

# Week 3

## Data rejection

- Is it okay to throw out data?
- Consider a data set

3.8, 3.5, 3.9, 3.4, 1.8

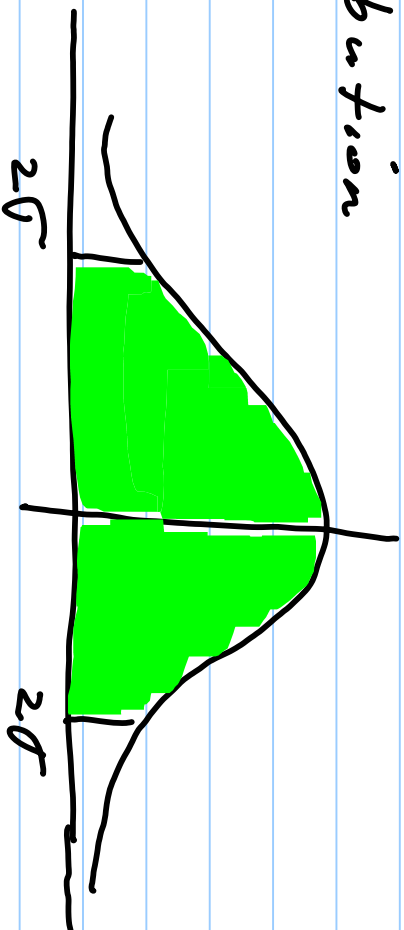
mean

$$\bar{X} = 3.4$$

Standard deviation

$$S_x = .8$$

If we assume a Gaussian Distribution



$$\begin{aligned}\text{Prob (outside } 2\sigma) &= 1 - \text{Prob (within } 2\sigma) \\ &= 1 - 0.95 \\ &= 0.05\end{aligned}$$

## Chauvenets Criterion

$N$  measurements

$$X_1, \dots, X_N$$

Consider one measurement  $X_{sus}$

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

Define  $N_3$  (expected number as deviant as  $X_{sus}$ )

So

4

$$n = N \times \text{Prob}(\text{outside } t_{\text{sus}}(\bar{V}))$$

If  $n < .5$  reject

From our data set

$$t_3 = \frac{1.8 - 3.4}{.8} = 2$$

Since

$$\text{Prob}(\text{outside } 2\sigma) = .05$$

$$n = 6 \times .05 = .3 \text{ reject!}$$

New

$$\bar{X} = \frac{3.8 + 3.5 + 3.9 + 3.9 + 3.4}{5} = 3.64 \text{ compares with } 3.4$$

Is this process a good idea?

- Can you identify a reason
- reject data at your peril
- Never a second time

## Problem of Combining Measurements <sup>6</sup>

Suppose you have two measurements of a lifetime

$$\text{Meas. A} = 10 \pm .5 \text{ sec.}$$

$$\text{Meas B} = 12 \pm 1 \text{ sec.}$$

Is it meaningful to combine these results?

?

Probability of value  $X_A$

$$\text{Prob}_x(x_A) \propto \frac{1}{\sigma_A} e^{-\frac{(x_A - \mu)^2}{2\sigma_A^2}}$$

and for value  $X_B$

$$\text{Prob}(x_B) = \frac{1}{\sigma_B} e^{-\frac{(x_B - \mu)^2}{2\sigma_B^2}}$$

Probability of  $X_A$  and  $X_B$

$$\begin{aligned} \text{Prob}_x(x_A, x_B) &= \text{Prob}_x(x_A) \text{Prob}(x_B) \\ &\propto \frac{1}{\sigma_A \sigma_B} e^{-\frac{x^2}{2}} \end{aligned}$$

$$X^2 = \left( \frac{X_A - X}{\sigma_A} \right)^2 + \left( \frac{X_B - X}{\sigma_B} \right)^2 \quad \delta$$

Find minimum

$$\frac{dX^2}{dX} = 0 = 2 \frac{X_A - X}{\sigma_A^2} + 2 \frac{X_B - X}{\sigma_B^2}$$

Solve for X

$$X = \left( \frac{X_A}{\sigma_A^2} + \frac{X_B}{\sigma_B^2} \right) / \left( \frac{1}{\sigma_A^2} + \frac{1}{\sigma_B^2} \right)$$



Define

$$w_A = \frac{1}{\sqrt{A^2}} \quad , \quad w_B = \frac{1}{\sqrt{B^2}}$$

Weighted average becomes

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

For our example  $X_A = 10$ ,  $X_B = 12$

Simple average = 11

Weighted average = 10.7

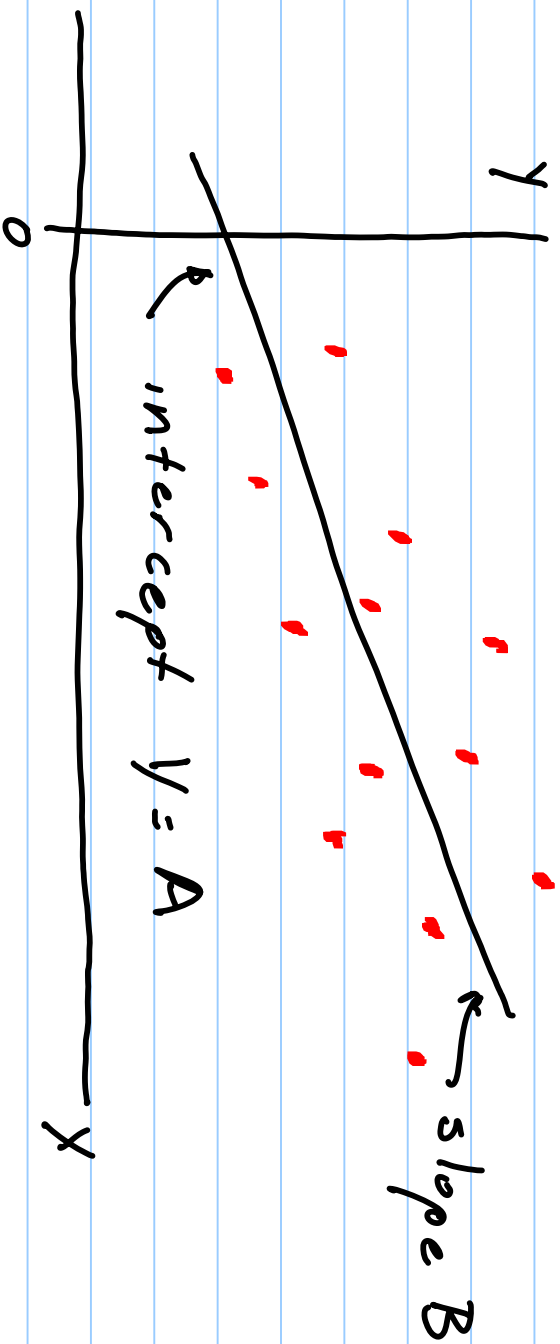
Karl Friedrich Gauss



Least Squares 1795

# Least Squares Fitting (Regression)

11



$$Y = A + BX$$

Characterise line by A & B

If we knew the constants  $A$  +  $B$

$$(\text{true value for } y_i) = A + Bx_i$$

$$P_{\text{Prob}_{A,B}}(y_i) \propto \frac{1}{\sigma_y} e^{-\frac{(y_i - A - Bx_i)^2}{2\sigma_y^2}}$$

For complete set of measurements

$$P_{\text{Prob}_{A,B}}(y_1, \dots, y_N) = P_{\text{Prob}_{A,B}}(y_1) \dots P_{\text{Prob}_{A,B}}(y_N) \\ \propto \frac{1}{\sigma_y^N} e^{-\chi^2/2}$$

where exponent

$$\chi^2 = \sum_{i=1}^N \frac{(y_i - A - Bx_i)^2}{\sigma_y^2}$$

Minimum value  $\chi^2$

$$\frac{\partial \chi^2}{\partial A} = \frac{-2}{\sigma_y^2} \sum_{i=1}^N (y_i - A - Bx_i) = 0$$

and

$$\frac{\partial \chi^2}{\partial B} = \frac{-2}{\sigma_y^2} \sum_{i=1}^N (y_i - A - Bx_i)x_i = 0$$

Two Equations

$$AN + B \sum x_i = \sum y_i$$

$$A \sum x_i^2 + B \sum x_i = \sum x_i y_i$$

Solving

$$A = \frac{\sum x_i^2 \sum y_i - \sum x_i \sum y_i}{\Delta}$$

$$B = \frac{N \sum x_i y_i - \sum x_i \sum y_i}{\Delta}$$

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

Uncertainty in Measurements of  $y$

- good guess

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

- better answer

$$V_y = \sqrt{\frac{1}{N-2} \sum_{i=1}^N (y_i - A - Bx_i)^2}$$

## Covariance

$$q_i = q(x_i, y_i) \quad i = 1, \dots, N$$

$N$  data pairs  $(x_1, y_1) \dots (x_N, y_N)$

$$\bar{q} \quad \text{mean}$$

$\sigma_q$  standard deviation

Make approximation

$$\begin{aligned} q_i &= q(x_i, y_i) \\ &\approx q(\bar{x}, \bar{y}) + \frac{\partial q}{\partial x}(x_i - \bar{x}) \\ &\quad + \frac{\partial q}{\partial y}(y_i - \bar{y}) \end{aligned}$$



Then

$$\bar{q} = \frac{1}{N} \sum_{i=1}^N q_i$$

$$= \frac{1}{N} \sum_{i=1}^N \left[ q(\bar{x}, \bar{y}) + \frac{\partial q}{\partial x} (x_i - \bar{x}) + \frac{\partial q}{\partial y} (y_i - \bar{y}) \right]$$

$$\text{Since } \sum_{i=1}^N (x_i - \bar{x}) = 0 = \sum_{i=1}^N (y_i - \bar{y}) = 0$$

$$\Rightarrow \bar{q} = q(\bar{x}, \bar{y})$$

$$V_{q_2}^2 = \frac{1}{N} \sum (q_i - \bar{q})^2$$

Since  $q_i \approx \bar{q} + \frac{\partial q}{\partial x} (x_i - \bar{x}) + \frac{\partial q}{\partial y} (y_i - \bar{y})$

Then

$$V_{q_2}^2 = \frac{1}{N} \sum \left[ \frac{\partial q}{\partial x} (x_i - \bar{x}) + \frac{\partial q}{\partial y} (y_i - \bar{y}) \right]^2$$

$$\begin{aligned}
 &= \left( \frac{\partial q}{\partial x} \right)^2 \frac{1}{N} \sum (x_i - \bar{x})^2 + \left( \frac{\partial q}{\partial y} \right)^2 \frac{1}{N} \sum (y_i - \bar{y})^2 \\
 &\quad + 2 \frac{\partial q}{\partial x} \frac{\partial q}{\partial y} \frac{1}{N} \sum (x_i - \bar{x})(y_i - \bar{y})
 \end{aligned}$$

Define Covariance

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

Then

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

## Correlation Coefficient

$$r = \frac{\overline{xy}}{\overline{x}\overline{y}}$$

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

$|r| \sim 1$  strong correlation

$|r| \sim 0$  weak correlation

— Summation —

# Poisson Distribution

- Event distribution

$Y$  Number of counts in given interval

$M = \bar{Y}$  mean number of counts in that interval

Prob ( $Y$  counts in a definite interval) =  $P_n(Y)$

$$= e^{-M} \frac{M^Y}{Y!}$$

$\bar{V} = M =$  average number  
of counts if we  
repeat experiment  
many times

$M =$  rate  $\times$  time  $\times$  RT

— Poisson Simulator —

## Current Example from Physics

- Proton decay

Consider  $N$  protons

$$N = N_0 e^{-\lambda t}$$

where  $\lambda = \frac{1}{\tau} = 10^{-33}/\text{year}$

$$e^{-\lambda t} \approx 1 - \lambda t$$

So

$$N \approx N_0 (1 - \lambda t)$$

21

Super Kamokanda Neutrino Detector

Volume  $\approx 7.5 \times 10^{33}$  protons

Then

$$N_0 - N = N_0 \lambda t = (7.5 \times 10^{33} \times 10^{-35}) \times 7.8$$

40% of tank covered by detectors so we expect to see  $\approx 3$  decays



Thus far no decays seen  
Probability of no observations

$$P(N) = e^{-N} \frac{N^N}{N!}$$

$$P(0) = e^{-3} \frac{3^0}{0!} \approx \underline{\underline{.05}}$$

Suggests that  $10^{33}$  years is too short!

## Properties of a Poisson Distribution

- Standard deviation

$$\sigma_Y^2 = \overline{(Y - \bar{Y})^2}$$

or

$$\sigma_Y^2 = \bar{Y} - (\bar{Y})^2$$

and

$$\bar{Y}^2 = N^2 + N$$

$$\sigma_Y = \sqrt{N}$$

---

---

Thus if we make a measurement of a number of events in time  $T$  the mean count is then

$$\bar{y} \pm \sqrt{\bar{y}}$$