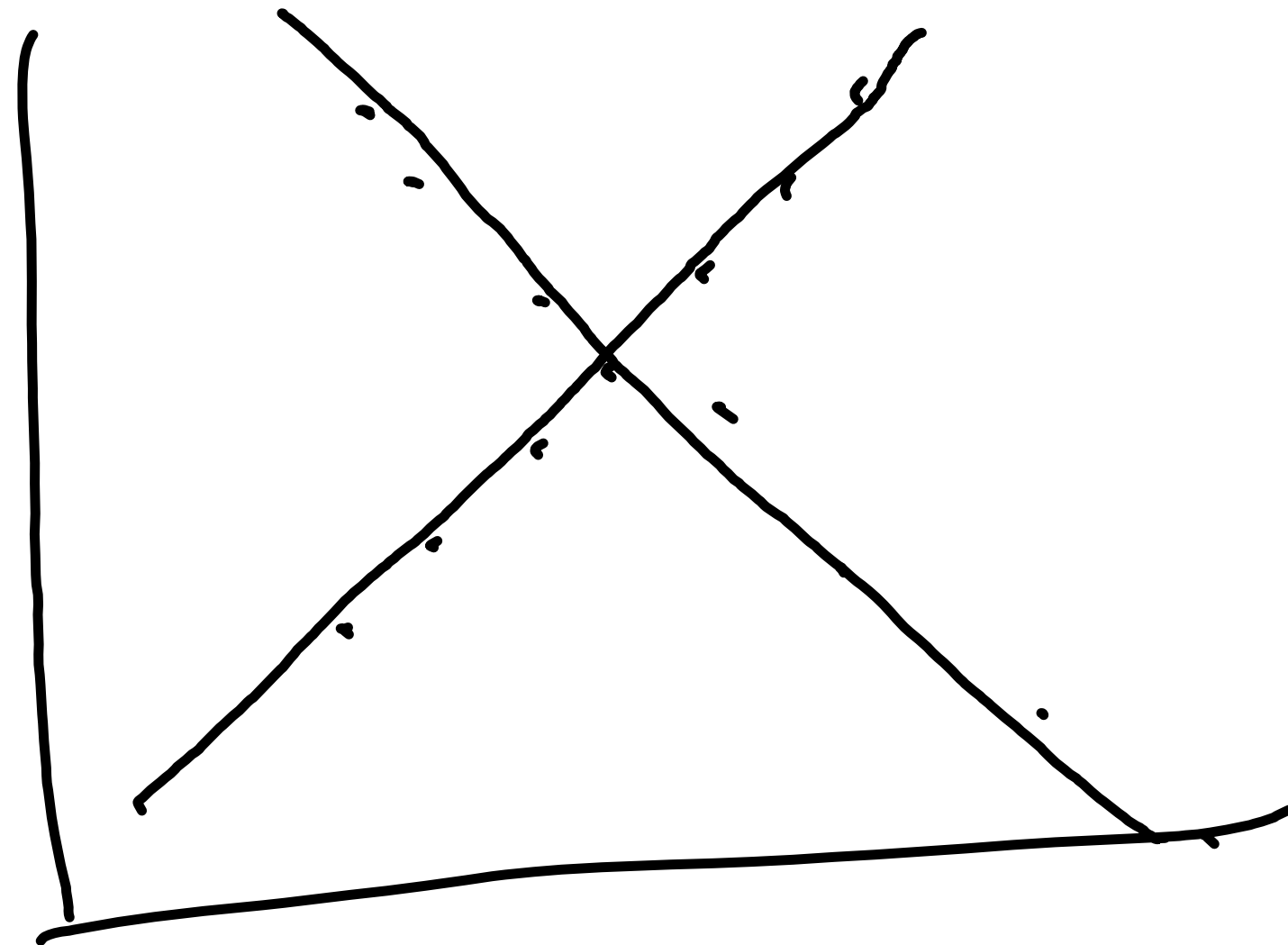
Really, what we want to understand the correlation

Correlation Coefficient

$$r = \frac{\sigma_{xy}}{\sigma_x \sigma_y}$$

$$-1 \leq r \leq 1$$

$$r = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum (x_i - \bar{x})^2 \sum (y_i - \bar{y})^2}}$$



Correlation Coefficient - perfect correlation

$$y_i = A + Bx_i$$

$$\bar{y} = A + B\bar{x}$$

$$y_i - \bar{y} = B(x_i - \bar{x})$$

$$r = \frac{B \sum (x_i - \bar{x})^2}{\sqrt{\sum (x_i - \bar{x})^2 B^2 \sum (x_i - \bar{x})^2}} = \frac{B}{|B|} = \pm 1$$

Correlation Coefficient revisit data

$$r = 0.8$$

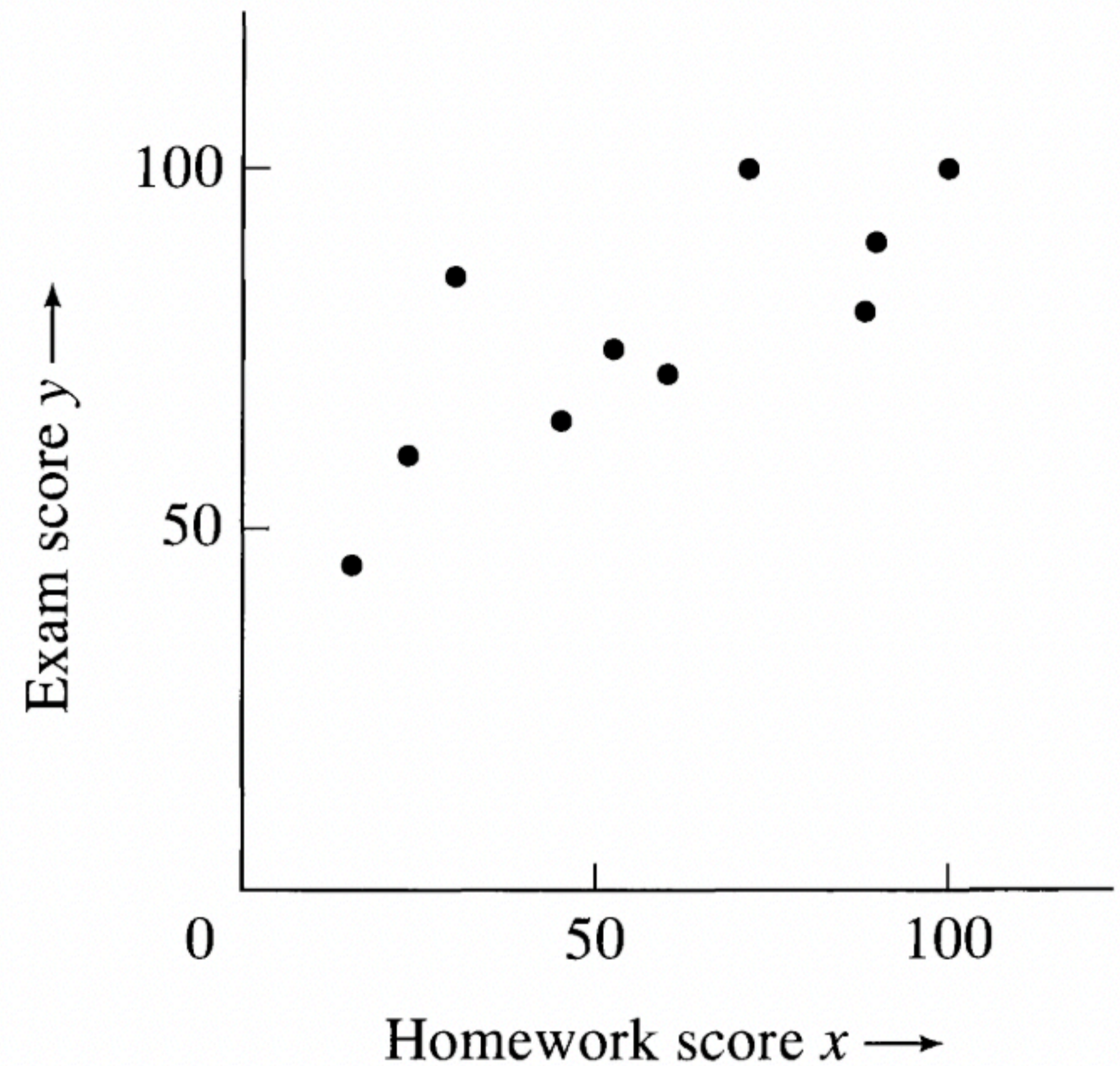


Table 9.3. Students' scores.

Student i	1	2	3	4	5	6	7	8	9	10
Homework x_i	90	60	45	100	15	23	52	30	71	88
Exam y_i	90	71	65	100	45	60	75	85	100	80

Find $r \rightarrow$ determine if significant at 95%.

Quantitative Significance of r

But how can we decide objectively what is good 'r'?

Suppose the two variables x and y are in reality uncorrelated; that is, in the limit of infinitely many measurements, the correlation coefficient r would be zero.

We can calculate the probability that r will exceed any specific value:

$$Prob_N(|r| \geq r_0)$$

$$Prob_N(|r| \geq 0.8)$$

Not straight forward calculation

Prob N measurements of two uncorrelated variables x and y would produce a correlation coefficient with

	r_0											
N	0	0.1	0.2	0.3	0.4	0.5	0.6	0.7	0.8	0.9	1	
3	100	94	87	81	74	67	59	51	41	29	0	
→ 6	100	85	70	56	43	31	21	12	6	1	0	
→ 10	100	78	58	40	25	14	7	2	0.5		0	
20	100	67	40	20	8	2	0.5	0.1			0	
50	100	49	16	3	0.4						0	

Given 10 random samples 0.5 %

→ $r = 0.80$

Back to our original question

Q: Given a set of measurements $(x_1, y_1) \dots (x_N, y_N)$, how well do they support the hypothesis that x and y are linearly related?

1. we can calculate correlation coefficient
2. we can find the probability of observing r with uncorrelated
3. if the value is sufficiently small, we support our hypothesis!

“significant” if the probability of obtaining a coefficient r with $|r| = |r|$ from uncorrelated variables is less than 5%. A correlation is sometimes called “highly significant” if the corresponding probability is less than 1%.

Main result: we have a quantitative measure of how improbable it is that they are uncorrelated.

Correlation Coefficient revisit data

$$r = 0.8$$

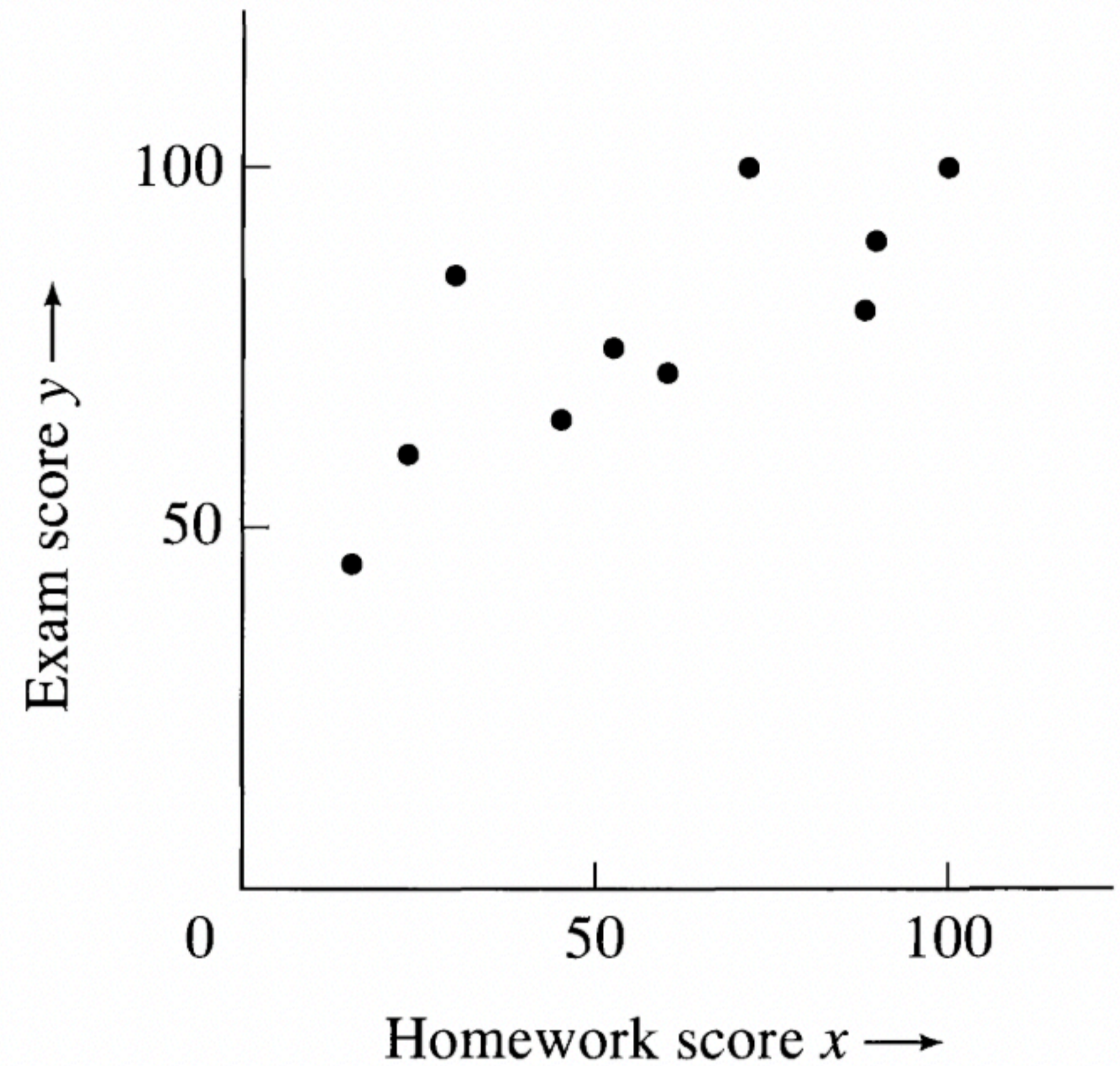


Table 9.3. Students' scores.

Student i	1	2	3	4	5	6	7	8	9	10
Homework x_i	90	60	45	100	15	23	52	30	71	88
Exam y_i	90	71	65	100	45	60	75	85	100	80

Prob₁₀ ($|r| > 0.8$)

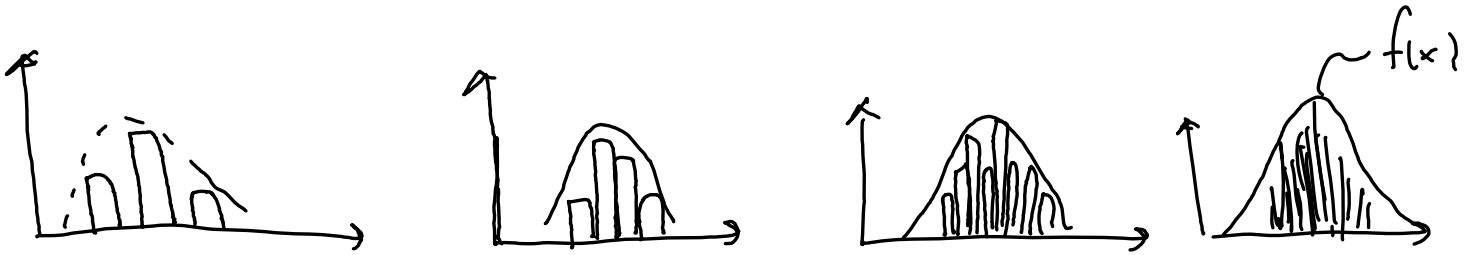
$n = 10$

Exam - 2 (Final) 3/19 @ 3-6pm

Exams will open @ 2:55pm
close @ 6:05pm

- include any external code/excel files.
- Show all your work
- organize the problems in order.
- box your answer !!

① Limiting Distributions



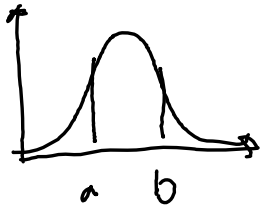
$N \rightarrow \infty$
the resulting distribution is called
"limiting distribution"

$f(x)$: Probability Density Function (PDF)

$$\int_{-\infty}^{\infty} f(x) dx = 1$$

$$\bar{x} = \int_{-\infty}^{\infty} x f(x) dx$$

$$\sigma_x^2 = \int_{-\infty}^{\infty} (x - \bar{x})^2 f(x) dx$$



$$\int_a^b f(x) dx$$

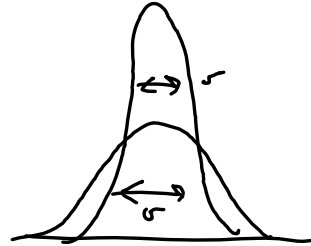
② Normal Distribution

→ most important limiting Distribution

$$G_{\mu, \sigma}(x) = \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{(x - \mu)^2}{2\sigma^2}}$$

μ : center parameter

σ : width parameter



$$\text{Prob}(\text{with } t\sigma) = \frac{1}{\sqrt{2\pi}} \int_{-t}^t e^{-z^2/2} dz$$

error function

we use a table for it.

③ Principle of Maximum Likelihood

$$\{x_1, \dots, x_N\}$$

We want to fit a normal distribution
maximizes:

$$\text{Prob}_{\mu, \sigma} (x_1, \dots, x_N) \propto \frac{1}{\sigma^N} e^{-\sum (x_i - \bar{X})^2 / 2\sigma^2}$$

parameters

minimizing $(x_i - \bar{X})^2 / 2\sigma^2$

(1.) Rejection of Data

$$t_{sus} = \frac{|X_{sus} - \bar{x}|}{\sigma}$$

Prob (outside $t_{sus} \cdot \sigma$) \rightarrow lookup on table.

$n = N \cdot \text{Prob}(\cdot) < 1/2 \Rightarrow$ outlier.

Chauvenet's criterion.

(5.) Weighted Average

Q: How to combine multiple measurements w/
different uncertainties?

$$X_{\text{wav}} = \frac{w_A X_A + w_B X_B}{w_A + w_B}$$

$$w_A = \frac{1}{\sigma_A^2}$$

$$w_B = \frac{1}{\sigma_B^2}$$

$$\Rightarrow X_{\text{wav}} = \frac{\sum w_i X_i}{\sum w_i}$$

$$w_i = \frac{1}{\sigma_i^2}$$

(6.) Least Squares

$$y = A + Bx$$

$$\chi^2 = \sum \frac{(y_i - A - Bx_i)^2}{\sigma_y^2}$$

$$\left\{ \begin{array}{l} \frac{\partial \chi^2}{\partial A} = -\frac{2}{\sigma_y^2} \sum (y_i - A - Bx_i) = 0 \\ \frac{\partial \chi^2}{\partial B} = -\frac{2}{\sigma_y^2} \sum x_i (y_i - A - Bx_i) = 0 \end{array} \right.$$

$$\left\{ \begin{array}{l} A \cdot N + B \sum x_i = \sum y_i \\ A \cdot \sum x_i + B \sum x_i^2 = \sum x_i y_i \end{array} \right\} \text{Normal Eq.}$$

$$A = \frac{\sum x^2 \sum y - \sum x \sum xy}{\Delta}$$

$$B = \frac{N \sum xy - \sum x \sum y}{\Delta}$$

$$\Delta = N \sum x^2 - (\sum x)^2$$

$$\sigma_y = \sqrt{\frac{1}{N-2} \sum (y_i - A - Bx_i)^2}$$

$$\sigma_A = \sigma_y \sqrt{\frac{\sum x^2}{\Delta}}$$

$$\sigma_B = \sigma_y \sqrt{\frac{N}{\Delta}}$$

Matrix Version :

$$y = A + Bx$$

$$Y = X \beta$$

β : parameter vector

X : data matrix

Y : true vector

$$\begin{Bmatrix} y_1 \\ \vdots \\ y_n \end{Bmatrix} = \begin{bmatrix} x_1 & 1 \\ \vdots & \vdots \\ x_n & 1 \end{bmatrix} \begin{bmatrix} B \\ A \end{bmatrix}$$

$n \times 1$

$n \times 2$

2×1

pseudo-inverse (pinv).

$$\beta^* = \left[\overbrace{(X^T X)^{-1} X^T} \right] Y$$

$$y = A \sin(x) + B$$

$$\begin{Bmatrix} y_1 \\ \vdots \\ y_N \end{Bmatrix} = \begin{bmatrix} \sin(x_1) & 1 \\ \vdots & \vdots \\ \sin(x_N) & \vdots \end{bmatrix} \begin{Bmatrix} A \\ B \end{Bmatrix}$$

$$y = A_0 + A_1 x + A_2 x^2 + A_3 x^3$$

$$\begin{Bmatrix} y_1 \\ \vdots \\ y_N \end{Bmatrix} = \begin{bmatrix} 1 & x_1 & x_1^2 & x_1^3 \\ \vdots & \vdots & \vdots & \vdots \\ 1 & x_N & x_N^2 & x_N^3 \end{bmatrix} \begin{Bmatrix} A_0 \\ A_1 \\ A_2 \\ A_3 \end{Bmatrix}$$

$$\begin{bmatrix} A_0^* \\ A_1^* \\ A_2^* \\ A_3^* \end{bmatrix} = \left[(X^T X)^{-1} X^T \right] Y$$