

Exercise Sheet 8

Exercise 1: How CNNs map input to suitable representation for classification (30 P)

Consider a C -class classification problem of a dataset whose samples $\mathbf{x} \in \mathcal{X}$. Suppose we have a CNN f solving the task consisting of the feature extractor module ϕ and the classification module ψ . The feature extractor performs a sequence of convolution (with activation function) and pooling operators, mapping the input \mathbf{x} into a two-dimensional vector $\mathbf{x} \mapsto \phi(\mathbf{x}) \in \mathbb{R}^d$, while the classification module $\psi \in \mathbb{R}^d \rightarrow \mathbb{R}^C$ is an affine transformation transforming the representation into a C -dimensional vector. In other words, the output of the CNN $f : \mathcal{X} \rightarrow \mathbb{R}^C$ is the composition

$$f(\mathbf{x}) = \psi(\phi(\mathbf{x})) \in \mathbb{R}^C,$$

where $f(\mathbf{x})_c = \mathbf{w}_c^\top \phi(\mathbf{x}) + b_c$. Using the softmax function, the probability of Class c given the input \mathbf{x} is therefore

$$p(c|\mathbf{x}) = \frac{\exp(f(\mathbf{x})_c)}{\sum_{c'} \exp(f(\mathbf{x})_{c'})}.$$

(a) Suppose $\forall c, c' \in \{1, \dots, C\} : \|\mathbf{w}_c\| = \|\mathbf{w}_{c'}\| = 1$ and $b_c = b_{c'}$. Denote $\mathbf{z} = \phi(\mathbf{x})$. Show that

$$p(c|\mathbf{x}) = \frac{\exp(\cos(\mathbf{w}_c, \mathbf{z})/\tau)}{\sum_{c'} \exp(\cos(\mathbf{w}_{c'}, \mathbf{z})/\tau)}.$$

where $\tau > 0$ is a quantity of \mathbf{z} .

Solution:

$$\begin{aligned} p(c|\mathbf{x}) &\propto \exp(\mathbf{w}_c^\top \mathbf{z} + b_c) \\ &= \frac{\exp(\mathbf{w}_c^\top \mathbf{z}) \exp(b_c)}{\sum_{c'} \exp(\mathbf{w}_{c'}^\top \mathbf{z}) \exp(b_{c'})} && \text{(by assumption that } b_c = b_{c'}) \\ &\propto \exp(\|\mathbf{w}_c\| \|\mathbf{z}\| \cos(\mathbf{w}_c, \mathbf{z})) && \text{(definition of the Euclidean inner product)} \\ &\propto \exp(\|\mathbf{w}_c\| \|\mathbf{z}\| \cos(\mathbf{w}_c, \mathbf{z})) && \text{(by assumption that } \|\mathbf{w}_c\| = 1) \\ \implies \tau &= \frac{1}{\|\mathbf{z}\|} \end{aligned}$$

(b) Geometrically, what information does the CNN leverage to determine whether the input \mathbf{x} belongs to Class c ? How does the quantity τ play role in prediction? What happens if $\tau \rightarrow \infty$ or $\tau \approx 0$.

Solution:

From (a), the angle is the *semantic* information the CNN uses to assign the label.

The parameter τ is known as the temperature and the norm $\|\mathbf{z}\|$ is thus the inverse of the temperature. The value of τ controls the entropy of the C -class categorical distribution $\text{Cat}(C)$. When the temperature is low ($\tau \rightarrow 0$), the entropy $H[\text{Cat}(C)]$ is small. When the temperature is high, $H[\text{Cat}(C)]$ is large and approaches $\log C$ when $\tau \rightarrow \infty$; that is $\text{Cat}(C)$ becomes more and more uniform.

In other words, $\|\mathbf{z}\|$ encodes uncertainty that the model has on the input \mathbf{x} . If \mathbf{x} has large $\|\mathbf{z}\|$, then τ is low (low entropy), implying that the model is more certain.

In practice, τ is a hyper-parameter and the actual temperature is scaled relative to $\|\mathbf{z}\|$.

Exercise 2: Auto-Encoders and PCA (20 P)

In this exercise, we would like to show an equivalence between linear autoencoders with tied weights (same

parameters for the encoder and decoder) and PCA. We consider the special case of an autoencoder with a single hidden unit. In that case, the autoencoder consists of the two layers:

$$s_i = \mathbf{w}^\top \mathbf{x}_i \quad (\text{encoder})$$

$$\hat{\mathbf{x}}_i = \mathbf{w} \cdot s_i \quad (\text{decoder})$$

where $\mathbf{w} \in \mathbb{R}^d$. We consider a dataset $\mathbf{x}_1, \dots, \mathbf{x}_N \in \mathbb{R}^d$ assumed to be centered (i.e. $\sum_i \mathbf{x}_i = 0$), and we define the objective that we would like to minimize to be the mean square error between the data and the reconstruction:

$$J(\mathbf{w}) = \frac{1}{N} \sum_{i=1}^N \|\mathbf{x}_i - \hat{\mathbf{x}}_i\|^2 \quad (1)$$

Furthermore, to make the objective closer to PCA, we can always rewrite the weight vector as $\mathbf{w} = \alpha \mathbf{u}$ where \mathbf{u} is a unit vector (of norm 1) and α is some positive scalar, and search instead for the optimal parameters \mathbf{u} and α .

(a) Show that the optimization problem can be equally rewritten as

$$\arg \min_{\alpha, \mathbf{u}} J(\mathbf{w}) = \arg \max_{\alpha, \mathbf{u}} \mathbf{u}^\top S \mathbf{u}$$

where $S = \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i \mathbf{x}_i^\top$, which is a common formulation of PCA.

$$\begin{aligned} \arg \min_{\alpha, \mathbf{u}} J(\mathbf{w}) &= \arg \min_{\alpha, \mathbf{u}} \frac{1}{N} \sum_i \|\mathbf{x}_i - \mathbf{w} \mathbf{w}^\top \mathbf{x}_i\|^2 \\ &= \arg \min_{\alpha, \mathbf{u}} \frac{1}{N} \sum_i \|\mathbf{x}_i - \alpha^2 \mathbf{u} \mathbf{u}^\top \mathbf{x}_i\|^2 \\ &= \arg \min_{\alpha, \mathbf{u}} \frac{1}{N} \sum_i (\|\mathbf{x}_i\|^2 - 2\alpha^2 \mathbf{x}_i^\top \mathbf{u} \mathbf{u}^\top \mathbf{x}_i + \alpha^4 \mathbf{x}_i^\top \mathbf{u} \mathbf{u}^\top \mathbf{u} \mathbf{u}^\top \mathbf{x}_i) \\ &= \arg \min_{\alpha, \mathbf{u}} \frac{1}{N} \sum_i (\|\mathbf{x}_i\|^2 - 2\alpha^2 \mathbf{u}^\top \mathbf{x}_i \mathbf{x}_i^\top \mathbf{u} + \alpha^4 \mathbf{u}^\top \mathbf{x}_i \mathbf{x}_i^\top \mathbf{u}) \\ &= \arg \min_{\alpha, \mathbf{u}} \frac{1}{N} \sum_i (-2\alpha^2 \mathbf{u}^\top \mathbf{x}_i \mathbf{x}_i^\top \mathbf{u} + \alpha^4 \mathbf{u}^\top \mathbf{x}_i \mathbf{x}_i^\top \mathbf{u}) \\ &= \arg \min_{\alpha, \mathbf{u}} \frac{1}{N} \sum_i (-2\alpha^2 + \alpha^4) (\mathbf{u}^\top \mathbf{x}_i \mathbf{x}_i^\top \mathbf{u}) \\ &= \arg \min_{\alpha, \mathbf{u}} (-2\alpha^2 + \alpha^4) \mathbf{u}^\top \left(\frac{1}{N} \sum_i \mathbf{x}_i \mathbf{x}_i^\top \right) \mathbf{u} \\ &= \arg \min_{\alpha, \mathbf{u}} (-2\alpha^2 + \alpha^4) \mathbf{u}^\top S \mathbf{u} \\ &= \arg \max_{\mathbf{u}} \mathbf{u}^\top S \mathbf{u} \end{aligned}$$

Exercise 3: Lecture Questions (15 P)

(a) Imagine that we want to learn a representation that is invariant to rotations with an autoencoder. How would you train the autoencoder? What would be the input and what would be the objective function?

We would train an autoencoder by presenting it with rotated images and letting it reconstruct the image without a rotation. This forces the autencoder to have the same representation for the same image which is rotated in two different angles. Given an image \mathbf{x}_i , the input would be randomly rotated images $r(\mathbf{x}_i)$ and the objective function $J(\theta) = \frac{1}{N} \sum_{i=1}^N \|\mathbf{x}_i - f_\theta(r(\mathbf{x}_i))\|^2$

(b) What is the purpose of skip connections in a U-Net?

The skip connections allow the model to pass low-level details to the decoder part directly. Only information that requires more layers (often abstract concepts) are therefore retained in the representation.

(c) Name three different applications of autoencoders.

- **Compression**
- **Anomaly Detection**
- **Denoising**
- **Segmentation**

Exercise 4: Programming (35 P)

Download the programming files on ISIS and follow the instructions.