

Exercise Sheet 12

Exercise 1: Building Neural Networks (20 + 20 P)

We consider the problem of learning decision functions in \mathbb{R}^2 where $x = (x_1, x_2)$ denotes the two-dimensional input vector. For this exercise, you only have access to neurons of the type

$$a_j = \sigma\left(b_j + \sum_i a_i w_{ij}\right)$$

where σ is the step function, i.e.

$$\sigma(t) = \begin{cases} 1 & \text{if } t > 0 \\ 0 & \text{if } t \leq 0 \end{cases}$$

(a) Construct a neural network that implements the decision boundary below

$$y = \begin{cases} 1 & \text{if } \max(x_1, x_2) > 2 \\ 0 & \text{if } \max(x_1, x_2) < 2 \end{cases}$$

Specifically, *draw* the neural network graph, and *indicate* for each neuron its weights and bias. (The exact behavior at the decision boundary does not need to be enforced.)

(b) Repeat the exercise for the decision function

$$y = \begin{cases} 1 & \text{if } \|x\|_1 > 2 \\ 0 & \text{if } \|x\|_1 < 2 \end{cases}$$

Exercise 2: Condition Number (20 + 10 P)

Consider a supervised dataset composed of inputs $\mathbf{x}_1, \dots, \mathbf{x}_N \in \mathbb{R}^d$ and respective targets $t_1, \dots, t_N \in \mathbb{R}$. Assume that $\frac{1}{N} \sum_{i=1}^N \mathbf{x}_i = \mathbf{0}$, i.e. the data is centered. Consider the homogeneous linear model $y = \mathbf{w}^\top \mathbf{x}$, with $\mathbf{w} \in \mathbb{R}^d$ a vector of parameters to be learned, and the regularized mean square objective:

$$\mathcal{E}(\mathbf{w}) = \alpha \|\mathbf{w}\|^2 + \frac{1}{N} \sum_{i=1}^N (\mathbf{w}^\top \mathbf{x}_i - t)^2$$

we would like to minimize.

(a) *Show* that the Hessian of the error function $\mathcal{E}(\mathbf{w})$ is given by the constant matrix:

$$H(\mathbf{w}) = 2(\Sigma + \alpha I)$$

where Σ is the covariance of the data.

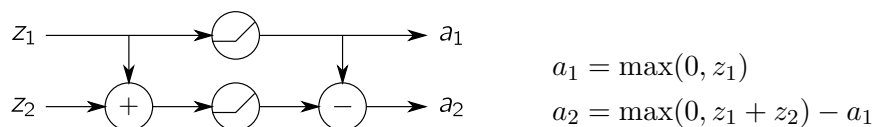
(b) *Show* that the condition number associated to this Hessian matrix is given by

$$c = \frac{\lambda_1 + \alpha}{\lambda_d + \alpha}$$

where $\lambda_1, \dots, \lambda_d$ are the eigenvalues of the matrix Σ sorted in decreasing order.

Exercise 3: Backpropagation (30 P)

Consider some portion of a neural network given by:



(a) Assuming that we know the error gradient $(\partial E / \partial a_1, \partial E / \partial a_2)$, *compute* the error gradient $(\partial E / \partial z_1, \partial E / \partial z_2)$.

Exercise 1: Building Neural Networks (20 + 20 P)

We consider the problem of learning decision functions in \mathbb{R}^2 where $x = (x_1, x_2)$ denotes the two-dimensional input vector. For this exercise, you only have access to neurons of the type

$$a_j = \sigma\left(b_j + \sum_i a_i w_{ij}\right)$$

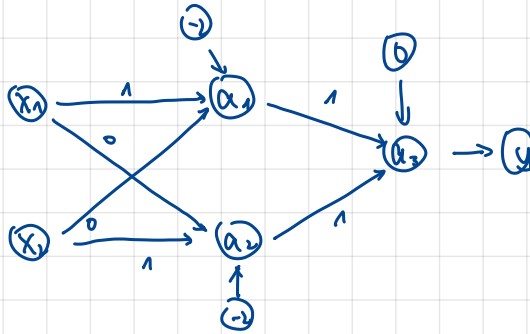
where σ is the step function, i.e.

$$\sigma(t) = \begin{cases} 1 & \text{if } t > 0 \\ 0 & \text{if } t \leq 0 \end{cases}$$

(a) Construct a neural network that implements the decision boundary below

$$y = \begin{cases} 1 & \text{if } \max(x_1, x_2) > 2 \\ 0 & \text{if } \max(x_1, x_2) < 2 \end{cases}$$

Specifically, *draw* the neural network graph, and *indicate* for each neuron its weights and bias. (The exact behavior at the decision boundary does not need to be enforced.)

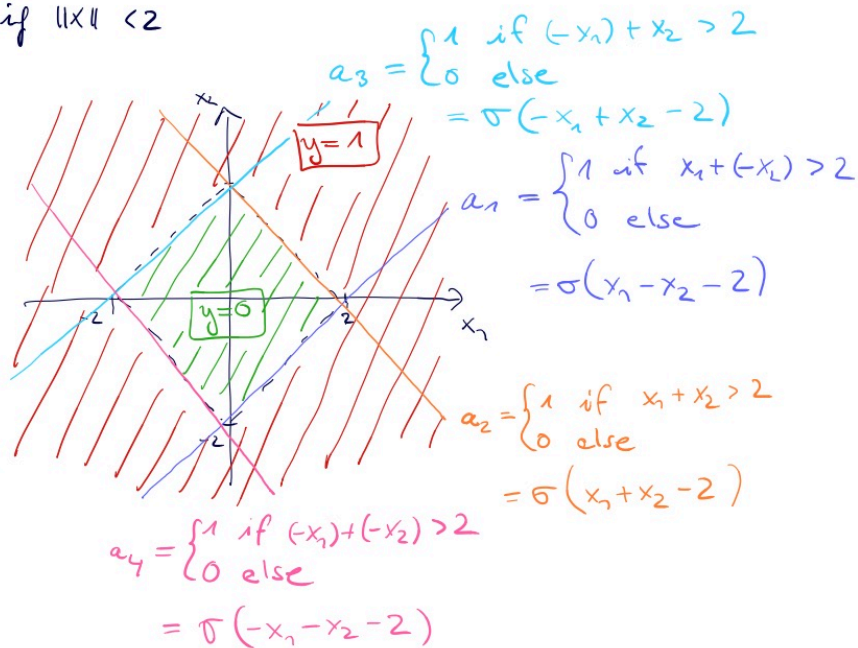


(b) Repeat the exercise for the decision function $\|x\|_1 = \sum |x_i|$

$$y = \begin{cases} 1 & \text{if } \|x\|_1 > 2 \\ 0 & \text{if } \|x\|_1 < 2 \end{cases}$$

b)

$$y = \begin{cases} 1 & \text{if } \|x\| > 2 \\ 0 & \text{if } \|x\| < 2 \end{cases}$$



$$a_1 = \sigma(x_1 - x_2 - 2)$$

$$\Rightarrow w_{11} = 1, w_{21} = -1, b_1 = -2$$

$$a_3 = \sigma(-x_1 + x_2 - 2)$$

$$\Rightarrow w_{13} = -1, w_{23} = 1, b_3 = -2$$

$$a_2 = \sigma(x_1 + x_2 - 2)$$

$$\Rightarrow w_{12} = 1, w_{22} = 1, b_2 = -2$$

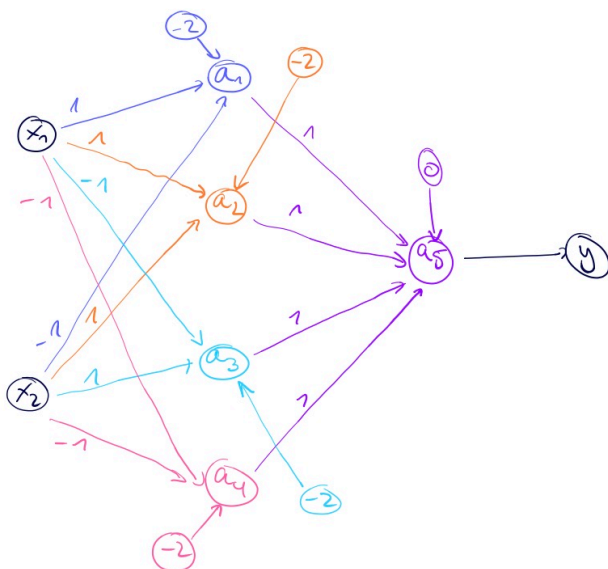
$$a_4 = \sigma(-x_1 - x_2 - 2)$$

$$\Rightarrow w_{14} = -1, w_{24} = -1, b_4 = -2$$

Again only one of $\{a_1, a_2, a_3, a_4\}$ has to be active for the condition to hold. Thus we construct g -as similarly as above:

$$a_5 = \sigma(a_1 + a_2 + a_3 + a_4)$$

$$\Rightarrow w_{15} = w_{25} = w_{35} = w_{45} = 1, b_5 = 0$$



Exercise 2: Condition Number (20 + 10 P)

Consider a supervised dataset composed of inputs $\mathbf{x}_1, \dots, \mathbf{x}_N \in \mathbb{R}^d$ and respective targets $t_1, \dots, t_N \in \mathbb{R}$. Assume that $\frac{1}{N} \sum_{i=1}^N \mathbf{x}_i = \mathbf{0}$, i.e. the data is centered. Consider the homogeneous linear model $y = \mathbf{w}^\top \mathbf{x}$, with $\mathbf{w} \in \mathbb{R}^d$ a vector of parameters to be learned, and the regularized mean square objective:

$$\mathcal{E}(\mathbf{w}) = \alpha \|\mathbf{w}\|^2 + \frac{1}{N} \sum_{i=1}^N (\mathbf{w}^\top \mathbf{x}_i - t_i)^2$$

we would like to minimize.

$$\Sigma^{-1} \mathbf{v}_i = \lambda_i^{-1} \mathbf{v}_i$$

(a) Show that the Hessian of the error function $\mathcal{E}(\mathbf{w})$ is given by the constant matrix:

$$H(\mathbf{w}) = 2(\Sigma + \alpha I)$$

where Σ is the covariance of the data.

$$(\mathbf{I} + \alpha \tilde{\Sigma}) \mathbf{v}_i$$

(b) Show that the condition number associated to this Hessian matrix is given by

$$c = \frac{\lambda_1 + \alpha}{\lambda_d + \alpha}$$

$$\mathbf{v}_i + \alpha \lambda_i^{-1} \mathbf{v}_i$$

where $\lambda_1, \dots, \lambda_d$ are the eigenvalues of the matrix Σ sorted in decreasing order.

$$1 + \alpha \lambda_i^{-1}$$

$$\begin{aligned} \text{a)} \quad \frac{\partial \mathcal{E}(\mathbf{w})}{\partial \mathbf{w}} &= 2\alpha \mathbf{w} + \frac{1}{N} \sum_{i=1}^N 2(\mathbf{w}^\top \mathbf{x}_i - t_i) \mathbf{x}_i \\ H(\mathbf{w}) = \frac{\partial^2 \mathcal{E}(\mathbf{w})}{\partial \mathbf{w}^2} &= 2\alpha + \frac{1}{N} \sum_{i=1}^N 2 \mathbf{x}_i \mathbf{x}_i^\top \\ &= 2(\alpha \mathbf{I} + \Sigma) \end{aligned}$$

$$\begin{aligned} \text{b)} \quad \lambda &= \mathbf{v}^\top H \mathbf{v} = 2 \mathbf{v}^\top (\alpha \mathbf{I} + \Sigma) \mathbf{v} \quad \text{with } \|\mathbf{v}\| = 1 \\ &\quad \uparrow \\ &\quad \text{for } H \\ &= 2\alpha \cdot \mathbf{v}^\top \mathbf{v} + 2 \mathbf{v}^\top \Sigma \mathbf{v} \\ &= 2\alpha \|\mathbf{v}\|^2 + 2 \mathbf{v}^\top \Sigma \mathbf{v} \\ &\quad \quad \quad \parallel \text{ for } \Sigma \\ &= 2(\alpha + \lambda_i) \end{aligned}$$

$$c = \frac{\lambda_{\max}}{\lambda_{\min}} = \frac{2(\alpha + \lambda_1)}{2(\alpha + \lambda_d)}$$

b) Let $\hat{\lambda}_1, \dots, \hat{\lambda}_d$ be the eigenvalues of $2(\Sigma + \alpha \mathbf{I})$ sorted in decreasing order. Then we have

$$c = \frac{\hat{\lambda}_1}{\hat{\lambda}_d}$$

We will now show that $\hat{\lambda}_i = k \cdot (\lambda_i + \alpha)$ for some constant k .

Let \mathbf{v}_i be an eigenvector of Σ with corresponding eigenvalue λ_i , i.e. $\Sigma \mathbf{v}_i = \lambda_i \mathbf{v}_i$. Then, \mathbf{v}_i is also an eigenvector of $2(\Sigma + \alpha \mathbf{I})$:

$$\begin{aligned} 2(\Sigma + \alpha \mathbf{I}) \mathbf{v}_i &= 2 \underbrace{\Sigma \mathbf{v}_i}_{= \lambda_i \mathbf{v}_i} + 2\alpha \underbrace{\mathbf{I} \mathbf{v}_i}_{= \mathbf{v}_i} \\ &= 2\lambda_i \mathbf{v}_i + 2\alpha \mathbf{v}_i \\ &= 2(\lambda_i + \alpha) \mathbf{v}_i \\ &\quad \quad \quad = \hat{\lambda}_i \end{aligned}$$

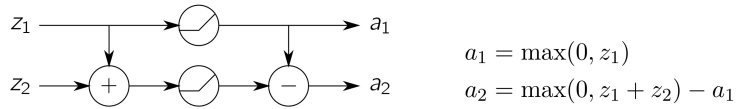
with eigenvalue $\hat{\lambda}_i = 2(\lambda_i + \alpha)$

Thus, we have

$$c = \frac{\hat{\lambda}_1}{\hat{\lambda}_d} = \frac{2(\lambda_1 + \alpha)}{2(\lambda_d + \alpha)} = \frac{\lambda_1 + \alpha}{\lambda_d + \alpha}$$

Exercise 3: Backpropagation (30 P)

Consider some portion of a neural network given by:



(a) Assuming that we know the error gradient $(\partial E / \partial a_1, \partial E / \partial a_2)$, compute the error gradient $(\partial E / \partial z_1, \partial E / \partial z_2)$.

$$\begin{aligned}\frac{\partial E}{\partial z_1} &= \frac{\partial E}{\partial a_1} \frac{\partial a_1}{\partial z_1} + \frac{\partial E}{\partial a_2} \frac{\partial (\max(0, z_1 + z_2) - a_1)}{\partial z_1} \\ &= \frac{\partial E}{\partial a_1} \frac{\partial a_1}{\partial z_1} + \frac{\partial E}{\partial a_2} \left(\frac{\partial \max(0, z_1 + z_2)}{\partial z_1} - \frac{\partial a_1}{\partial z_1} \right) \\ &= \frac{\partial E}{\partial a_1} 1_{z_1 > 0} + \frac{\partial E}{\partial a_2} \left(1_{z_1 + z_2 > 0} - 1_{z_1 > 0} \right)\end{aligned}$$

$$\frac{\partial E}{\partial z_2} = \frac{\partial E}{\partial a_2} \frac{\partial a_2}{\partial z_2} = \frac{\partial E}{\partial a_2} \left(1_{z_1 + z_2 > 0} - 1_{z_1 > 0} \right)$$