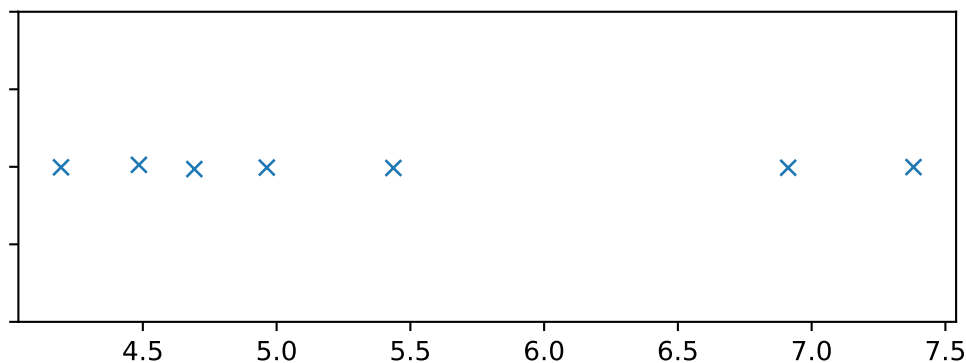


# 1 Consolidation assignment

## 1. Complete the survey about the project phase on StudOn.

2. “Seven scientists”. Assume that you are trying to infer the numerical value of some natural constant,  $\mu$ , which is unknown to you. You ask seven (or  $N$ ) scientists to use their measurement devices in order to give you an estimate of  $\mu$ . We assume that the  $i$ th scientist’s device returns the value  $y_i = \mu + \epsilon_i$ , where each  $\epsilon_i$  is a measurement error distributed according to a Gaussian random variable  $N(0, \sigma_i)$ .

Consider measurements  $y = [4.96334778, 5.43644897, 7.38067265, 4.69244271, 6.91227063, 4.48530689, 4.19409183]$ . (see figure below)



Write down explicitly the expression for the *likelihood*  $\mathbb{P}(y|\mu, \{\sigma_i\})$ , i.e. the probability density of a measurement  $y$  given a fixed set of parameters  $\mu$  and  $\{\sigma_i\}$ .

- First assume that all devices have the same measurement uncertainty  $\sigma_i = 1$ . What is the maximum likelihood estimator for  $\mu$  (i.e. the value  $\hat{\mu}_{ML} = \arg\max_{\mu} \mathbb{P}(y|\mu, \{\sigma_i\})$ )?
  - Now we assume that we don't know each scientist's measuring uncertainty  $\sigma_i$ . What is the joint maximum likelihood estimator for  $\mu$  and all  $\sigma_i$ , i.e. which values of  $\mu$  and  $\sigma_i, i = 1, \dots, N$  minimize the likelihood  $\mathbb{P}(y|\mu, \{\sigma_i\})$ ? What goes wrong?
  - Elaborate on a possible solution for the problem in b).
3. Let  $X \sim N(\mu, \sigma^2)$ ,  $\epsilon \sim N(0, \gamma^2)$  and  $Y = a \cdot X + \epsilon$ . In class we derived the following:

$$Y|X \sim N(a \cdot X, \gamma^2)$$

For simplicity we assume that  $\mu = 0$ . This means that

$$f_X(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \cdot \exp\left(-\frac{x^2}{2\sigma^2}\right)$$

$$f_{Y|X}(y|x) = \frac{1}{\sqrt{2\pi\gamma^2}} \cdot \exp\left(-\frac{(y - ax)^2}{2\gamma^2}\right).$$

Calculate  $f_{X,Y}(x, y) = f_{Y|X}(y|x) \cdot f_X(x)$  and thus derive the distribution of  $(X, Y)$ . Hint: The density for a 2-d multivariate Gaussian distribution  $(W, Z) \sim N(\vec{m}, C)$  is

$$f_{W,Z}(w, z) = \frac{1}{\sqrt{2\pi \det C}} \exp \left( -\frac{\left( \begin{pmatrix} w \\ z \end{pmatrix} - \vec{m} \right)^T \cdot C^{-1} \cdot \left( \begin{pmatrix} w \\ z \end{pmatrix} - \vec{m} \right)}{2} \right)$$

4. (optional) If you're interested in Simpson's paradox, you can read more about it on <https://pwacker.com/simpson.html>

## 2 Preparation assignment

We consider the setting of linear regression with a polynomial regression function: Given data  $z \in \mathbb{R}^{2 \times N}$  (i.e. each of the  $N$  data points  $z_i$  has two coordinates which we abbreviate by  $z_i = (x_i, y_i)$ ).

We assume that the data was generated by noisy evaluation of a polynomial with  $n$  unknown coefficients, i.e.

$$y_i = a_0 + a_1 \cdot x_i + a_2 \cdot x_i^2 + \dots + a_{n-1} \cdot x_i^{n-1} + \varepsilon_i$$

where  $\varepsilon_i \sim N(0, \sigma^2)$  are independent noise terms (with known variance  $\sigma^2$ ).

- Write Matlab/Octave code that can do the following: Given parameters  $\vec{a} = (a_0, \dots, a_{n-1})$ , measurements positions  $x_i$  and noise parameter  $\sigma$ , generate data  $y$  as above.
- Write the measurement process in vectorial form

$$\vec{y} = M \cdot \vec{a} + \vec{\varepsilon}.$$

## 3 Notes / Insights from class