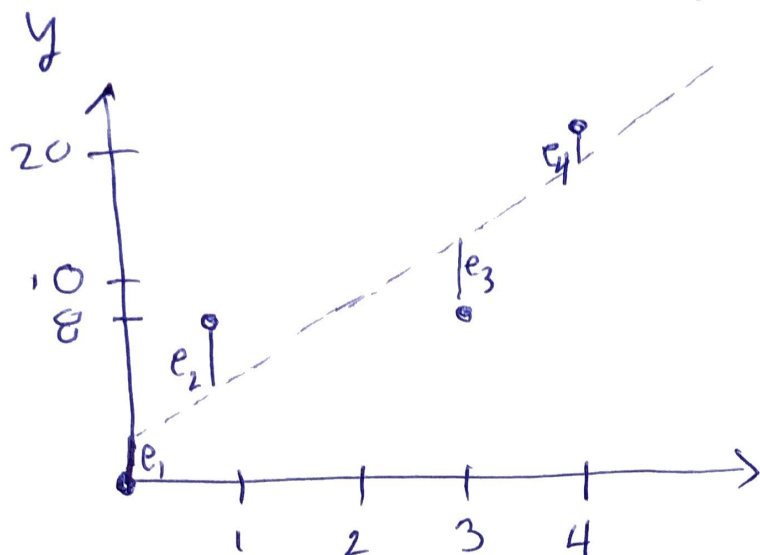


Data :  $(0,0), (1,8), (3,8), (4,20)$

Section 12.1  
In textbook



Line of best fit

$$y = C + Dt$$

is the one ~~in~~ in which  $C$  and  $D$  are chosen such that

$$e_1^2 + e_2^2 + e_3^2 + e_4^2$$

is minimized.

We had

$$\begin{bmatrix} e_1 \\ e_2 \\ e_3 \\ e_4 \end{bmatrix} = \begin{bmatrix} 0 - (C + D \cdot 0) \\ 8 - (C + D \cdot 1) \\ 8 - (C + 3D) \\ 20 - (C + 4D) \end{bmatrix} = \begin{bmatrix} -C \\ 8 - C - D \\ 8 - C - 3D \\ 20 - C - 4D \end{bmatrix}$$

So we can solve the problem by ~~correcting~~ <sup>regarding</sup> it as a calculus problem:

$$f(C, D) = C^2 + (8 - C - D)^2 + (8 - C - 3D)^2 + (20 - C - 4D)^2$$

We can minimize  $f$  by finding critical points, eg when

$$\nabla f = \begin{bmatrix} \frac{\partial f}{\partial C} \\ \frac{\partial f}{\partial D} \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

Because the equation is a sum of squares, it's concave up and so any critical point is a minimum.

## A change of perspective

Suppose now that we assume that there is a 'ground truth' which is that there is a linear relationship between independent and dependent variables. This is the simplest mathematical relationship

$$Y = C + Dt + \text{Error}$$

For example, assume that

+ 10 ppm CO<sub>2</sub>  
increases average global temp. by 0.1°C.

Any deviations  $\vec{e} = Y - (C + Dt)$  are due to random measurement noise.

Usually it is assumed that  $e_1, \dots, e_n$  are independent ~~normal~~ random variables, representing noise.

Then the observations are  $Y_i = C + Dt_i + \epsilon_i$   
or

$$\begin{bmatrix} Y_1 \\ \vdots \\ Y_n \end{bmatrix} = \begin{bmatrix} 1 & t_1 \\ \vdots & \vdots \\ 1 & t_n \end{bmatrix} \begin{bmatrix} C \\ D \end{bmatrix} + \begin{bmatrix} \epsilon_1 \\ \vdots \\ \epsilon_n \end{bmatrix}$$

$$\vec{Y} = A\vec{x} + \vec{\epsilon} \quad (*)$$

$$= \vec{C} + D\vec{t} + \vec{\epsilon}, \text{ where } \vec{C} = \begin{bmatrix} C \\ D \end{bmatrix}, \vec{t} = \begin{bmatrix} t_1 \\ \vdots \\ t_n \end{bmatrix}$$

Usually,  $\epsilon_1, \dots, \epsilon_n \sim N(0, \sigma^2)$ .

This eqn (\*) is called the simple linear regression model. The expected value of  $Y$  is a linear function of  $\bar{x}$ .

Eg if  $\sigma^2 = 1$ , what if we sampled a point at  $t = 2$ ? What is ~~prob~~ prob that the measurement is ~~2 away from the~~ between 7 and 11? \$

$$P[7 \leq Y_2 \leq 11] = P[7 \leq 1 + 4 \cdot 2 + \epsilon_2 \leq 11] = P[\epsilon_2 \leq 2] \approx .98$$

Example Suppose we have a relationship between

$x$  = applied stress

(explaining or  
Predictor  
variable)

$y$  = time-to-failure

(response  
variable)

Suppose that the observed relationship is

$$y = 65 - 1.2x$$

with ~~errors~~ deviations ~~having~~  ~~$N(0, 8)$~~  ~~distribution~~  
from this being random with Gaussien with  
 $\sigma = 8$ .

Then for any <sup>fixed</sup> value  $x^*$  of stress, time  
to failure has normal distribution with  
mean  $\mu = 65 - 1.2x^*$  and  $\sigma = 8$

What is the Mean time to failure when  
applied stress is  $x = 20$ .

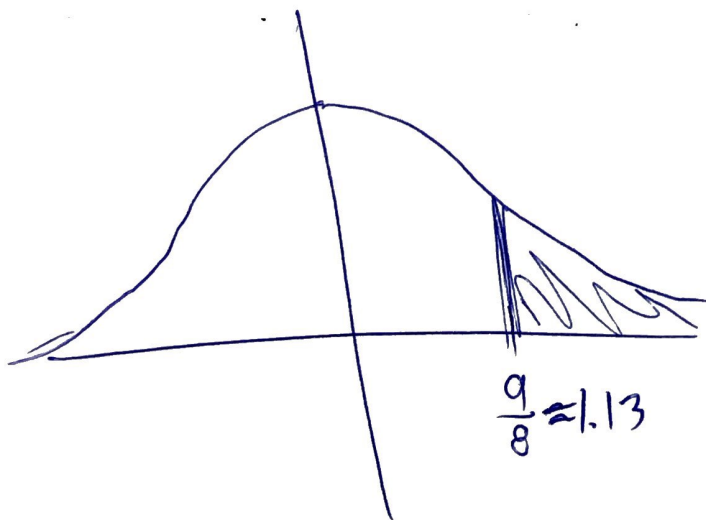


What is the probability that the time to failure is  $\geq 50$  when the applied stress is  $x=20$ ?  
When  $x=2$   $Y \sim N(41, 8)$  s.

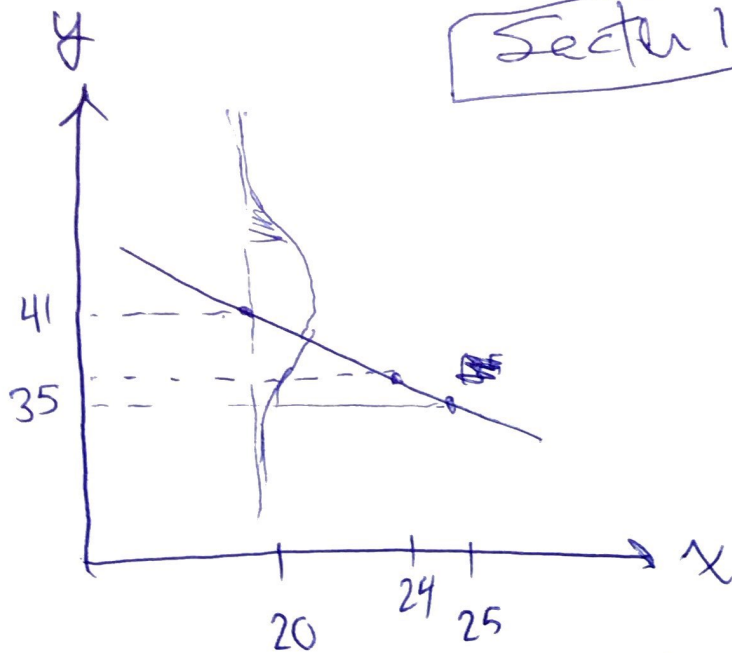
$$\cancel{P[Y \geq 50 | x=20]}$$

$$\cancel{P} \quad P[Y \geq 50] \stackrel{\text{standardize}}{=} P\left[\frac{Y-41}{8} \geq \frac{9}{8}\right]$$

$$\stackrel{=}{=} P\left[Z \geq \frac{9}{8}\right] = .13$$



# Section 12.1



~~What is~~ Let  ~~$U$~~  and  ~~$V$~~  be the  
time until failures when  ~~$x=25$~~  and  ~~$x=24$~~   
~~and  $x=25$~~

~~What is  $P[U \geq V]$ ?~~  $V \sim N(36.2, 8)$   
 $U \sim N(35, 8)$

What is  $P[U \geq \text{ ~~$V$~~ }]$ ?

$$P[U - \text{ ~~$V$~~ } \geq 0] = P\left[\frac{(U - V) - (-1.2)}{\sqrt{128}} \geq \frac{-1.2}{\sqrt{128}}\right]$$

Normal  
with mean  $36.2 - 35 = 1.2$

$Z$

$$\begin{aligned} \text{Var}(U - V) &= \text{Var}(U) + \text{Var}(V) \\ &= 8^2 + 8^2 = 128 \end{aligned}$$

$$\text{So } \sigma = \sqrt{128} \approx 11.3$$

$$\begin{aligned} &\approx P[Z \geq .11] \\ &= .45 \end{aligned}$$