



UNIVERSITETET  
I OSLO

DEPARTMENT OF PHYSICS

STK4051 - COMPUTATIONAL STATISTICS

---

## Mandatory Assignment Part 1/2

---

*Author:*  
Jonas Semprini Næss

Date: 1st May 2024

---

## Table of Contents

<b>List of Figures</b>	<b>ii</b>
<b>List of Tables</b>	<b>ii</b>
<b>Feedback</b>	<b>1</b>
<b>Problem 1 - (Lp-regularization)</b>	<b>1</b>
a.) . . . . .	1
b.) . . . . .	2
c.) . . . . .	3
d.) . . . . .	4
e.) . . . . .	4
f.) . . . . .	5
<b>Problem 2 - (EM-algorithm)</b>	<b>6</b>
a.) . . . . .	6
b.) . . . . .	7
c.) . . . . .	7
d.) . . . . .	9
e.) . . . . .	10
f.) . . . . .	11
g.) . . . . .	12
<b>Problem 3 - (Bayesian Approach)</b>	<b>13</b>
a.) . . . . .	13
b.) . . . . .	15
c.) . . . . .	16
<b>Problem 4 - (Combinatorial optimization)</b>	<b>16</b>
a.) . . . . .	17
b.) . . . . .	18
c.) . . . . .	19
d.) . . . . .	20
<b>Problem 5 - (Stochastic gradient decent, SGD)</b>	<b>21</b>
a.) . . . . .	21

---

b.) . . . . .	22
c.) . . . . .	22
<b>Appendix</b>	<b>24</b>
<b>A GitHub repository</b>	<b>24</b>
<b>B Source Code</b>	<b>24</b>

### List of Figures

1	Plots of $f_{p,\gamma}(\beta_i)$ where $f(\beta) = \beta$ is added for reference . . . . .	3
2	Plots of the inverse of $f_{p,\gamma}(\beta_i)$ where $f(\beta) = \beta$ is added for reference . . . . .	3
3	Finding the solution to $f_{p,\gamma}(\beta_i) = 0$ through the bisection method . . . . .	4
4	Scatter plot of bootstrap estimations for $\tau^2$ and $p$ . . . . .	11
5	Contour plot of $L(\theta y)$ in the grid $(p, \tau^2) \in [0.8, 1] \times [50, 130]$ . . . . .	13
6	Estimated $\beta$ for $y \in [-5, 5]$ . . . . .	16
7	Plots of the most optimal tour and cooling rate for the latest run. . . . .	17
8	Plots of the most optimal tour and evolution of distance with respect to iterations. . . . .	18
9	Various plots of accuracy measurements . . . . .	22

### List of Tables

1	Residuals for penalized regression . . . . .	5
2	Residuals for Alternative Estimation . . . . .	5
3	Iterative estimation of $p$ and $\tau^2$ where $i$ denotes the number of iterations . . . . .	10
4	$n!$ for $n$ up to 20 . . . . .	19
5	$\frac{(n-1)!}{2}$ for $n$ up to 20 . . . . .	19

---

## Feedback

The mandatory assignment has been engaging, exciting, challenging, and at times quite time-consuming. I have gained a deeper insight into various numerical methods that can be used in statistics, as well as a new perspective on other nuances of theory and methods that I have already touched upon. However, unfortunately, I did not have the opportunity to complete the entire second part of the assignment, since the assignment as a whole (both part 1 and 2) require its time and diligence complete to a satisfactory level. Nevertheless, I hope that what has been answered in these PDF's provide, and demonstrates the amount of work that has been put in.

### Problem 1 - (Lp-regularization)

We will in this exercise consider the problem where the number of parameters we will estimate are of the same size as the number of data. In this case it is common to use regularization to impose additional constraints on the model. In a simplified regression model, we have data on the form:

$$y_i = \beta_i + \varepsilon_i, \quad i = 1, 2, \dots, n$$

Where  $y_i$  is the data,  $\beta_i$  is the parameter and  $\varepsilon_i$  is an error term, we will assume that the error term has variations according to a normal distribution with mean zero and unit variance, i.e.  $\varepsilon_i \sim \mathcal{N}(0, 1^2)$

a.)

*Derive the maximum likelihood estimator for  $\beta_i$ ,  $i = 1, 2, \dots, n$*

**Solution:**

By definition we observe that when predicting or modelling  $y_i$  we are conditioned on the current estimate of  $\beta_i$ . Moreover, it is understood that the according probability of determining  $y_i$  is given by looking at the probability of  $\varepsilon_i$ , namely  $y_i - \beta_i$ . From this we recall that the probability density function of an arbitrary standard normal distributed variable  $X$  is given as

$$f(x) = \frac{1}{\sqrt{2\pi}} e^{-x^2/2}$$

where in our case we turn to  $\varepsilon_i$

$$f(\varepsilon_i) = \frac{1}{\sqrt{2\pi}} e^{-\varepsilon_i^2/2} \tag{1}$$

for  $\varepsilon_i \in (-\infty, \infty)$ . Hence, can we express the likelihood as

$$L(\beta|y) = \prod_{i=1}^n f(\varepsilon_i) = \prod_{i=1}^n f(y_i - \beta_i) \tag{2}$$

---

and consequently the log-likelihood

$$\begin{aligned}
\ell(\beta|y) &= \log(L(\beta|y)) \\
&= \log\left(\prod_{i=1}^n f(y_i - \beta_i)\right) \\
&= \sum_{i=1}^n \log(f(y_i - \beta_i)) \\
&= \sum_{i=1}^n \log\left(\frac{1}{\sqrt{2\pi}} e^{-(y_i - \beta_i)^2/2}\right) \\
&= \sum_{i=1}^n -\log(\sqrt{2\pi}) - \frac{1}{2}(y_i - \beta_i)^2
\end{aligned}$$

whereby we achieve the maximum likelihood estimator for  $\beta_i$  when differentiating  $\ell$  w.r.t  $\beta_i$  and solvig the resulting equation when equal to zero. That yields

$$\frac{\partial \ell(\beta|y)}{\partial \beta_i} = \frac{2(y_i - \beta_i)}{2} = y_i - \beta_i = 0 \quad (3)$$

meaning the maximum likelihood estimator for  $\beta_i$  is thus

$$\hat{\beta}_i = y_i, \quad i = 1, 2, \dots, n \quad (4)$$

■

This is reasonable considering we are working with a model where the parameter space is of equal size to the data samples ( $i$ ), and additionally does not demand any further transformations.

**b.)**

*Show that:*

$$f_{p,\gamma}(\beta_i) = \beta_i + \gamma \cdot \text{sign}(\beta_i) |\beta_i|^{p-1} \quad (5)$$

**Solution:**

When now turning to penalized least squares we recall from (a) that we already have an expression for the log-likelihood which by direct insertion into our model expression gives

$$\min_{\beta} \left\{ \log(\sqrt{2\pi}) + \frac{1}{2}(y_i - \beta_i)^2 + \frac{\gamma}{p} \|\beta\|_p^p \right\}. \quad (6)$$

Similarly as in (a.) we differentiate the expression w.r.t  $\beta_i$

$$\begin{aligned}
\frac{\partial}{\partial \beta_i} \left( \log(\sqrt{2\pi}) + \frac{1}{2}(y_i - \beta_i)^2 + \frac{\gamma}{p} \|\beta\|_p^p \right) &= \left( -(y_i - \beta_i) + \frac{\gamma}{p} p |\beta_i|^{p-1} \frac{\partial |\beta_i|}{\partial \beta_i} \right) \\
&= -(y_i - \beta_i) + \gamma |\beta_i|^{p-1} \frac{\partial |\beta_i|}{\partial \beta_i} \\
&= -(y_i - \beta_i) + \gamma |\beta_i|^{p-1} \text{sign}(\beta_i)
\end{aligned}$$

and set it equal to zero

$$-(y_i - \beta_i) + \gamma |\beta_i|^{p-1} \text{sign}(\beta_i) = 0 \quad (7)$$

$$\Updownarrow \quad (8)$$

$$y_i = \beta_i + \gamma \cdot \text{sign}(\beta_i) |\beta_i|^{p-1} \quad (9)$$

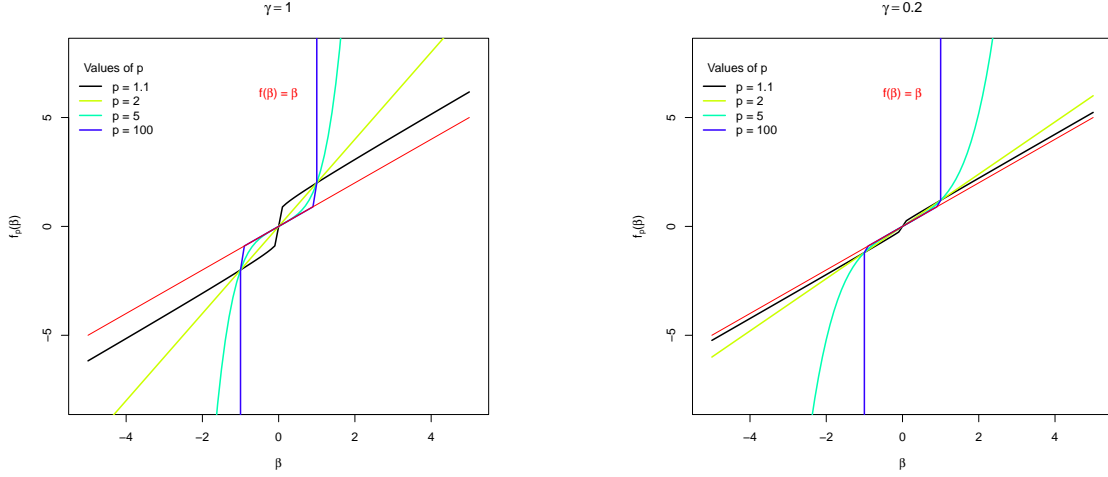
which is what we wanted to show. ■

c.)

For  $\gamma = 1$  and  $\gamma = 0.2$ , plot the function  $f_{p,\gamma}(\beta)$  for  $p = 1.1, 2, 5$  and  $100$  on the square  $[-5, 5] \times [-5, 5]$ , plot also the inverse function by flipping the order of the arguments in the plotting function. Give an interpretation of the results.

**Solution:**

The implementation can be found in Appendix B, from code listing (1). It produces the following plots.

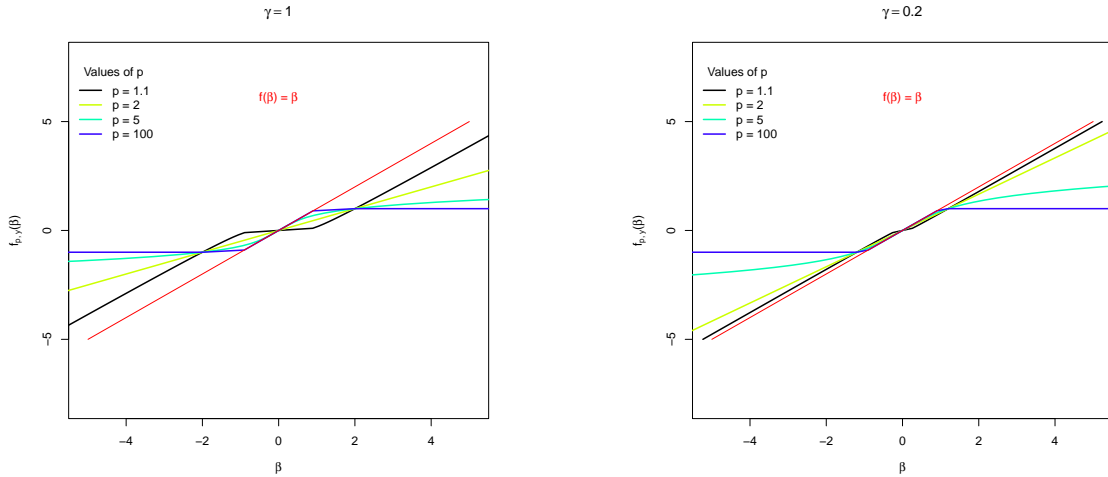


(a)  $f_{p,\gamma}(\beta_i)$  for  $\gamma = 1$

(b)  $f_{p,\gamma}(\beta_i)$  for  $\gamma = 0.2$

Figure 1: Plots of  $f_{p,\gamma}(\beta_i)$  where  $f(\beta) = \beta$  is added for reference

with the respective inverse functions



(a)  $f_{p,\gamma}^{-1}(\beta_i)$  for  $\gamma = 1$

(b)  $f_{p,\gamma}^{-1}(\beta_i)$  for  $\gamma = 0.2$

Figure 2: Plots of the inverse of  $f_{p,\gamma}(\beta_i)$  where  $f(\beta) = \beta$  is added for reference

The constant  $\gamma$  determines the impact of the penalty term on the function  $f$ . When  $\gamma$  equals zero, the penalty term is eliminated, resulting in a straight line. As  $\gamma$  increases, the regularization effect becomes more pronounced.

The  $p$ -value influences the type of  $L_p$ -norm used for regularization. Higher  $p$ -values lead to a greater penalty for  $\beta$  values larger than 1 in absolute value, while smaller  $p$ -values result in a lesser penalty, as raising  $\beta$  to the  $p$ th power tends toward zero for large  $p$ -values.

Overall, higher  $p$ -values result in a heavier penalty for large  $\beta$  values, and the  $\gamma$  parameter allows control over the overall magnitude of this penalty term.

d.)

Implement a function which finds the root of expression (3). Use the method of bisection (book page 23). The root of (3) will be the estimator  $\hat{\beta}_{\gamma,p}(y)$ . What are good starting values for upper and lower bounds (in general)? Test the algorithm for  $\gamma = 1$ , and  $p = 1.1$ ,  $p = 2$ , and  $p = 100$ . Evaluate the results for  $y$  in the interval  $[-5, 5]$ , and plot the results in the same form as in c.)

**Solution:**

The implementation can be found in Appendix B, from code listing (2). It produces the following plot:

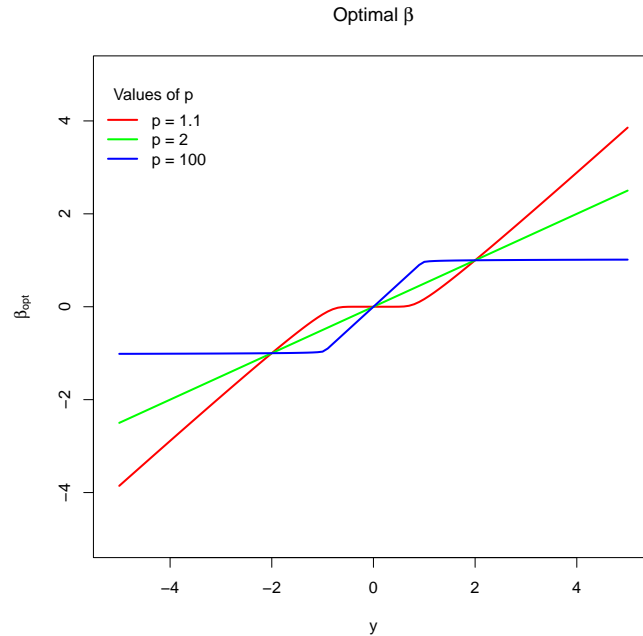


Figure 3: Finding the solution to  $f_{p,\gamma}(\beta_i) = 0$  through the bisection method

Sensible starting values for upper and lower bounds for the bisection method would be  $\mathcal{O}_\beta = |\beta|$  and  $\Omega_\beta = -|\beta|$  respectively. Meaning then for some  $\beta_0 \in (\Omega_\beta, \mathcal{O}_\beta)$  we could be more assured to reach a solution.

e.)

Use the data set `sparseDataWithErrors.dat` and perform estimation of parameters  $\beta_i, i = 1, \dots, n$ , from the data  $y_i, i = 1, \dots, n$ . Compare the result to the ground truth in terms of residual sum of squares. Do the estimation for penalized regression, i.e. compute  $\hat{\beta}_{\gamma,p}$  with  $\gamma = 1, p = 1.1, p = 2$  and  $p = 100$ , compare also to the MLE estimator. For  $p = 100$ , try also to estimate the parameters  $\beta_i$  by using the residuals, i.e.  $\hat{\beta}_{1,100}^{Alt} = y - \hat{\beta}_{1,100}$ . Comment on the results. How does MLE measure up?

---

**Solution:**

The implementation can be found in Appendix B, from code listing (3). The estimation gives:

$p$	$\hat{\beta}_{1,100}$
1.1	273.9674
2.0	1256.5316
100.0	3786.6181

Table 1: Residuals for penalized regression

$p$	$\hat{\beta}_{1,100}^{\text{Alt}}$
1.1	3590.8957
2.0	1256.5316
100.0	251.9418

Table 2: Residuals for Alternative Estimation

The residual sum of squares of the MLE,  $\text{RSS}(\hat{\theta}_{\text{MLE}})$ , gives 998.4152. Comparing the results, we find that the most accurate estimate for  $\hat{\beta}_{1,100}^{\text{Alt}}$  occurs at  $p = 100$ , yielding the lowest residual error. Following closely, the second-best estimate for  $\hat{\beta}_{1,100}$  comes at  $p = 1.1$ .

f.)

In the linear regression problem for  $p = n$ , we have data on the form:

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$$

Where  $\mathbf{y}$  is the  $(n \times 1)$  data vector,  $\boldsymbol{\beta}$  is the  $(n \times 1)$  parameter vector,  $\mathbf{X}$  is a  $(n \times n)$  design matrix, and  $\boldsymbol{\varepsilon}$  is the  $(n \times 1)$  error vector, with  $\boldsymbol{\varepsilon} \sim \mathcal{N}(\mathbf{0}, I)$ . How can you use the ADMM algorithm together with the solution to the problem above to derive a solution to the penalized regression problem?

$$\min_{\boldsymbol{\beta}} \left\{ -\ell(\boldsymbol{\beta}|\mathbf{y}, \mathbf{X}) + \frac{\gamma}{p} \|\boldsymbol{\beta}\|_p^p \right\}$$

**Solution:**

To derive a solution through the ADMM algorithm for the objective function

$$\min_{\boldsymbol{\beta}} \left\{ -\ell(\boldsymbol{\beta}|\mathbf{y}, \mathbf{X}) + \frac{\gamma}{p} \|\boldsymbol{\beta}\|_p^p \right\}$$

where  $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$  and  $\boldsymbol{\varepsilon}$  is a standard normal distributed variable, one can proceed as follows

1. **Objective Function:** Start by defining the objective function, where as we know  $-\ell(\boldsymbol{\beta}|\mathbf{y}, \mathbf{X})$  represents the negative log-likelihood term, but is additionally conditioned on the data  $\mathbf{X}$ :

$$\min_{\boldsymbol{\beta}} \left\{ -\ell(\boldsymbol{\beta}|\mathbf{y}, \mathbf{X}) + \frac{\gamma}{p} \|\boldsymbol{\beta}\|_p^p \right\}$$

2. **Introduction of Auxiliary Variable:** Introduce an auxiliary variable  $\boldsymbol{\lambda}$  such that the problem becomes separable:

$$\begin{aligned} \min_{\boldsymbol{\beta}} \left\{ -\ell(\boldsymbol{\beta}|\mathbf{y}, \mathbf{X}) + \frac{\gamma}{p} \|\boldsymbol{\lambda}\|_p^p \right\} \\ \text{subject to } \boldsymbol{\beta} = \boldsymbol{\lambda} \iff \boldsymbol{\beta} - \boldsymbol{\lambda} = \mathbf{0} \end{aligned}$$



- 
3. **Augmented Lagrangian:** We form the augmented Lagrangian function, combining the objective function and the constraint using a scaled dual variable  $\mathbf{u}$  (or the Lagrange multiplier):

$$L_\rho(\boldsymbol{\beta}, \boldsymbol{\lambda}, \mathbf{u}) = -\ell(\boldsymbol{\beta} | \mathbf{y}, \mathbf{X}) + \frac{\gamma}{p} \|\boldsymbol{\lambda}\|_p^p + \mathbf{u}^T (\boldsymbol{\beta} - \boldsymbol{\lambda}) + \frac{\rho}{2} \|\boldsymbol{\beta} - \boldsymbol{\lambda}\|_2^2$$

4. **ADMM Iterations:** We alternate between updating  $\boldsymbol{\beta}$ ,  $\boldsymbol{\lambda}$ , and  $\mathbf{u}$  until convergence:

- **Update  $\boldsymbol{\beta}$ :**

$$\boldsymbol{\beta}^{(k+1)} = \arg \min_{\boldsymbol{\beta}} \left\{ -\ell(\boldsymbol{\beta} | \mathbf{y}, \mathbf{X}) + \frac{\rho}{2} \|\boldsymbol{\beta} - \boldsymbol{\lambda}^{(k)} + \mathbf{u}^{(k)}\|_2^2 \right\}$$

- **Update  $\boldsymbol{\lambda}$ :**

$$\boldsymbol{\lambda}^{(k+1)} = \arg \min_{\boldsymbol{\lambda}} \left\{ \frac{\gamma}{p} \|\boldsymbol{\lambda}\|_p^p + \frac{\rho}{2} \|\boldsymbol{\beta}^{(k+1)} - \boldsymbol{\lambda} + \mathbf{u}^{(k)}\|_2^2 \right\}$$

- **Update  $\mathbf{u}$ :**

$$\mathbf{u}^{(k+1)} = \mathbf{u}^{(k)} + \rho(\boldsymbol{\beta}^{(k+1)} - \boldsymbol{\lambda}^{(k+1)})$$

5. **Convergence:** We check for convergence based on the magnitude of the primal and dual residuals.

### Problem 2 - (EM-algorithm)

We will now assume that the data  $Y_i, i = 1, \dots, n$ , are independent and identically distributed according to the mixture distribution:

$$f(y_i) = p \cdot \phi(y_i; 0, 1^2) + (1 - p) \cdot \phi(y_i; 0, \tau^2 + 1^2)$$

where  $\phi(y_i; \mu, \sigma^2)$  is the normal density with mean  $\mu$  and variance  $\sigma^2$ . We will now consider estimation of the parameters  $\theta = (p, \tau^2)$ .

a.)

*Give an expression for the likelihood of  $\theta$ .*

**Solution:**

We recall that the likelihood function of a parameter  $\theta$  is given by

$$L(\theta | y_i) = \prod_{i=1}^n f(y_i) \tag{10}$$

where  $f(y_i)$  refers to the probability mass function of  $y_i$ . Hence by direct insertion we have that

$$\begin{aligned} L(\theta | y_i) &= \prod_{i=1}^n f(y_i) \\ &= \prod_{i=1}^n [p \cdot \phi(y_i; 0, 1^2) + (1 - p) \cdot \phi(y_i; 0, \tau^2 + 1^2)] \end{aligned}$$

and we are done. ■

---

b.)

Introduce the variable  $C_i$  which identify which of the two modes  $y_i$  belongs to. Give an expression for the complete log-likelihood using the pairs  $(C_i, y_i)$ ,  $i = 1, \dots, n$

**Solution:**

Given the Likelihood function of our mixture distribution (i.e 2a):

$$L(\theta|y_i) = \prod_{i=1}^n [p \cdot \phi(y_i; 0, 1^2) + (1-p) \cdot \phi(y_i; 0, \tau^2 + 1^2)]$$

We aim to derive the complete log-likelihood function when introducing the variable  $C_i$ , which identifies which of the two modes  $y_i$  belongs to. The indicator variable  $C_i$  is a binary variable, defined accordingly

$$C_i = \begin{cases} 0, & \text{if } y_i \text{ belongs to the first mode,} \\ 1, & \text{if } y_i \text{ belongs to the second mode.} \end{cases}$$

Furthermore, is it understood that  $C_i$  takes on the given values with probabilities

$$\mathbb{P}(C_i) = p \iff C_i = 1 \iff y_i \text{ belongs to the second mode}$$

$$\mathbb{P}(C_i) = 1 - p \iff C_i = 0 \iff y_i \text{ belongs to the first mode}$$

From here we start of by converting the Likelihood function to a log-likelihood

$$\ell(\theta|y) = \log(L(\theta|y)) \tag{11}$$

$$= \log \left( \prod_{i=1}^n [p \cdot \phi(y_i; 0, 1^2) + (1-p) \cdot \phi(y_i; 0, \tau^2 + 1^2)] \right) \tag{12}$$

$$= \sum_{i=1}^n \log [p \cdot \phi(y_i; 0, 1^2) + (1-p) \cdot \phi(y_i; 0, \tau^2 + 1^2)] \tag{13}$$

Observing that for each  $y_i$ , only one term is chosen from the probability density function, based on the mode for that  $y_i$ . Meaning thus that we interpret  $C_i$  as an indicator function. By logarithmically transforming both terms, we obtain the complete log-likelihood function for the observed and missing data pairs  $(C_i, Y_i)$ ,  $i = 1, \dots, n$ , expressed as:

$$\ell(\theta|y_i, C_i) = \sum_{i=1}^n [C_i \cdot \log(p \cdot \phi(y_i; 0, 1^2)) + (1 - C_i) \cdot \log((1-p) \cdot \phi(y_i; 0, \tau^2 + 1^2))]$$

And we are done. ■

c.)

Give an expression for  $Q(\theta | \theta^{(t)})$ , what is the interpretation of  $Q(\theta | \theta^{(t)})$ . Derive the estimates for  $\theta = (p, \tau^2)$ , using  $Q(\theta | \theta^{(t)})$ .

**Solution:**

To rigidly express  $Q(\theta | \theta^{(t)})$ , we start off by looking at the conditional probabilities of  $C_i$  given the observation on  $y_i$ .

The probability that  $C_i = 0$  given the observation on  $y_i$  is:

$$\mathbb{P}(C_i = 0 | y_i, \theta^t) = \frac{(1-p) \cdot \phi(y_i; 0, \tau^2 + 1^2)}{p \cdot \phi(y_i; 0, 1^2) + (1-p) \cdot \phi(y_i; 0, \tau^2 + 1^2)} \tag{14}$$

---

Analogously is the probability that  $C_i = 1$  given the observation on  $y_i$  expressed as:

$$\mathbb{P}(C_i = 1|y_i, \theta^t) = \frac{p \cdot \phi(y_i; 0, 1^2)}{p \cdot \phi(y_i; 0, 1^2) + (1-p) \cdot \phi(y_i; 0, \tau^2 + 1^2)} \quad (15)$$

Henceforth we proceed to compute the expected value of the joint log-likelihood function of the complete dataset, denoted by  $Q(\theta|\theta^{(t)})$

$$Q(\theta|\theta^{(t)}) = \mathbb{E} \left\{ \log L(\theta|y) \mid y, \theta^{(t)} \right\} \quad (16)$$

$$\begin{aligned} &= \sum_{i=1}^n \mathbb{P}(C_i = 1|y_i, \theta^{(t)}) \log(p \cdot \phi(y_i; 0, 1^2)) \\ &\quad + \mathbb{P}(C_i = 0|y_i, \theta^{(t)}) \log((1-p) \cdot \phi(y_i; 0, \tau^2 + 1^2)) \end{aligned} \quad (17)$$

The Q-function computes the expected complete log-likelihood using both the observed data and the current parameter estimates. By maximizing the Q-function, we adjust the parameters to maximize the expected complete log-likelihood.

To derive the estimates for  $\theta$  we need to maximize  $Q(\theta|\theta^{(t)})$  with respect to  $\theta$  itself. By ease of notation we define the maximizer to be  $\theta^{(t+1)} = (p^{(t+1)}, \tau^{2(t+1)})$ .

Initially, we look at the derivative of  $Q$  with respect to  $p$ :

$$\frac{\partial}{\partial p} Q(\theta|\theta^{(t)}) = \frac{\partial}{\partial p} \left[ \sum_{i=1}^n \mathbb{P}(C_i = 1|y_i, \theta^{(t)}) \log(p \cdot \phi(y_i; 0, 1^2)) + \mathbb{P}(C_i = 0|y_i, \theta^{(t)}) \log((1-p) \cdot \phi(y_i; 0, \tau^2 + 1^2)) \right] \quad (18)$$

$$= \sum_{i=1}^n \left[ \frac{\mathbb{P}(C_i = 1|y_i, \theta^{(t)})}{p} - \frac{\mathbb{P}(C_i = 0|y_i, \theta^{(t)})}{1-p} \right] \quad (19)$$

then multiply by  $p(1-p)$

$$\sum_{i=1}^n (1-p) \mathbb{P}(C_i = 1|y_i, \theta^{(t)}) - p \left( \mathbb{P}(C_i = 0|y_i, \theta^{(t)}) \right)$$

which by expansion of the first term yields

$$\sum_{i=1}^n \mathbb{P}(C_i = 1|y_i, \theta^{(t)}) - p \left( \mathbb{P}(C_i = 1|y_i, \theta^{(t)}) \right) - p \left( \mathbb{P}(C_i = 0|y_i, \theta^{(t)}) \right)$$

and through factoring for  $p$

$$\sum_{i=1}^n \mathbb{P}(C_i = 1|y_i, \theta^{(t)}) - p \left( \mathbb{P}(C_i = 1|y_i, \theta^{(t)}) + \mathbb{P}(C_i = 0|y_i, \theta^{(t)}) \right)$$

whereby setting the expression equal to zero conclusively gives

$$\sum_{i=1}^n \mathbb{P}(C_i = 1|y_i, \theta^{(t)}) - p \left( \mathbb{P}(C_i = 1|y_i, \theta^{(t)}) + \mathbb{P}(C_i = 0|y_i, \theta^{(t)}) \right) = 0$$

$\Updownarrow$

$$p \left( \sum_{i=1}^n \mathbb{P}(C_i = 1|y_i, \theta^{(t)}) + \mathbb{P}(C_i = 0|y_i, \theta^{(t)}) \right) = \sum_{i=1}^n \mathbb{P}(C_i = 1|y_i, \theta^{(t)})$$

---

meaning

$$p = \frac{\sum_{i=1}^n \mathbb{P}(C_i = 1 | y_i, \theta^{(t)})}{\sum_{i=1}^n \mathbb{P}(C_i = 1 | y_i, \theta^{(t)}) + \mathbb{P}(C_i = 0 | y_i, \theta^{(t)})}$$

which says that the updating rule on  $p$  is

$$p^{(t+1)} = \frac{1}{n} \sum_1^n \frac{p^{(t)} \cdot \phi(y_i; 0, 1^2)}{p^{(t)} \cdot \phi(y_i; 0, 1^2) + (1 - p^{(t)}) \cdot \phi(y_i; 0, \tau^2 + 1^2)}. \quad (20)$$

Similarly, we proceed by differentiating  $Q$  with respect to  $\tau^2$ :

$$\frac{\partial}{\partial \tau^2} Q(\theta | \theta^{(t)}) = \sum_{i=1}^n (1 - p) \left( \frac{y_i^2}{(\tau^2 + 1)^2} - \frac{1}{\tau^2 + 1} \right) \quad (21)$$

<sup>1</sup> and equating to zero

$$\begin{aligned} \sum_{i=1}^n (1 - p) \left( \frac{y_i^2}{(\tau^2 + 1)^2} - \frac{1}{\tau^2 + 1} \right) &= 0 \\ \sum_{i=1}^n (1 - p) \left( \frac{y_i^2 - (\tau^2 + 1)}{(\tau^2 + 1)^2} \right) &= 0 \\ \Downarrow \\ \sum_{i=1}^n (1 - p) (y_i^2 - (\tau^2 + 1)) &= 0 \\ \Updownarrow \\ \sum_{i=1}^n (1 - p) y_i^2 &= \sum_{i=1}^n (1 - p) (\tau^2 + 1) \\ \Downarrow \\ \tau^2 + 1 &= \frac{\sum_{i=1}^n (1 - p) y_i^2}{\sum_{i=1}^n 1 - p} \end{aligned}$$

thus yielding

$$\tau^2 = \frac{\sum_{i=1}^n (1 - p) y_i^2}{\sum_{i=1}^n 1 - p} - 1$$

Meaning the updating rule on  $\tau^2$  is then

$$\tau^{2(t+1)} = \frac{\sum_{i=1}^n (1 - p) y_i^2}{\sum_{i=1}^n 1 - p} - 1 \quad (22)$$

and we are done. ■

**d.)**

*Implement the solution you derived in c) as a function and apply it to the data: `sparseDataWithErrors.dat`*

---

<sup>1</sup>This result comes from the fact that the Gaussian density function  $\phi(\cdot; \mu, \sigma^2)$  has the derivative:

$$\frac{d}{d\sigma^2} \phi(x; \mu, \sigma^2) = \frac{(x - \mu)^2}{\sigma^4} \phi(x; \mu, \sigma^2) - \frac{1}{\sigma^2} \phi(x; \mu, \sigma^2)$$

---

*What is a good initialization?*

**Solution:**

A natural assumption on initialization is to assume that the modes are equal and that the data points are evenly distributed between them. Hence set  $p = 0.5$  and  $\tau^2 = 0$ . More generally we could assume an initialization of  $p \in (0, 1)$  (due to it being a probability), and  $\tau^2 \neq 1$  (to avoid identical distribution).

The implementation can be found in Appendix B, from code listing (4). The code produces the following result:

$$(p, \tau^2) = (0.946004, 75.005977)$$

where a section of the iteration show the following progression

$i$	$p$	$\tau^2$
1	0.5	3
5	0.862 589	30.258 267
10	0.944 501	73.045 733
15	0.945 983	74.977 963
20	0.946 004	75.005 588
30	0.946 004	75.005 977

Table 3: Iterative estimation of  $p$  and  $\tau^2$  where  $i$  denotes the number of iterations

As presented above we see that even for quite low values of  $p$  and  $\tau^2$ , the algorithm converges quickly within around 30 iterations, which of course can be reduced even more given that the starting values are closer to the realised ones.

**e.)**

*Compute a bootstrap estimate of the uncertainty of the two parameters, using the function from d. Sample  $B = 1000$  times and display scatter plot of the values.*

**Solution:**

The implementation can be found in Appendix B, from code listing (5). The code produces the following plot:

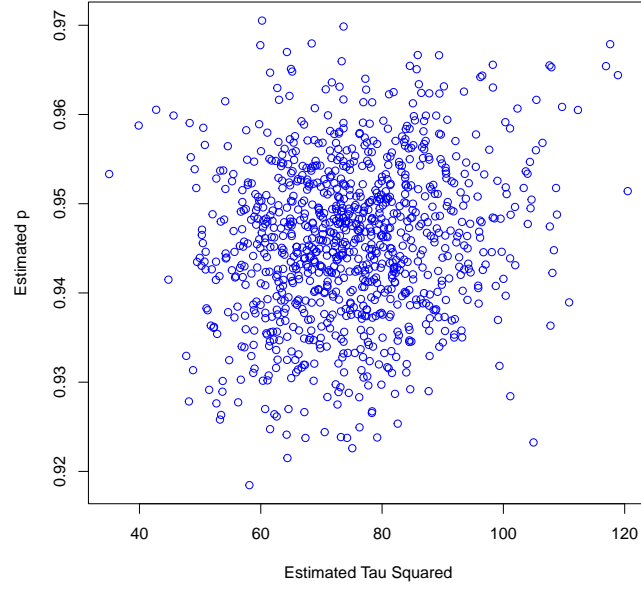


Figure 4: Scatter plot of bootstrap estimations for  $\tau^2$  and  $p$

whereby constructing 95 % confidence intervals for  $p$  and  $\tau^2$  yield

$$\mathbb{P}_p (0.9278386 \leq p \leq 0.9631131) = 95\% \quad (23)$$

$$\mathbb{P}_{\tau^2} (51.4844 \leq \tau^2 \leq 103.3594) = 95\% \quad (24)$$

It is observed that there is notably higher uncertainty associated with  $\hat{\tau}^2$  compared to  $\hat{p}$ , which is understandable given that the predicaments of the 95 % confidence interval for  $p$  admit values of  $p > 0.925$ , implying that fewer than 75 samples belong to mode 2. Moreover, the width of  $\mathbb{P}_{\tau^2}$  naturally tells us that there are clear indications of instability in the estimate of  $\tau^2$ .

f.)

*Compute the observed information matrix. How can you use the observed information matrix to give an uncertainty estimate for  $\theta$ ? Compare the result to e.*

**Solution:**

To approach this computation it is necessary to define exactly what we want to compute, correctly. We are working with a "true" parameter  $\theta$ , a consistent estimate  $\hat{\theta}$  (from 2c and 2d), the expected (fisher) information  $I(\theta)$  at  $\theta$  and the observed (fisher) information  $\mathcal{J}(\theta)$  at  $\theta$ .

What then needs to be established is that all of these quantities are asymptotically equivalent, respectively. This means that given enough samples (data), the observed fisher information will in probability converge to the expected fisher information and likewise the consistent estimate to the

---

true parameter. More mathematically this says that

$$\mathcal{J}(\theta) = \frac{1}{n} \sum_{i=1}^n \frac{\partial^2}{\partial \theta^2} \ell(\theta|y) \quad (\text{Observed Fisher information}) \quad (25)$$

$$\lim_{n \rightarrow \infty} \mathcal{J}(\theta) = \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n \frac{\partial^2}{\partial \theta^2} \ell(\theta|y) \quad (26)$$

$$= \mathbb{E}_\theta \left( \frac{\partial^2}{\partial \theta^2} \ell(\theta|y) \right) \quad (\text{Expected Fisher information}). \quad (27)$$

The assured convergence follows directly from the law of large numbers.

Since we are primarily looking at using the full dataset in this exercise, there will therefore not be a distinct difference between the two measures and can hence express the observation matrix accordingly

$$\begin{aligned} \mathcal{J}(\theta) &= -\nabla \nabla \ell(\theta) \\ &= -H_\ell \end{aligned}$$

where  $H_\ell$  is the Hessian matrix of the log-likelihood with respect to the parameter  $\theta$ . In our case the Hessian is simply a  $2 \times 2$  matrix, with elements

$$H_\ell = \begin{pmatrix} \frac{\partial^2 \ell(\theta)}{\partial p^2} & \frac{\partial^2 \ell(\theta)}{\partial p \partial \tau^2} \\ \frac{\partial^2 \ell(\theta)}{\partial \tau^2 \partial p} & \frac{\partial^2 \ell(\theta)}{\partial (\tau^2)^2} \end{pmatrix}. \quad (28)$$

By construction of  $\ell(\theta|y)$  it is too tedious to compute  $H_\ell$  analytically, and we will numerically approximate it instead. The implementation can be found in Appendix B, from code listings (6). Where the resulting observation matrix becomes

$$\mathcal{J}(\theta) = \begin{pmatrix} 13693.156468 & -1.896655 \\ -1.896655 & 0.004318338 \end{pmatrix}$$

The Negative values admitted in  $\mathcal{J}(\theta)$ <sup>2</sup> for  $\frac{\partial^2 \ell(\theta)}{\partial \tau^2 \partial p}$  indicate potential anomalies in the data, where  $p$  and  $\tau^2$  might show significant correlation. This implies that specific care in interpretation of the diagonal elements, which represent mean square error (due to proportionality), is recommended.

The related covariance matrix, or the inverse of the observation matrix, further establishes the aforementioned:

$$\mathcal{J}^{-1}(\theta) = \text{Cov}(\theta) = \begin{pmatrix} 7.775973 \cdot 10^{-5} & 0.03415281 \\ 0.03415281 & 246.57082498 \end{pmatrix}$$

**g.)**

*Compute the likelihood from 2a, in a dense grid  $(p, \tau^2) \in [0.8, 1] \times [50, 130]$ . Normalize it by dividing by the maximum value and plot it in a contour plot. Use contours lines  $[0.01 \ 0.1 \ 0.5 \ 0.95]$ , mark the ML estimator in the plot. Compare the results to e and f.*

**Solution:**

The implementation can be found in Appendix B, from code listings (7). It produces the following plot

---

<sup>2</sup>I am not completely certain that this is correct, so if not I would fully appreciate a more thorough explanation and result

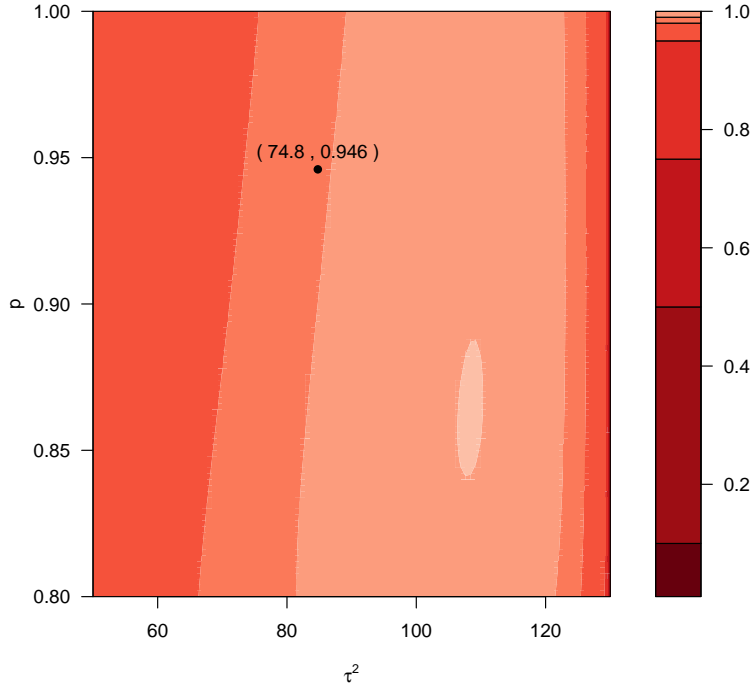


Figure 5: Contour plot of  $L(\theta|y)$  in the grid  $(p, \tau^2) \in [0.8, 1] \times [50, 130]$

which compares well to the estimates from our bootstrap for the MLE.

### Problem 3 - (Bayesian Approach)

Assuming that the data  $Y_i$  are on the same form as in Exercise 1, i.e.

$$y_i = \beta_i + \varepsilon_i, i = 1, \dots, n,$$

where the parameters are as described for (1). Assume also that  $y_i$ , follows the mixture distribution as in Exercise 2,

$$f(y_i) = p \cdot \phi(y_i; 0, 1^2) + (1 - p) \cdot \phi(y_i; 0, \tau^2 + 1^2)$$

a.)

Use a Bayesian interpretation of  $\beta_i$ , (e.g. assume  $\beta_i$  is random) and argue that  $\mathbb{P}(\beta_i = 0 \mid C_i = 0) = 1$ , and that  $f(\beta_i \mid C_i = 1) = \phi(\beta_i; 0, \tau^2)$ . Give an expression for  $\mathbb{P}(\beta_i = 0 \mid y_i = y)$ .

**Solution:**

- **Result 1:**  $\mathbb{P}(\beta_i = 0 \mid C_i = 0) = 1$

By construction of our model we have that

$$Y_i = \beta_i + \varepsilon_i \tag{29}$$



whereby looking at the conditionals with respect to  $C_i = 0$  yields

$$Y_i|(C_i = 0) = \beta_i|(C_i = 0) + \varepsilon_i|(C_i = 0) \quad (30)$$

But the  $\varepsilon_i$  are by assumption independent of the class/mode, and moreover standard normal distributed variables (i.e  $\varepsilon_i \sim \mathcal{N}(0, 1)$ ) meaning thus that both the right and left hand sides of Eq.(30) must follow a standard normal distribution with mean  $\beta_i|(C_i = 0)$  and variance 1. Since  $\beta_i$  and  $\varepsilon_i$  are independent it follows that  $\mathbb{P}(\beta_i = 0|C_i = 0) = 1$ .

• **Result 2:**  $f(\beta_i | C_i = 1) = \phi(\beta_i; 0, \tau^2)$

In like manner we have that

$$Y_i|(C_i = 1) = \beta_i|(C_i = 1) + \varepsilon_i|(C_i = 1) \quad (31)$$

where again  $\varepsilon_i$  are assumed to be independent of the class/mode. Since we now are conditioning on the class/mode  $C_i = 1$ , and know that  $\beta_i$  and  $\varepsilon_i$  are independent, it follows that the only distribution satisfying Eq.(31) is  $\beta_i|(C_i = 1) \sim \mathcal{N}(0, \tau^2)$  since

$$\begin{aligned} f(y_i|C_i = 1) &= f(\beta_i + \varepsilon_i|C_i = 1) \Rightarrow \boxed{f(y_i|C_i = 1) \sim \mathcal{N}(0, \tau^2 + 1)} \\ &= f(\beta_i|C_i = 1) + f(\varepsilon_i|C_i = 1) \\ &= f(\beta_i|C_i = 1) + \mathcal{N}(0, 1) \\ &= \mathcal{N}(0, \tau^2) + \mathcal{N}(0, 1) \\ &= \mathcal{N}(0, \tau^2 + 1). \end{aligned}$$

To express  $\mathbb{P}(\beta_i = 0 | y_i = y)$  we rewrite it as

$$\mathbb{P}(\beta_i = 0 | y_i = y) = \mathbb{P}(\beta_i = 0|C_i = 0, y_i = y)\mathbb{P}(C_i = 0|y_i = y) \quad (32)$$

$$+ \mathbb{P}(\beta_i = 0|C_i = 1, y_i = y)\mathbb{P}(C_i = 1|y_i = y) \quad (33)$$

where we have used the law of total probability<sup>3</sup>. Further, we apply Bayes' rule on the complete expression Eq.(33)

$$\begin{aligned} \mathbb{P}(\beta_i = 0 | y_i = y) &= \mathbb{P}(\beta_i = 0|C_i = 0, y_i = y) \cdot \frac{\mathbb{P}(y_i = y|C_i = 0) \cdot \mathbb{P}(C_i = 0)}{\mathbb{P}(y_i = y)} \\ &+ \mathbb{P}(\beta_i = 0|C_i = 1, y_i = y) \cdot \frac{\mathbb{P}(y_i = y|C_i = 1) \cdot \mathbb{P}(C_i = 1)}{\mathbb{P}(y_i = y)} \\ &= \frac{\mathbb{P}(y_i = y|\beta_i = 0, C_i = 0) \cdot \mathbb{P}(\beta_i = 0|C_i = 0) \cdot \mathbb{P}(C_i = 0)}{\mathbb{P}(y_i = y)} \\ &+ \frac{\mathbb{P}(y_i = y|\beta_i = 0, C_i = 1) \cdot \mathbb{P}(\beta_i = 0|C_i = 1) \cdot \mathbb{P}(C_i = 1)}{\mathbb{P}(y_i = y)} \end{aligned}$$

where we then assume  $y_i = \beta_i + \varepsilon_i$  and independence on  $\beta_i$  and  $\varepsilon_i$ , giving

$$\begin{aligned} \mathbb{P}(\beta_i = 0 | y_i = y) &= \mathbb{P}(\varepsilon_i = y) \cdot \frac{1 \cdot \mathbb{P}(C_i = 0)}{\mathbb{P}(y_i = y)} + 0 \cdot \frac{\mathbb{P}(y_i = y|\beta_i = 0, C_i = 1) \cdot \mathbb{P}(C_i = 1)}{\mathbb{P}(y_i = y)} \\ &= \mathbb{P}(\varepsilon_i = y) \cdot \frac{\mathbb{P}(C_i = 0)}{\mathbb{P}(y_i = y)} \\ &= \phi(y; 0, 1) \cdot \frac{p}{p \cdot \phi(y; 0, 1^2) + (1 - p) \cdot \phi(y; 0, \tau^2 + 1^2)}. \end{aligned}$$

Where we have used that, since  $\mathbb{P}(\beta_i = 0 | C_i = 0) = 1$  it follows that  $\mathbb{P}(\beta_i = 0 | C_i = 1) = 0$ .<sup>4</sup> ■

<sup>3</sup> $P(A) = \sum_i P(A|B_i) \cdot P(B_i)$

<sup>4</sup>Also a consequence of the fact that when  $C_i = 1$ ,  $\beta$  is a continuous distribution.

**NB!** I want the reader to notice that there has been a natural shift in interpretation of the probabilities for the modes  $C_i = 0$  and  $C_i = 1$  in comparison to problem 2. There it was stated that we assume  $\mathbb{P}(C_i = 1) = p$  and  $\mathbb{P}(C_i = 0) = 1 - p$ , which in problem 3 now was considered opposite. This does not by any means change the result.

b.)

Argue that the estimator for  $\beta_i$  (i.e. the conditional expectation of the parameter given the data), is  $\mathbb{P}(C_i = 1 | y_i = y) \mathbb{E}(\beta_i | y_i = y, C_i = 1)$ . Plot this expression as a function of  $y$  in the interval  $[-5, 5]$ , and compare the result with 1.d. Use values:  $p = 0.9$ , and  $\tau^2 = 80$ . Hint:  $\mathbb{E}(\beta_i | y_i = y, C_i = 1) = \frac{\tau^2}{\tau^2 + 1} y$ .

**Solution:**

To proceed, we notice that the problem can be rephrased to showing that

$$\mathbb{E}(\beta_i | y_i) = \mathbb{P}(C_i = 1 | y_i = y) \mathbb{E}(\beta_i | y_i = y, C_i = 1) \quad (34)$$

which is a direct calculation when working with a continuous distribution.

$$\begin{aligned} \mathbb{E}(\beta_i | y_i) &= \int_{-\infty}^{\infty} \beta_i [f(\beta_i | y_i)] d\beta_i \\ &= \int_{-\infty}^{\infty} \beta_i [f(\beta_i | y_i = y, C_i = 0) \mathbb{P}(C_i = 0 | y_i = y) + f(\beta_i | y_i = y, C_i = 1) \mathbb{P}(C_i = 1 | y_i = y)] d\beta_i \\ &= \int_{-\infty}^{\infty} \beta_i [f(\beta_i | y_i = y, C_i = 0) \mathbb{P}(C_i = 0 | y_i = y)] d\beta_i \\ &\quad + \int_{-\infty}^{\infty} \beta_i [f(\beta_i | y_i = y, C_i = 1) \mathbb{P}(C_i = 1 | y_i = y)] d\beta_i \\ &= \mathbb{P}(C_i = 0 | y_i = y) \left[ \int_{-\infty}^{\infty} \beta_i f(\beta_i | y_i = y, C_i = 0) d\beta_i \right] \\ &\quad + \mathbb{P}(C_i = 1 | y_i = y) \left[ \int_{-\infty}^{\infty} \beta_i f(\beta_i | y_i = y, C_i = 1) d\beta_i \right] \\ &= \mathbb{P}(C_i = 0 | y_i = y) \left[ \int_{-\infty}^{\infty} \beta_i \cdot 1 d\beta_i \right] + \mathbb{P}(C_i = 1 | y_i = y) \mathbb{E}(\beta_i | y_i = y, C_i = 1) \\ &= 0 + \mathbb{P}(C_i = 1 | y_i = y) \mathbb{E}(\beta_i | y_i = y, C_i = 1) \\ &= \mathbb{P}(C_i = 1 | y_i = y) \mathbb{E}(\beta_i | y_i = y, C_i = 1) \end{aligned}$$

And we are done. ■

The implementation can be found in Appendix B, from code listings (8), and produces the following plot:

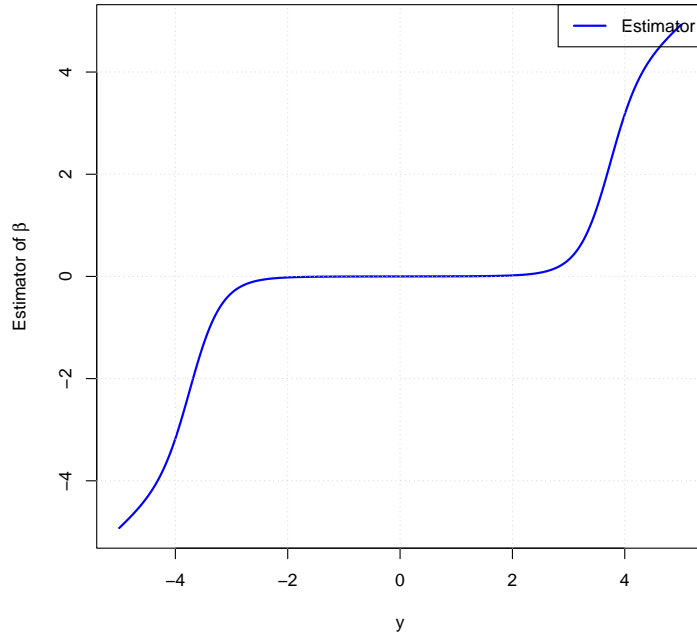


Figure 6: Estimated  $\beta$  for  $y \in [-5, 5]$

The estimator for the  $\beta$ -value shares similarities with one of the lines seen in Figure (2a), of Exercise 1d). This resemblance is notable especially for  $p = 1.1$ , where the optimal  $\beta$  tends to be 0 for  $y \in (-1, 1)$ , though the central region of Figure (6) is wider. An assumption on the trend of Figure (6) could suggest that increasing  $\gamma$  from 1, in Exercise 1d), would generate an optimal  $\beta$  more similar to the one seen above.

c.)

For the dataset in `sparseDataWithErrors.dat`, evaluate the estimator using the values  $p = 0.9$ , and  $\tau^2 = 80$ . Compare the result with 1.e, also in terms of residual sum of squares.

**Solution:**

The implementation can be found in Appendix B, from code listings (9). It produces the following output:

$$e_{\hat{\beta}_{\text{mix}}} = 108.0825$$

In comparison to the estimates from Exercise 1.e, the  $\hat{\beta}_{\text{mix}}$  estimator proves to be the most optimal, by a clear margin.

#### Problem 4 - (Combinatorial optimization)

In this problem we will consider a logistics problem. Assume that you work in a company which want to distribute their product to 20 cities, starting from its home city. In the file `optimaltransport.dat` you will find the locations of all 21 cities the distances between the cities is the time it takes to travel between them. The first city in the file is the home city. We want to find the loop which reduces the traveling time going from the home city, cycling through all cities and returning to the home city. You can modify the scripts found on the home page of the course to solve the problem.

a.)

Write an optimization using simulated annealing to find the optimal solution. Define your neighborhood using mathematical notation. Define your cooling schedule. Argue that your proposed algorithm can reach all possible states.

**Solution:**

In this implementation we have defined the neighbourhood

$$\mathcal{N}(\theta) = \{\theta^*; \exists i, j \text{ such that } \theta_i^* = \theta_j, \theta_j^* = \theta_i \wedge \theta_k^* = \theta_k, \forall k \neq i, j\}$$

which describes the action of swapping or switching to elements/entries in a sequence. For our problem this simply states the action of swapping to towns/cities in our heuristic. More rigidly, we could also define it as follows:

Let  $\theta$  be the tour represented as an ordered sequence of cities:  $\theta = (t_1, t_2, \dots, t_n)$ , where  $t_i$  represents the  $i$ -th town in the tour.

The neighborhood  $\mathcal{N}(\theta)$  of tour  $\theta$  obtained by swapping cities  $i$  and  $j$  is defined as:

$$\mathcal{N}(\theta) = \{\theta^* \mid \theta^* \text{ is obtained by swapping towns } i \text{ and } j \text{ in } \theta\}$$

Which in more syntactic language means:

$$\mathcal{N}(\theta) = \{\theta^* \mid \theta^* = (t_1, \dots, t_{i-1}, t_j, t_{i+1}, \dots, t_{j-1}, t_i, t_{j+1}, \dots, t_n)\}$$

where  $\theta^*$  is a permutation of the original tour  $\theta$  with cities  $i$  and  $j$  swapped.

The algorithm explores all possible states because it can move between any two tours by swapping adjacent cities. This connected structure ensures that every potential tour configuration can be constructed in a finite number of steps.

The code of the implementation can be found in Appendix B, from code listings (10) and produces the following plots

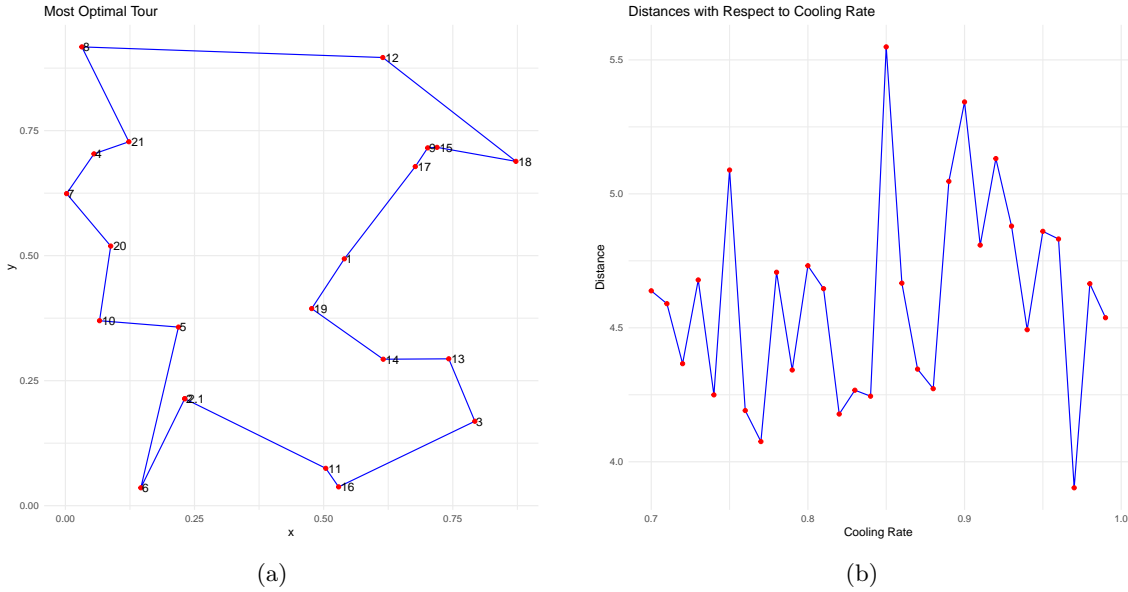


Figure 7: Plots of the most optimal tour and cooling rate for the latest run.

where the all time best tour  $\theta^b$  gave

$$\theta^b = \{17, 9, 15, 18, 12, 8, 21, 4, 7, 20, 10, 5, 2, 6, 11, 16, 3, 13, 14, 19, 1\}$$

$$\text{dist}(\theta^b) = 3.770697$$

$$\alpha = 0.77 \quad (\text{Cooling rate})$$

b.)

*Implement also a TABU search for the optimal path. In what way does both the simulated annealing and the TABU algorithm differ from a steepest decent algorithm?*

**Solution:**

The code of the implementation can be found in Appendix B, from code listings (11) and produces the following plots

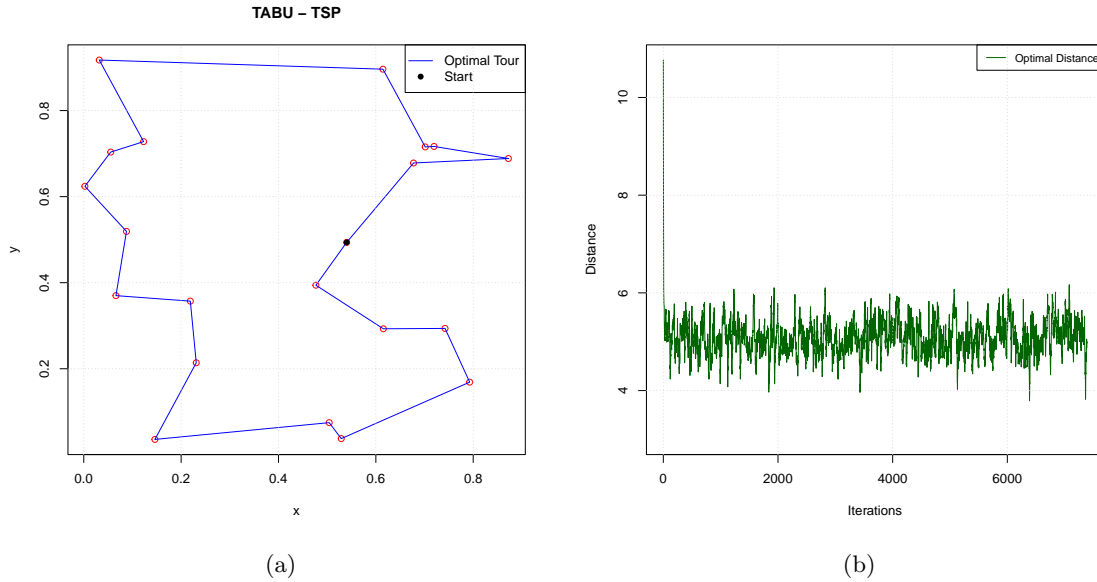


Figure 8: Plots of the most optimal tour and evolution of distance with respect to iterations.

where the all time best tour  $\theta^b$  gave

$$\theta^b = \{1, 14, 13, 3, 16, 11, 6, 2, 5, 7, 4, 21, 12, 9, 8, 15, 18, 10, 17, 20, 19\}$$

$$\text{dist}(\theta^b) = 3.581809$$

Simulated Annealing (SA) and Tabu Search (TS) algorithms differ from steepest descent by prioritizing exploration over exploitation. Steepest descent follows a determined search where it solely relies on the local information of the gradient. In contrast, SA and TS adopt more exploratory rules, therefore inclined to observe new territories in the solution space.

Mathematically, steepest descent iterates through solutions following the gradient direction  $\nabla f(x_n)$  to minimize the objective function  $f(x)$ . SA introduces randomness, with the probability of accepting worse solutions given by  $\mathbb{P}(\Delta E, T) = e^{-\frac{\Delta E}{T}}$ , where  $\Delta E$  is the change in energy and  $T$  is the temperature parameter. TS incorporates a TABU list, preventing revisits to recently explored solutions.

---

While steepest descent relies on deterministic information, SA and TS take a heuristic approach, utilizing randomness and memory to explore a wider solution space. This allows them to break free from local optima encountered by steepest descent and potentially find more optimal solutions elsewhere.

c.)

*If your manager asks you whether the solution you propose is the optimal one, how should you reply? How can you improve confidence in your result?*

**Solution:**

To present our boss with a sensible answer it is initially important to elaborate clearly on how complex the Traveling Salesman problem (TSP) becomes just by minute increases in quantity of data.

The TSP can be mathematically formulated as follows: Given  $n$  cities represented by vertices  $V = \{1, 2, \dots, n\}$  and the distances or costs between them represented by a distance matrix  $C = [c_{ij}]$ , where  $c_{ij}$  denotes the cost of traveling from city  $i$  to city  $j$ , the objective is to find the shortest Hamiltonian cycle that visits each city exactly once and returns to the starting city. A Hamiltonian cycle can be represented as a permutation  $\pi = (\pi_1, \pi_2, \dots, \pi_n)$  of the cities.

The number of possible tours in the TSP is given by the factorial of the number of cities  $n$ . This is because there are  $n$  choices for the first city,  $n - 1$  choices for the second city,  $n - 2$  choices for the third city, and so on, resulting in  $n \times (n - 1) \times (n - 2) \times \dots \times 1 = n!$  possible tours.

$n$	$n!$
1	1
5	120
10	3,628,800
15	1,307,674,368,000
20	2,432,902,008,176,640,000

Table 4:  $n!$  for  $n$  up to 20

However, each tour also has a mirror image (i.e., traversing the tour in reverse), which is considered equivalent. Therefore, the number of unique tours is half of  $(n - 1)!$ , leading to  $\frac{(n-1)!}{2}$  unique tours.

$n$	$\frac{(n-1)!}{2}$
3	1
5	12
10	181,440
15	653,837,184,000
20	1,216,451,004,088,320,000

Table 5:  $\frac{(n-1)!}{2}$  for  $n$  up to 20

The combinatorial explosion of possible tours makes it computationally infeasible to exhaustively search for the optimal solution, especially for large values of  $n$ . This is where heuristic algorithms like Simulated Annealing and TABU Search become influential. They provide efficient methods for exploring the solution space, although without the guarantee of finding the most optimal solution.

To assure our manager, we propose a couple of strategies. Firstly, we could offer a visual representation of the suggested route. By plotting the cities on a map and drawing lines connecting

---

nearby cities, followed by longer lines connecting these regions, the manager can visually assess the proposed route's reasonability. This method offers a straightforward way to confirm that the suggested route is likely close to the best possible route.

Secondly, we could demonstrate our heuristic approaches using a simpler scenario involving, for instance, 10 cities. In this case, we can compare our results with the exact solution obtained through brute-force evaluation of all possible routes. Additionally, we could rerun our algorithm using different initial conditions and present the outcomes in a table to illustrate the variability across runs.

Furthermore, we could evaluate our heuristic methods against alternative approaches, such as the Lin-Kernighan algorithm. The Lin-Kernighan algorithm shows promise as one of the most efficient algorithms when presented with TSP (especially in cases where it is symmetric) and has a time-complexity of  $\mathcal{O}(n^{\lfloor \frac{p}{2} \rfloor})$  where  $p$  represents the backtracking depth. This comparison would provide insights into the effectiveness of our methods relative to established techniques.

d.)

*The company considers buying an additional lorry and want to find out how much this will reduce the distribution time. If the two lorries start at the same time the distribution time is defined as the time it takes until both lorries have returned to the home city.*

*Suggest a modification to your simulated annealing algorithm such that you can have two cycles which both must start and end in the home city. Define a neighborhood for this setup and argue that the algorithm can reach all possible states using your proposed neighborhood. You do not need to implement the solution.*

**Solution:**

Let  $T = \{1, 2, \dots, n\}$  denote the number of towns/cities in our database and let  $H$  represent the home/starting town.

We then allocate  $p$  towns to  $L_1$  (Lorry 1) and the remaining  $n - p$  towns to  $L_2$  (Lorry 2). This implies that  $|T_1| = p$  and  $|T_2| = n - p$  where  $T_1, T_2$  correspond to the cycles of  $L_1$  and  $L_2$  respectively.

For  $T_1$ , the neighborhood  $\mathcal{N}_1(\theta)$  consists of all possible solutions that can be obtained by performing 2-swaps within  $T_1$  (i.e 4a). Mathematically, let  $\theta = (t_1, t_2, \dots, t_p)$  be a solution for  $T_1$ , where  $t_i \in T_1$  represents the city at position  $i$  in the cycle. Then, the neighborhood  $\mathcal{N}_1(\theta)$  is defined as:

$$\mathcal{N}_1(\theta) = \{\theta^* : \theta^* = (t_1^*, t_2^*, \dots, t_p^*), t_i^* \in T_1, \text{ and } t_i^* \neq t_j^* \text{ for } i \neq j\}$$

Similarly, for  $T_2$ , the neighborhood  $\mathcal{N}_2(\theta)$  consists of all possible solutions that can be obtained by performing a 2-swaps within  $T_2$ :

$$\mathcal{N}_2(\theta) = \{\theta' : \theta' = (t'_{p+1}, t'_{p+2}, \dots, t'_n), t'_i \in T_2, \text{ and } t'_i \neq t'_j \text{ for } i \neq j\}$$

Additionally, we introduce an inter-cycle 2-swap, where cities are exchanged between  $T_1$  and  $T_2$ . Mathematically, let  $\theta^s = (t_1^s, t_2^s, \dots, t_n^s)$  be a solution for the entire problem, where  $t_i^s$  represents the city at position  $i$  in the combined cycles. Then, the neighborhood  $\mathcal{N}_{\text{inter}}(\theta^s)$  is defined as:

$$\mathcal{N}_{\text{inter}}(\theta^s) = \{\theta^p : \theta^p = (t_1^p, t_2^p, \dots, t_n^p), t_i^p \in T, \text{ and } t_i^p \neq t_j^p \text{ for } i \neq j\}$$

The nested neighborhood, along with the stochastic effect of the individual cycle 2-swaps (all towns have a probability of getting swapped to a different cycle), ensures the algorithm explores all feasible states within the solution space.

---

### Problem 5 - (Stochastic gradient decent, SGD)

In this exercise we will try out the minibatch approach for optimizing neural nets. We will consider a simple 1D situation. The input  $x$  and output  $y$  are both one dimensional. We will use a model architecture with one hidden layer, and a width of 50. As the activation function,  $\sigma(x)$ , you should use the ReLu, which is defined as:

$$\sigma(x) = \begin{cases} x & \text{for } x > 0 \\ 0 & \text{for } x \leq 0 \end{cases}$$

Define also the derivative of the activation function to be:

$$\sigma'(x) = \begin{cases} 1 & \text{for } x > 0 \\ 0 & \text{for } x \leq 0 \end{cases}$$

The estimator have the form (also called the architecture):

$$f(x) = \sum_{i=1}^{50} \beta_i \sigma(\alpha_i x_i + \alpha_{0,i}) + \beta_0$$

We will minimize the squared sum of errors:

$$\text{sse} = \sum_{i=1}^N \left( y_i - \hat{f}(x_i) \right)^2$$

a.)

*Adapt the formulas from the lecture slides or the SGD note of Geir Storvik such that these can be used for the specific problem formulation defined above. Present the formulas that are relevant for implementing SGD on the current problem. How many parameters are there in the architecture?*

#### Solution:

When dealing with neural networks, it is crucial to optimize the magnitude of the cost or loss function associated with the architecture. This optimization process involves constructing and examining the backpropagation steps, which are responsible for training our model. Backpropagation essentially estimates the parameters of interest in the given model derived from the information of the gradient. This in turn tells us how we should update the parameters.

By this reasoning we proceed to calculate the gradient of our cost function

$$\mathcal{C}(\theta) = \text{sse} = \sum_{i=1}^N \left( y_i - \hat{f}(x_i) \right)^2$$



with respect to all influential parameters  $\beta_0, \beta_i, \alpha_i, \alpha_{0,i}$ . This yields

$$\frac{\partial}{\partial \beta_0} \mathcal{C}(\theta) = -2 \sum_{i=1}^N (y_i - f(x_i)) \quad (35)$$

$$\frac{\partial}{\partial \beta_i} \mathcal{C}(\theta) = -2 \sum_{i=1}^N (y_i - f(x_i)) \sigma(\alpha_i x + \alpha_{0,i}) \quad (36)$$

$$\frac{\partial}{\partial \alpha_i} \mathcal{C}(\theta) = -2 \sum_{i=1}^N (y_i - f(x_i)) \beta_i \sigma'(\alpha_i x + \alpha_{0,i}) x_i \quad (37)$$

$$\frac{\partial}{\partial \alpha_{0,i}} \mathcal{C}(\theta) = -2 \sum_{i=1}^N (y_i - f(x_i)) \beta_i \sigma'(\alpha_i x + \alpha_{0,i}). \quad (38)$$

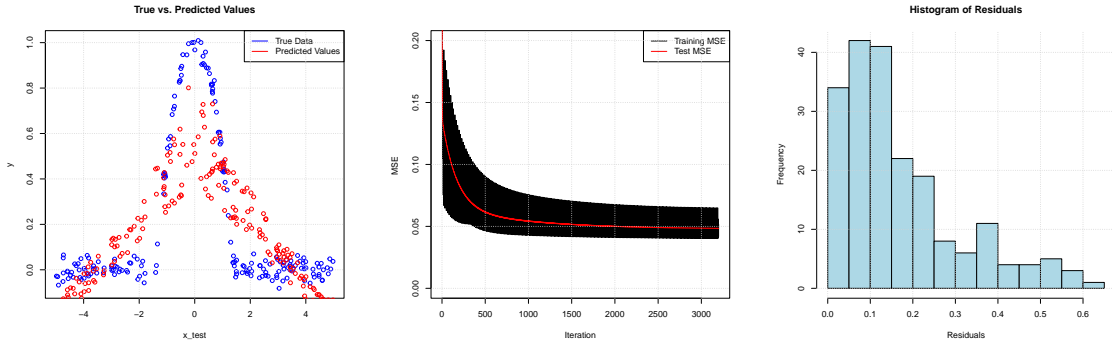
Since we are presented with a network of  $N = 50$  nodes there are hence 50 values of  $\beta_i$ , 50  $\alpha_i$ , 50  $\alpha_{0,i}$  and additionally a singular  $\beta_0$ , giving  $50 + 50 + 50 + 1 = 151$  parameters.

b.)

Implement the SGD for the architecture above and test it on the data found in `functionEstimationNN.dat`. Test different alternatives for the learning rate and record the test error for each epoch. Use a batch size of 50.

**Solution:**

The code of the implementation can be found in Appendix B, from code listings (12) and produces the following plots



(a) Point plot, True vs Predicted

(b) Evolution of MSE

(c) Freq - Histogram, residuals

Figure 9: Various plots of accuracy measurements

c.)

Discuss the choices made and comment on the results from b.

**Solution:**

The results in (b.) show clear promise of an architecture that yield both feasible and decent estimates. We can be assured that for the given data, and presumably data structured similarly in 1D, the network will provide an "accurate"<sup>5</sup> description of its reality. But there are naturally

<sup>5</sup>Accurate is in every sense of its meaning interpreted loosely, especially in statistics, but one might be advised to understand it in this particular instance as a precise regression estimate.

---

a few catches to elaborate on.

Having a neural network with only one layer of 50 nodes (or any number of nodes  $> 1$ ) may not be advised for several reasons. Firstly, the complexity of many real-world problems often require more than one layer to capture the intricate and oftentimes stochastic patterns in the data effectively. A single layer might not have enough capacity to learn complex relationships between inputs and outputs. Additionally, a larger number of nodes in a single layer can lead to overfitting, where the network memorizes the training data instead of learning how to generalise patterns, especially if the dataset is not sufficiently large or sparse.

The ReLU activation function, while being simplistic and capable of diminishing the vanishing gradient problem, may not always be the best choice. Its main limitation is the "dying ReLU" problem, where neurons can become inactive during training. Furthermore, ReLU lacks a bounded output range, which can also lead to exploding activations. It is not suitable for models requiring negative activations, such as some generative models. Alternative functions like Leaky ReLU or ELU might be better in such cases. In general, for overall performance in a neural network, functions like ADAM, tanh, Sigmoid and RMSprop might be preferred.

Lastly, a neural network's sensitivity to learning rate comes from its crucial role in controlling the size of parameter updates during training. If the learning rate is too high, the network may fail to converge or oscillate. Conversely, if it is too low, convergence might be slow. Thus, selecting an optimal learning rate is vital for stable and efficient training. In this respective example we observed the prominent sensitivity which the model had towards its parameters during training, notably for higher learning rates. It was found that on many occasions the gradients would simply explode due to the instability of the hyper-parameters which led to more meticulous grid-searches on the optimal balance between initial learning rate and learning rate decay.

---

## Appendix

### A. GitHub repository

### B. Source Code

Exercise 1 c): Plot the function  $f_{\gamma,p}(\beta)$  for  $\gamma = 1$  and  $\gamma = 0.2$

```
1
2 # Define parameters
3 beta <- seq(-5, 5, length = 101)
4 gamma_list <- c(1, 0.2)
5 p_list <- c(1.1, 2, 5, 100)
6 colors <- c("black", rainbow(length(p_list) - 1, start = 0.2, end = 0.7))
7
8 # Define the function
9 f <- function(beta, gamma, p) {
10   beta + gamma * sign(beta) * abs(beta) ^ (p - 1)
11 }
12
13 # Plotting function
14 plot_function <- function(gamma, title) {
15   plot(beta, rep(0, length(beta)), type = 'n', ylim = c(-8, 8),
16         xlim = c(-5.1, 5.1), xlab = expression(beta), lwd = 2, main =
17         title, ylab = expression(paste("f["p, y]", "(", beta, ")")))
18   for (i in 1:length(p_list)) {
19     lines(beta, f(beta, gamma, p_list[i]), col = colors[i], lwd = 2)
20   }
21   legend('topleft', legend = paste('p = ', p_list), col = colors, lwd = 2,
22         bty = 'n', title = "Values of p", xjust = 0, yjust = 1, inset = c(0,
23         0.05))
24   lines(beta, beta, col = 'red')
25   text(mean(range(beta)), max(beta) * 1.1, labels = expression(paste("f",
26         "(", beta, ") ", " = ", beta)), pos = 3, col = 'red')
27 }
28
29 # Set the size of the plot
30 options(repr.plot.width = 12, repr.plot.height = 8)
31
32 # Plotting for each value of gamma
33 eq1 <- bquote(bold(gamma == .(gamma_list[1])))
34 eq2 <- bquote(bold(gamma == .(gamma_list[2])))
35
36 plot_function(gamma_list[1], eq1)
37 plot_function(gamma_list[2], eq2)
38
39 # Save figures as PDF
40 # pdf(file = "plots.pdf")
41 # plot_function(gamma_list[1], eq1)
42 # plot_function(gamma_list[2], eq2)
43 # dev.off()
```

Listing 1: Exercise 1c

Exercise 1 d): Root of bisection method for  $f_{\gamma,p}(\beta)$

```

1 # Define the function
2 f <- function(beta, gamma, p) {
3   beta + gamma * sign(beta) * abs(beta) ^ (p - 1)
4 }
5
6 # Bisection method for finding the root of an expression
7 bisection_method <- function(y, gamma, p, init_beta1, init_beta2, eps =
8   1e-10, max_iter = 1000) {
9   func <- function(beta) f(beta, gamma, p) - y
10  check_roots <- func(init_beta1) * func(init_beta2)
11  if (check_roots == 0) return(ifelse(func(init_beta1) == 0, init_beta1,
12    init_beta2))
13  if (check_roots > 0) stop('The initial guesses have the same sign of
14    the function (not acceptable)')
15  iter <- 0
16  while (iter < max_iter) {
17    mid_beta <- (init_beta1 + init_beta2) / 2
18    if (func(mid_beta) * func(init_beta1) < 0) init_beta2 <- mid_beta
19    else if (func(mid_beta) * func(init_beta2) < 0) init_beta1 <- mid_beta
20    else return(mid_beta)
21    iter <- iter + 1
22    if (abs(func(mid_beta)) < eps) return(mid_beta)
23  }
24  return(mid_beta)
25 }
26
27 # Test bisection method
28 gamma_val <- 1
29 p_values <- c(1.1, 2, 100)
30 y_values <- seq(-5, 5, length = 101)
31 init_beta1_val <- -5
32 init_beta2_val <- 4.9
33
34 # Plotting function
35 plot_function <- function(p) sapply(y_values, function(y_i)
36   bisection_method(y_i, gamma_val, p, init_beta1_val, init_beta2_val))
37
38 # Plotting
39 plot_colors <- rainbow(length(p_values))
40 plot_legend <- paste('p =', p_values)
41
42 pdf("bisection_plots.pdf")
43 plot(y_values, plot_function(p_values[1]), type='l', col =
44   plot_colors[1], ylim = c(-5, 5),
45   xlim = c(-5.1, 5.1), lwd = 2, main = expression(paste("Optimal", " ",
46     beta)), xlab = 'y', ylab = expression(beta["opt"]))
47 for (i in 2:length(p_values)) lines(y_values, plot_function(p_values[i]),
48   col = plot_colors[i], lwd = 2)
49 legend('topleft', legend = plot_legend, col = plot_colors, lwd = 2, bty =
50   'n', title = "Values of p", xjust = 0, yjust = 1, inset = c(0, 0.05))
51 dev.off()

```

Listing 2: Exercise 1d

Exercise 1 e): Estimation of parameters  $\beta_i$ , sparseDataWithErrors.dat

```

1 # Define the function
2 f <- function(beta, gamma_val, p) {
3   beta + gamma_val * sign(beta) * abs(beta) ^ (p - 1)
4 }
5
6 # Bisection method for finding the root of an expression
7 bisection_method <- function(y, gamma, p, init_beta1, init_beta2, eps =
8   1e-10, max_iter = 1000) {
9   func <- function(beta) f(beta, gamma, p) - y
10  check_roots <- func(init_beta1) * func(init_beta2)
11  if (check_roots == 0) return(ifelse(func(init_beta1) == 0, init_beta1,
12    init_beta2))
13  if (check_roots > 0) stop('The initial guesses have the same sign of
14    the function (not acceptable)')
15  iter <- 0
16  while (iter < max_iter) {
17    mid_beta <- (init_beta1 + init_beta2) / 2
18    if (func(mid_beta) * func(init_beta1) < 0) init_beta2 <- mid_beta
19    else if (func(mid_beta) * func(init_beta2) < 0) init_beta1 <- mid_beta
20    else return(mid_beta)
21    iter <- iter + 1
22    if (abs(func(mid_beta)) < eps) return(mid_beta)
23  }
24  return(mid_beta)
25 }
26
27 # Load data
28 sparseData <- read.table("data/sparseDataWithErrors.ascii", header =
29   FALSE)
30 colnames(sparseData) <- c('betaGT', 'y')
31
32 # Collecting the data from the dataframe
33 y <- sparseData$y
34 betaGT <- sparseData$betaGT
35
36 # We want to estimate beta_optimal for penalized regression
37 gamma_val <- 1
38 p_values <- c(1.1, 2, 100)
39 init_beta1_val <- -100
40 init_beta2_val <- 101
41 residuals <- c()
42 beta_p <- list()
43 # Finding the residuals for different p-values
44 for (p in p_values){
45   estimation <- c()
46   for (y_i in y){
47     estimation <- c(estimation, bisection_method(y_i, gamma_val, p,
48       init_beta1_val, init_beta2_val))
49   }
50
51   # comparing the estimation vs the ground truth
52   beta_p <- c(beta_p, list(estimation))
53   residuals <- c(residuals, sum((estimation - betaGT)^2))
54 }
55
56 residuals <- data.frame(p_values, residuals)
57
58 # Comparison to the MLE estimator, just equal to y
59 beta_MLE <- y
60 residual_MLE <- sum((beta_MLE - betaGT)^2)
61
62 # Estimate beta by using the residuals
63 y_list3 <- rep(list(y), 3)
64
65 beta_alternative <- list()
66 residuals_alternative <- rep(NA, 3)
67 for (i in 1:3) {

```

```

63  beta_alternative <- c(beta_alternative, list(unlist(y_list3[i]) -
        unlist(beta_p[i])))
64  residuals_alternative[i] <- sum((unlist(beta_alternative[i]) -
        betaGT)^2)
65  }
66
67  residuals_alternative_table <- data.frame(p_values, residuals_alternative)
68
69  # Print residuals
70  print(residuals)
71  print(residual_MLE)
72  print(residuals_alternative_table)

```

Listing 3: Exercise 1e

### Exercise 2 d): EM-algorithm, $Q(\theta|\theta^{(t)})$ , sparseDataWithErrors.dat

```

1  # Define the function to calculate p_i
2  calculate_pi <- function(y_i, squared_tau, p) {
3    num <- dnorm(y_i, mean = 0, sd = 1) * p
4    denom <- num + dnorm(y_i, mean = 0, sd = sqrt(squared_tau + 1)) * (1 -
        p)
5    num / denom
6  }
7
8  # Define the EM algorithm function
9  EM_algorithm <- function(y, p_init, squared_tau_init, printing = FALSE,
        eps = 1e-10, max_iter = 1000) {
10   iter <- 1
11   converged <- FALSE
12   p_prev <- p_init
13   squared_tau_prev <- squared_tau_init
14
15   while (!(converged) && iter < max_iter) {
16     p_new <- mean(sapply(y, function(y_i) calculate_pi(y_i,
        squared_tau_prev, p_prev)))
17     squared_tau_new <- sum(sapply(y, function(y_i) {
18       pi_val <- calculate_pi(y_i, squared_tau_prev, p_prev)
19       (1 - pi_val) * y_i^2
20     }) / sum(sapply(y, function(y_i) 1 - calculate_pi(y_i,
        squared_tau_prev, p_prev))) - 1
21
22     # Check for convergence
23     if (abs(p_new - p_prev) + abs(squared_tau_new - squared_tau_prev) <
        eps) {
24       converged <- TRUE
25     }
26
27     if (printing) {
28       print(sprintf('iter: %2.0f, %f %f', iter, p_prev, squared_tau_prev))
29     }
30
31     # Update parameters
32     p_prev <- p_new
33     squared_tau_prev <- squared_tau_new
34
35     iter <- iter + 1
36   }
37
38   c(p_new, squared_tau_new)
39 }

```

```

40
41 # Read the sparse data
42 sparseData <- read.table("data/sparseDataWithErrors.ascii", header =
  FALSE)
43 colnames(sparseData) <- c('betaGT', 'y')
44
45 # Extract data from the dataframe
46 y <- sparseData$y
47
48 # Initialize parameters
49 squared_tau_init <- 3
50 p_init <- 0.5
51
52 # Run the EM algorithm and collect results
53 values <- EM_algorithm(y, p_init, squared_tau_init, printing = TRUE)
54 p_opt <- values[1]
55 squared_tau_opt <- values[2]

```

Listing 4: Exercise 2d

### Exercise 2 e): Bootstrap estimate

```

1 # Set seed for reproducibility
2 set.seed(123)
3
4 # Read sparse data
5 sparse_data <- read.table("data/sparseDataWithErrors.ascii", header =
  FALSE)
6 colnames(sparse_data) <- c('betaGT', 'y')
7
8 # Extract data
9 y <- sparse_data$y
10
11 # Define function to calculate p_i
12 calculate_pi <- function(y_i, tau_squared, p) {
13   numerator <- dnorm(y_i, mean = 0, sd = 1) * p
14   denominator <- numerator + dnorm(y_i, mean = 0, sd = sqrt(tau_squared +
    1)) * (1 - p)
15   numerator / denominator
16 }
17
18 # Define EM algorithm function
19 EM_algorithm <- function(y, p_init, tau_squared_init, eps = 1e-10,
  max_iter = 1000) {
20   iter <- 1
21   converged <- FALSE
22   p_prev <- p_init
23   tau_squared_prev <- tau_squared_init
24
25   while (!(converged) && iter < max_iter) {
26     p_new <- mean(sapply(y, function(y_i) calculate_pi(y_i,
      tau_squared_prev, p_prev)))
27     tau_squared_new <- sum(sapply(y, function(y_i) {
28       pi_val <- calculate_pi(y_i, tau_squared_prev, p_prev)
29       (1 - pi_val) * y_i^2
30     }) / sum(sapply(y, function(y_i) 1 - calculate_pi(y_i,
      tau_squared_prev, p_prev))) - 1
31
32     # Check for convergence
33     if (abs(p_new - p_prev) + abs(tau_squared_new - tau_squared_prev) <
      eps) {

```

```

34     converged <- TRUE
35   }
36
37   # Update parameters
38   p_prev <- p_new
39   tau_squared_prev <- tau_squared_new
40
41   iter <- iter + 1
42 }
43
44 c(p_new, tau_squared_new)
45 }
46
47 # Initialization
48 tau_squared_init <- mean(y)
49 p_init <- 0.5
50
51 # Run EM algorithm
52 values <- EM_algorithm(y, p_init, tau_squared_init)
53 p_opt <- values[1]
54 tau_squared_opt <- values[2]
55
56 # Bootstrap method
57 B <- 1000 # Increase to 1000
58 n <- length(y)
59 p_simulated <- numeric(B)
60 tau_squared_simulated <- numeric(B)
61
62 for (b in 1:B) {
63   y_sample <- sample(y, n, replace = TRUE)
64   values <- EM_algorithm(y_sample, p_init, tau_squared_init)
65   p_simulated[b] <- values[1]
66   tau_squared_simulated[b] <- values[2]
67   if(b %% 100 == 0){
68     print(b)
69   }
70 }
71
72 # Plot results and save to PDF with blue color
73 pdf("bootstrap_results.pdf")
74 plot(tau_squared_simulated, p_simulated, xlab = "Estimated Tau Squared",
75      ylab = "Estimated p", col = "blue")
76 dev.off()
77
78 # Bias
79 cat("Bias:\n")
80 print(mean(p_simulated) - p_opt)
81 print(mean(tau_squared_simulated) - tau_squared_opt)
82
83 # Standard error
84 cat("Standard error:\n")
85 print(sd(p_simulated))
86 print(sd(tau_squared_simulated))
87
88 # 95% confidence interval for p_simulated
89 p_conf_int <- quantile(p_simulated, c(0.025, 0.975))
90 cat("95% Confidence interval for p_simulated:\n")
91 print(p_conf_int)
92
93 # 95% confidence interval for tau_squared_simulated
94 tau_squared_conf_int <- quantile(tau_squared_simulated, c(0.025, 0.975))
95 cat("95% Confidence interval for tau_squared_simulated:\n")
96 print(tau_squared_conf_int)

```

Listing 5: Exercise 2e



---

Exercise 2 f): Observed information:  $-H_\ell$

```
1 set.seed(123)
2
3 # Read sparse data
4 sparse_data <- read.table("data/sparseDataWithErrors.ascii", header =
5   FALSE)
6
7 # Extract data
8 y <- sparse_data$y
9
10 # Define the log-likelihood function
11 log_likelihood <- function(params, y) {
12   p <- params[1]
13   tau_squared <- params[2]
14
15   log_term <- sum(log(p * dnorm(y, mean = 0, sd = 1) + (1 - p) * dnorm(y,
16     mean = 0, sd = sqrt(tau_squared + 1))))
17   return(-log_term) # Note: We negate the log-likelihood for minimization
18 }
19
20 # Set the parameters and data
21 params <- c(p = 0.6, tau_squared = 3)
22
23 # Define a function to compute the negative log-likelihood
24 neg_log_likelihood <- function(params) {
25   return(log_likelihood(params, y))
26 }
27
28 # Use the `optim()` function to estimate the observed Fisher information
29 observed_fisher_information <- function(params, y) {
30   hessian <- optim(params, neg_log_likelihood, method = "BFGS", hessian =
31     TRUE, control = list(trace = 0))$hessian
32   fisher_info <- hessian
33   return(fisher_info)
34 }
35
36 # Compute the observed Fisher information
37 fisher_info <- observed_fisher_information(params, y)
38 inv_fisher <- solve(fisher_info)
39 print(fisher_info)
40 print(inv_fisher)
```

Listing 6: Exercise 2f

## Exercise 2 g): Contour of likelihood: $\ell(\theta|y)$

```

1 # Read sparse data and rename columns
2 data_df <- read.table("data/sparseDataWithErrors.ascii", header = FALSE)
3 colnames(data_df) <- c('beta_GT', 'y')
4
5 # Extract data
6 y <- data_df$y
7
8 # Define the log-likelihood function
9 likelihood <- function(y, p, tau_sq) {
10   fy <- p * dnorm(y, mean = 0, sd = 1) + (1 - p) * dnorm(y, mean = 0, sd
11     = sqrt(tau_sq + 1))
12   return(sum(log(fy)))
13 }
14
15 # Define the grid for p and tau_sq
16 p_grid <- seq(0.8, 1, length.out = 101)
17 tau_sq_grid <- seq(50, 130, length.out = 101)
18
19 # Calculate the log-likelihood values on the grid
20 z <- outer(p_grid, tau_sq_grid, Vectorize(function(p, tau_sq)
21   likelihood(y, p, tau_sq)))
22
23 # Normalize the log-likelihood values
24 z <- max(z) / z
25
26 # Define contour levels
27 levels_ <- c(0.01, 0.1, 0.5, 0.75, 0.95, 0.98, 0.99, 0.9999, 1)
28
29 # Plot the contour plot with the reversed custom red color palette
30 pdf("contour_plot.pdf")
31 #contour(tau_sq_grid, p_grid, z, levels = levels_, xlab = "squared_tau",
32   ylab = "p", col = red_palette(length(levels_)))
33 #filled.contour(tau_sq_grid, p_grid, z, ylab = 'p', xlab =
34   expression(tau^2), levels = levels_, col =
35   red_palette(length(levels_)))
36
37 # Find the ML estimator and mark it in the plot
38 ML_idx <- which(z == max(z), arr.ind = TRUE)
39 points(tau_sq_grid[ML_idx[2]], p_grid[ML_idx[1]], col = "black", pch = 16)
40
41 # Add label for ML estimator coordinates
42 text(tau_sq_grid[ML_idx[2]], p_grid[ML_idx[1]], labels = paste("(",
43   tau_sq_grid[ML_idx[2]], ",", p_grid[ML_idx[1]], ")"), pos = 3)
44 dev.off()

```

Listing 7: Exercise 2g

---

**Exercise 3 b): Plot of conditional expectation:**

$$\mathbb{E}(\beta_i | y_i) = \mathbb{P}(C_i = 1 | y_i = y) \mathbb{E}(\beta_i | y_i = y, C_i = 1)$$

```
1 # Function to calculate P(C_i=1 | y_i=y)
2 calc_cond_prob <- function(y, p, tau) {
3   num <- p * dnorm(y, 0, 1)
4   den <- num + (1 - p) * dnorm(y, 0, tau + 1)
5   p_0 <- num / den
6   return(1 - p_0)
7 }
8
9 # Function to calculate E(beta | y_i=y)
10 calc_expectation <- function(y, tau) {
11   e <- y * tau / (tau + 1)
12   return(e)
13 }
14
15 # Function to calculate beta estimator given y_i=y
16 calc_beta_hat <- function(y, p, tau) {
17   p_c1 <- calc_cond_prob(y, p, tau)
18   e <- calc_expectation(y, tau)
19   beta_hat <- e * p_c1
20   return(beta_hat)
21 }
22
23 # Generate sequence of y values
24 y_values <- seq(-5, 5, length.out = 501)
25 p <- 0.90
26 tau <- 80
27
28 # Calculate beta estimator for each y value
29 beta_hat_y <- calc_beta_hat(y_values, p, tau)
30
31 # Plotting
32 pdf("estimator_plot.pdf") # Open PDF device
33 plot(y_values, beta_hat_y, type = "l", col = "blue", lwd = 2, xlab = "y",
34      ylab = expression(paste("Estimator of", " ", beta)), ylim =
35      range(beta_hat_y))
36 legend("topright", legend = "Estimator", col = "blue", lty = 1, lwd = 2)
37 grid()
38 dev.off() # Close PDF device
```

Listing 8: Exercise 3b

**Exercise 3 c) Evaluation of estimator: sparseDataWithErrors.dat,  $p = 0.9$  and  $\tau^2 = 80$** 

```
1 # Read the sparse data
2 sparseData <- read.table("data/sparseDataWithErrors.ascii", header =
3   FALSE)
4 colnames(sparseData) <- c('betaGT', 'y')
5 betaGT <- sparseData$betaGT
6 # Extract data from the dataframe
7 y <- sparseData$y
8
9 # Function to calculate P(C_i=1 | y_i=y)
10 calc_cond_prob <- function(y, p, tau) {
11   num <- p * dnorm(y, 0, 1)
12   den <- num + (1 - p) * dnorm(y, 0, tau + 1)
13   p_0 <- num / den
14   return(1 - p_0)
```

```

14 }
15
16 # Function to calculate E(beta | y_i=y)
17 calc_expectation <- function(y, tau) {
18   e <- y * tau / (tau + 1)
19   return(e)
20 }
21
22 # Function to calculate beta estimator given y_i=y
23 calc_beta_hat <- function(y, p, tau) {
24   p_c1 <- calc_cond_prob(y, p, tau)
25   e <- calc_expectation(y, tau)
26   beta_hat <- e * p_c1
27   return(beta_hat)
28 }
29
30 # Generate sequence of y values
31 p <- 0.90
32 tau <- 80
33
34 # Calculate beta estimator for each y value
35 beta_hat_y <- calc_beta_hat(y, p, tau)
36
37 residuals = sum((beta_hat_y - betaGT)^2)
38
39 print(residuals)

```

Listing 9: Exercise 3c

#### Exercise 4 a): Simulated annealing for Travelling salesman problem (TSP)

```

1 library(ggplot2)
2
3 # Read data
4 optimalTransport <- read.table("data/optimalTransport.ascii", header =
5   FALSE)
6
7 colnames(optimalTransport) <- c('x', 'y')
8
9 set.seed(124)
10 # Define neighborhood function (swap two cities)
11 swap_cities <- function(tour) {
12   n <- length(tour)
13   i <- sample(2:(n - 1), 1)
14   j <- sample(setdiff(1:n, i), 1)
15   new_tour <- tour
16   new_tour[c(i, j)] <- new_tour[c(j, i)]
17   return(new_tour)
18 }
19
20 # Simulated annealing algorithm
21 simulated_annealing <- function(distances, initial_temp, cooling_rate,
22   max_iter) {
23   # Define objective function (total traveling time)
24   total_time <- function(tour, distances) {
25     total <- 0
26     n <- length(tour)
27     for (i in 1:(n - 1)) {
28       total <- total + distances[tour[i], tour[i + 1]]
29     }
30     total <- total + distances[tour[n], tour[1]] # Return to starting
31     city
32     return(total)
33   }
34 }

```

```

29 }
30
31 n <- nrow(distances)
32 current_tour <- sample(1:n)
33 best_tour <- current_tour
34 best_time <- total_time(current_tour, distances)
35 current_temp <- initial_temp
36
37 for (iter in 1:max_iter) {
38   new_tour <- swap_cities(current_tour)
39   # Calculate total time
40   delta <- total_time(new_tour, distances) - total_time(current_tour,
41     distances)
42
43   if (delta < 0 || runif(1) < exp(-delta / current_temp)) {
44     current_tour <- new_tour
45     current_time <- total_time(new_tour, distances)
46     if (current_time < best_time) {
47       best_tour <- new_tour
48       best_time <- current_time
49     }
50   }
51
52   current_temp <- current_temp * cooling_rate
53 }
54
55 return(list(best_tour = best_tour, best_time = best_time))
56 }
57
58 # Calculate distances matrix (Euclidean distance between cities)
59 distances <- as.matrix(dist(cbind(optimalTransport$x,
60   optimalTransport$y)))
61
62 # Set parameters
63 initial_temp <- 10
64 cooling_rates <- seq(0.7, 0.99, by = 0.01) # Range of cooling rates
65 max_iter <- 10000
66
67 # Store results
68 results <- list()
69
70 # Run simulated annealing for each cooling rate
71 for (cooling_rate in cooling_rates) {
72   result <- simulated_annealing(distances, initial_temp, cooling_rate,
73     max_iter)
74   best_tour <- result$best_tour
75   best_time <- result$best_time
76   results[[as.character(cooling_rate)]] <- list(tour = best_tour,
77     total_time = best_time)
78 }
79
80 # Find the best cooling rate
81 best_result <- which.min(sapply(results, function(x) x$total_time))
82 best_cooling_rate <- as.numeric(names(results)[best_result])
83
84 # Output best tour and total time for the best cooling rate
85 cat("Best Cooling Rate:", best_cooling_rate, "\n")
86 cat("Best Tour:", results[[as.character(best_cooling_rate)]]$tour, "\n")
87 cat("Best Total Time:",
88   results[[as.character(best_cooling_rate)]]$total_time, "\n")
89
90 # Plot the most optimal trip
91 optimal_path <- c(results[[as.character(best_cooling_rate)]]$tour,
92   results[[as.character(best_cooling_rate)]]$tour[1])
93 cities <- optimalTransport[optimal_path, ]
94
95 pdf("opt_annealing.pdf")
96 ggplot() +
97   geom_path(data = cities, aes(x = x, y = y), color = "blue") +
98   geom_point(data = cities, aes(x = x, y = y), color = "red") +

```

```

93   geom_text(data = cities, aes(x = x, y = y, label = rownames(cities)),
94             hjust = -0.2, vjust = 0.5) +
95   ggtitle("Most Optimal Tour") +
96   theme_minimal()
97   dev.off()
98   # Plot distances with respect to cooling rates
99   data <- data.frame(cooling_rates = as.numeric(names(results)),
100                    total_time = sapply(results, function(x) x$total_time))
101   pdf("opt_cooling.pdf")
102   ggplot(data, aes(x = cooling_rates, y = total_time)) +
103     geom_line(color = "blue") +
104     geom_point(color = "red") +
105     labs(x = "Cooling Rate", y = "Distance") +
106     ggtitle("Distances with Respect to Cooling Rate") +
107     theme_minimal()
108
109   dev.off()
110   #Random seed
111   # Best Cooling Rate: 0.77
112   # > cat("Best Tour:", results[[as.character(best_cooling_rate)]]$tour,
113         "\n")
114   # Best Tour: 17 9 15 18 12 8 21 4 7 20 10 5 2 6 11 16 3 13 14 19 1
115   # > cat("Best Total Time:",
116         results[[as.character(best_cooling_rate)]]$total_time, "\n")
117   # Best Total Time: 3.770697
118
119   #Seed(124)
120   # Best Cooling Rate: 0.97
121   # > cat("Best Tour:", results[[as.character(best_cooling_rate)]]$tour,
122         "\n")
123   # Best Tour: 2 11 16 3 13 14 19 1 17 9 15 18 12 8 21 4 7 20 10 5 6
124   # > cat("Best Total Time:",
125         results[[as.character(best_cooling_rate)]]$total_time, "\n")
126   # Best Total Time: 3.902836

```

Listing 10: Exercise 4a

#### Exercise 4 b): TABU - Search for Travelling salesman problem (TSP)

```

1   set.seed(2000)
2   # Read city data
3   city_data <- read.table("data/optimalTransport.ascii", header = FALSE)
4   colnames(city_data) <- c("x", "y")
5
6   # Function to calculate total journey distance
7   calculate_journey_distance <- function(route, distances) {
8     total <- sum(distances[cbind(route[-length(route)], route[-1])])
9     total <- total + distances[route[length(route)], route[1]] # Return to
10    starting point
11    return(total)
12  }
13
14  # Function to perform TABU search
15  find_optimal_route <- function(distances, max_iterations,
16                                tabu_list_length) {
17    num_cities <- nrow(distances) # Number of cities
18    theta <- c(1, sample(2:num_cities, 1)) # Random order and returns to the
19    starting city
20    d <- as.vector(as.matrix(distances))

```

```

19  initial_distance <- calculate_journey_distance(theta, distances) #
    Total distance of the initial route
20  optimal_distance <- initial_distance # Initial optimal distance
21  distance_sequence <- initial_distance # Sequence of distances
22  num_swaps <- num_cities * (num_cities - 1) / 2 - (num_cities - 1) #
    Number of iterations for neighbor search
23  search_table <- combn(2:num_cities, 2) # Table of neighbor search
24  tabu_table <- numeric(0) # Tabu table
25  for (iteration in 1:max_iterations) {
26    optimal_distance2 <- optimal_distance + 10
27    for (i in 1:num_swaps) {
28      if (!(i %in% tabu_table)) {
29        swap_index1 <- search_table[1, i]
30        swap_index2 <- search_table[2, i]
31        new_route <- theta
32        new_route[c(swap_index1, swap_index2)] <- theta[c(swap_index2,
    swap_index1)]
33        new_distance <- calculate_journey_distance(new_route, distances)
34        if (new_distance < optimal_distance2) {
35          optimal_distance2 <- new_distance
36          swap_index <- i
37        }
38      }
39    }
40    theta[c(search_table[1, swap_index], search_table[2, swap_index])] <-
    theta[c(search_table[2, swap_index], search_table[1, swap_index])]
41    optimal_distance <- optimal_distance2
42    distance_sequence <- c(distance_sequence, optimal_distance)
43    tabu_table <- c(tabu_table, swap_index)
44    if (length(tabu_table) > tabu_list_length) {
45      tabu_table <- tabu_table[-1]
46    }
47    if (optimal_distance < initial_distance) {
48      optimal_route <- theta
49      initial_distance <- optimal_distance
50    }
51  }
52  return(list(optimal_route = optimal_route, shortest_distance =
    initial_distance, distance_sequence = distance_sequence))
53 }
54
55
56 # Calculate distances between cities
57 city_distances <- as.matrix(dist(city_data))
58
59 # Set parameters for TABU search
60 max_iterations <- 1000
61 tabu_list_length <- 40
62
63 # Run TABU search
64 result <- find_optimal_route(city_distances, max_iterations,
    tabu_list_length)
65 optimal_route <- result$optimal_route
66 shortest_distance <- result$shortest_distance
67 distance_sequence <- result$distance_sequence
68
69 # Print the results
70 cat("Shortest distance:", shortest_distance, "\n")
71 cat("Optimal route:\n")
72 print(optimal_route)
73 cat("Number of iterations to find optimal route:",
    which.min(distance_sequence), "\n")
74
75 pdf("opt_tabu.pdf")
76 plot(city_data, main = "TABU - TSP", col = "red")
77 lines(city_data[optimal_route, 1], city_data[optimal_route, 2], col =
    "blue")
78 points(city_data[1, ], col = "black", pch = 16)
79 legend("topright", legend = c("Optimal Tour", "Start"), col = c("blue",
    "black"), pch = c(NA, 16), lty=c(1, NA))
80 grid()

```

```

81 dev.off()
82
83 # Plot time series of the shortest distance found over iterations
84 optimal_route_iteration <- which.min(distance_sequence)
85 reasonable_plot_range <- min(optimal_route_iteration + 1000,
86                               max_iterations)
86 pdf("distance_iter.pdf")
87 plot.ts(distance_sequence[1:reasonable_plot_range], xlab = "Iterations",
88          ylab = "Distance", col = "darkgreen", ylim = c(3,
89                max(distance_sequence)))
88 legend("topright", legend = c("Optimal Distance"), col = c("darkgreen"),
89        lty = 1, cex = 0.8)
89 grid()
90 dev.off()
91
92
93
94
95 # Seed 2000
96
97 # > cat("Shortest distance:", shortest_distance, "\n")
98 # Shortest distance: 3.776257
99 # > cat("Optimal route:\n")
100 # Optimal route:
101 #   > print(optimal_route)
102 # [1] 1 19 13 3 16 11 6 2 5 10 20 7 4 21 8 12 9 15 18 17 17 1
103 # > cat("Number of iterations to find optimal route:",
104       which.min(distance_sequence), "\n")
104 # Number of iterations to find optimal route: 927
105
106 #Random seed
107
108 # > cat("Shortest distance:", shortest_distance, "\n")
109 # Shortest distance: 3.732432
110 # > cat("Optimal route:\n")
111 # Optimal route:
112 #   > print(optimal_route)
113 # [1] 1 14 3 18 18 17 17 9 9 15 12 8 21 4 7 20 10 6 2 5 19 1
114 # > cat("Number of iterations to find optimal route:",
115       which.min(distance_sequence), "\n")
115 # Number of iterations to find optimal route: 900
116
117 #All time best random seed
118
119 # > cat("Shortest distance:", shortest_distance, "\n")
120 # Shortest distance: 3.581809
121 # > cat("Optimal route:\n")
122 # Optimal route:
123 #   > print(optimal_route)
124 # [1] 1 14 13 3 16 11 6 2 5 7 4 21 12 9 8 15 18 10 17 20
125 #   19
125 # > cat("Number of iterations to find optimal route:",
126       which.min(distance_sequence), "\n")
126 # Number of iterations to find optimal route: 292

```

Listing 11: Exercise 4b

## Exercise 5 b): Stochastic Gradient Descent

```

1 # Load data
2 dataset <- read.table("data/functionEstimationNN.ascii", header = FALSE,
3                       col.names = c("x", "y"))

```



```

3
4 # Sample a subset of the data
5 set.seed(123) # For reproducibility
6 subset_size <- 1000
7 subset <- dataset[sample(nrow(dataset), subset_size), ]
8
9 # Train-test split
10 train_ratio <- 0.8
11 train_indices <- sample(1:nrow(subset), train_ratio * nrow(subset))
12 train_set <- subset[train_indices, ]
13 test_set <- subset[-train_indices, ]
14
15 # Extract train and test data
16 x_train <- train_set$x
17 y_train <- train_set$y
18 x_test <- test_set$x
19 y_test <- test_set$y
20
21 # Network parameters
22 n_inputs <- 1
23 n_hidden <- 50
24 n_outputs <- 1
25
26 # Initialize weights and biases
27 weights_hidden <- matrix(rnorm(n_inputs * n_hidden, sd =
28   sqrt(2/n_inputs)), nrow = n_inputs, ncol = n_hidden)
29 bias_hidden <- runif(n_hidden, min = 0, max = 1)
30 weights_output <- matrix(rnorm(n_hidden * n_outputs, sd =
31   sqrt(2/n_hidden)), nrow = n_hidden, ncol = n_outputs)
32 bias_output <- runif(n_outputs, min = 0, max = 1)
33
34 # Hyperparameters
35 learning_rate <- 0.001
36 learning_rate_decay <- 1e-4
37 batch_size <- 50
38 n_epochs <- 200
39
40 # Train network with SGD
41 SSE_seq <- MSE_seq <- MSE_test_seq <- NULL
42 iter <- 0
43 for (epoch in 1:n_epochs) {
44   learning_rate <- learning_rate / (1 + epoch * learning_rate_decay)
45   x_shuffled <- x_train
46   y_shuffled <- y_train
47   n_batches <- ceiling(length(x_train) / batch_size)
48   batches <- split(1:length(x_train), rep(1:n_batches, each = batch_size,
49     length.out = length(x_train)))
50
51   for (batch in batches) {
52     iter <- iter + 1
53     x_batch <- x_shuffled[batch]
54     y_batch <- y_shuffled[batch]
55
56     # Forward propagation
57     hidden_output <- pmax(x_batch %*% weights_hidden + bias_hidden, 0) #
58     ReLU activation
59     y_pred_batch <- hidden_output %*% weights_output + bias_output
60
61     # Loss calculation
62     sse_batch <- sum((y_pred_batch - y_batch)^2)
63     mse_batch <- sse_batch / length(batch)
64     SSE_seq <- c(SSE_seq, sse_batch)
65     MSE_seq <- c(MSE_seq, mse_batch)
66
67     # Backward propagation
68     d_y_sse_batch <- -2 * (y_batch - y_pred_batch) / length(batch)
69     d_weights_output <- t(hidden_output) %*% d_y_sse_batch
70     d_bias_output <- colSums(d_y_sse_batch)
71     d_hidden <- d_y_sse_batch %*% t(weights_output) * (hidden_output > 0)
72     d_weights_hidden <- t(x_batch) %*% d_hidden
73     d_bias_hidden <- colSums(d_hidden)

```

```

70
71     # Update weights and biases
72     weights_hidden <- weights_hidden - learning_rate * d_weights_hidden
73     bias_hidden <- bias_hidden - learning_rate * d_bias_hidden
74     weights_output <- weights_output - learning_rate * d_weights_output
75     bias_output <- bias_output - learning_rate * d_bias_output
76
77     # Forward propagation for test data
78     hidden_output_test <- pmax(x_test %% weights_hidden + bias_hidden, 0)
79     y_pred_test <- hidden_output_test %% weights_output + bias_output
80     mse_test <- sum((y_pred_test - y_test)^2) / length(y_test)
81     MSE_test_seq <- c(MSE_test_seq, mse_test)
82 }
83
84 # Print MSE for current epoch
85 cat("Epoch:", epoch, "| MSE:", MSE_seq[length(MSE_seq)], "| Learning
    Rate:", learning_rate, "\n")
86 }
87
88 # Plot MSE for training and test data
89 # Plot MSE for training and test data with y-axis scaled from 0 to 2
90 pdf("MSE_train_test.pdf")
91 plot(MSE_seq, type = "l", xlab = "Iteration", ylab = "MSE ", col =
    "black", ylim = c(0, 0.2))
92 lines(MSE_test_seq, type = "l", col = "red")
93 legend("topright", legend = c("Training MSE", "Test MSE"), col =
    c("black", "red"), lty = 1)
94 grid()
95 dev.off()
96
97 pdf("True_v_Predict.pdf")
98 # Plot true vs. predicted values
99 plot(x_test, y_test, col = "blue", main = "True vs. Predicted Values",
    xlab = "x_test", ylab = "y")
100 points(x_test, y_pred_test, col = "red")
101 legend("topright", legend = c("True Data", "Predicted Values"), col =
    c("blue", "red"), lty = 1)
102 grid()
103 dev.off()
104 # Calculate residuals
105 residuals <- abs(y_test - y_pred_test)
106
107 # # Plot residuals vs. predicted values
108 # plot(y_pred_test, residuals, col = "blue", main = "Residuals vs.
    Predicted Values", xlab = "Predicted Values", ylab = "Residuals")
109 # abline(h = 0, col = "red")
110
111 # Histogram of residuals
112 pdf("Residual_histo.pdf")
113 hist(residuals, col = "lightblue", main = "Histogram of Residuals", xlab
    = "Residuals", ylab = "Frequency")
114 grid()
115 dev.off()
116 # Scatterplot of residuals
117 # plot(y_pred_test, residuals, col = "blue", main = "Scatterplot of
    Residuals", xlab = "Predicted Values", ylab = "Residuals")
118 # abline(h = 0, col = "red")

```

Listing 12: Exercise 5b