# Exercise A.

1. (a) $E = \eta \|w\|_F^2 + \sum_{i=1}^{N} \|x_i - ws_i\|^2 + \lambda \|s_i\|_1 \quad \forall_{i=1}^{N} \quad s_i \geq 0$

$$\frac{\partial E}{\partial w} = 2\eta w + \sum_{i=1}^{N} 2(x_i - ws_i)(-s_i) = 0.$$

1. (b) $\frac{\partial E}{\partial s_i} = 2(x_i - ws_i)(-w) + \lambda \, \text{sign}(s_i) = 0.$

$\text{sign}(s_i) = 1 \quad \text{as} \quad s_i \geq 0.$

$$\frac{\partial E}{\partial s_i} = 2(x_i - ws_i)(-w) + \lambda \mathbb{1} = 0.$$

ML2 E54 T2a)

$$\|r_i\|^2 \overset{!}{=} \|S_i\|_1 \qquad \text{with } S_i = g(r_i)$$

$$\Rightarrow \|r_i\|^2 = \|g(r_i)\|_1 \qquad \text{with } \|g(r_i)\|_1 = \sum_{i=1}^{N} |g(r_i)|$$

$$\Rightarrow \sum_{i=1}^{N} r_i^2 = \sum_{i=1}^{N} |g(r_i)| \qquad \text{and } \|r_i\|^2 = \left(\sqrt{\sum_{i=1}^{N} r_i^2}\right)^2 = \sum_{i=1}^{n} r_i^2$$

$$\Rightarrow g: r_i \rightarrow r_i^2$$

$$(g: \mathbb{R}^h \rightarrow \mathbb{R}^h)$$

b) (1) $\ell 1$ norm is not applicable to gradient descent as it is not differentiable at $0$. We can use a sparsity (or smoothing) parameter $\epsilon: \sqrt{x^2 + \epsilon}$ to be able to perform the gradient descent. The $\ell 2$ norm is differentiable at each point.

(2) The $\ell 2$ norm penalties discourage sparsity because of the diminishing returns when elements move closer to $0$ but this is what we wish for in using the encoder.

T3 a) $$\frac{\partial E}{\partial v} = \sum_{i=1} \frac{\partial E}{\partial r_i} \cdot \underbrace{\frac{\partial r_i}{\partial v}}_{x_i} = \sum_{i=1} \underbrace{\frac{\partial E}{\partial g(\cdot)}}_{} \cdot \underbrace{\frac{\partial g(\cdot)}{\partial r_i}}_{g'(\cdot)} \cdot x_i = \sum \frac{\partial E}{\partial S_i} \cdot \underbrace{\frac{\partial S_i}{\partial g(\cdot)}}_{\Rightarrow \frac{\partial g(\cdot)}{\partial g(\cdot)} = 1} \cdot g'(\cdot) \cdot x_i$$

$$= \frac{\partial E}{\partial \hat{x}_i} \cdot \underbrace{\frac{\partial \hat{x}_i}{\partial S_i}}_{W} \cdot g'(\cdot) \cdot x_i$$

with $$E = \gamma \|W\|_F^2 + \sum_{i=1}^{n} \|x_i - W g(r_i)\|^2 + \lambda \|r_i\|^2$$

and $$\frac{\partial E}{\partial \hat{x}_i} = x_i - \hat{x}_i \qquad \text{where } \hat{x}_i = W g(\cdot)$$

$$\Rightarrow \frac{\partial E}{\partial v} = \sum_{i=1} (x_i - \hat{x}_i) \cdot W \cdot g'(\cdot) \cdot x_i$$

Exercise 3b
Advantages:
- it is often faster and simpler to obtain sparse representations via autoencoders
- huge reduction in parameters(example: in case of natural images)
- sparsifying non-linearity
- the estimate of the expectation E [$h_j$(x; W, b)] is very noisy in direct optimization but autoencoder denoises the data
- easier to implement
- faster to optimize (no need to keep track of source codes)
- good initial guess for si, optimize from there --> save iterations of source optimization
- can be trained by backpropagation

Disadvantages:
- 2 layers to train
- (more parameters)
- no control of regularization
- bad encoder --> bad decoder