

# Lecture 11. Kohonen maps, autoencoders, transfer learning, generative adversarial networks

Lecture 11. Kohonen maps, autoencoders, transfer learning, generative adversarial networks

Alex Avdyushenko

Kazakh-British Technical University

November 10, 2022



Good afternoon, dear students. Today we are going to talk about another interesting topic in deep learning, it is autoencoders and generative adversarial networks or simply GANs.

# Lecture 11. Kohonen maps, autoencoders, transfer learning, generative adversarial networks

## └ Five-minutes block

- What are the main disadvantages of the encoder-decoder architecture?
- Write down the attention model formula  $Attn(q, K, V)$
- Describe two criteria for BERT training

As always we start with five-minutes questions on the previous lecture. Please, write answers or send photos with them directly to me in private messages here in Teams or may be in Telegram, but choose only one option please :)

# Lecture 11. Kohonen maps, autoencoders, transfer learning, generative adversarial networks

## └ Formulation of the clustering problem

Given:

$X^T = \{x_1, \dots, x_\ell\}$  — training set of objects  $x_i \in \mathbb{R}^n$

$\rho : X \times X \rightarrow [0, \infty)$  — distance function between objects

Find:

$Y$  is the set of clusters, for example, given by their centers  $w_y \in \mathbb{R}^n$

Let the clustering algorithm is Winner Takes All (WTA):

$$a(x) = \arg \min_{y \in Y} \rho(x, w_y)$$

Optimization criterion: average intracluster distance

$$Q(w; X^T) = \sum_{i=1}^{\ell} \rho^2(x_i, w_{a(x_i)}) \rightarrow \min_{w_y, y \in Y}$$

But we start from afar with the formulation of the clustering problem. So you are given a training set of objects, essentially consisting of number features each.

# Lecture 11. Kohonen maps, autoencoders, transfer learning, generative adversarial networks

## └ Kohonen Network — Two Layer Neural Network



Gradient step in stochastic gradient descent:

$$w_j = w_j + \eta(x_i - w_j)[a(x_i) - y]$$

If  $x_i$  belongs to the cluster  $y$ , then  $w_j$  is shifted towards  $x_i$

T. Kohonen. Self-organized formation of topologically correct feature maps. 1982.

Actually we can write the same operations in the form of special two-layer neural network. And of course you use stochastic gradient descent for optimization of our network.

# Lecture 11. Kohonen maps, autoencoders, transfer learning, generative adversarial networks

## └ Stochastic Gradient Descent

Input: sample  $X^t$ , learning rate  $\eta$ , parameter  $\lambda$

Output: cluster centers  $w_1, \dots, w_c \in \mathbb{R}^n$

□ initialize centers  $w_j, y \in Y$

□ functional estimate:

$$Q = \sum_{i=1}^{\ell} \rho^2(x_i, w_{\phi(x_i)})$$

□ repeat

▪ select object  $x_i$  from  $X^t$  (e.g. random)

▪ compute cluster:  $y = \arg \min_{x \in X^t} \rho(x_i, w_y)$

▪ make a gradient step:  $w_y := w_y + \eta(x_i - w_y)$

▪ evaluate functional:  $Q = \lambda \eta^2[x_i, w_y] + (1 - \lambda)Q$

□ while the value of  $Q$  and/or the weights of  $w_y$  do not converge

Ok, here we have written down all the steps of SGD in this case.

# Lecture 11. Kohonen maps, autoencoders, transfer learning, generative adversarial networks

## └ Hard and soft competition

### WTA, Winner Takes All

$$w_j = w_j + \eta(x_i - w_j)[a(x_i) = y_j], y \in Y$$

#### WTA rule disadvantages

- slow convergence rate
- some cluster centers may never be selected

#### Soft competition rule (WTM, Winner Takes Most)

$$w_j = w_j + \eta(x_i - w_j)K(p(x_i, w_j)), y \in Y$$

where the kernel  $K(p)$  is a nonnegative non-increasing function. Now the centers of all clusters are shifted towards  $x_i$ , but the farther from  $x_i$ , the smaller the shift

Now let's compare two rules of competition. The first one is the Winner Takes All rule, we used it on previous slides. The second one is a little bit softer, the Winner Takes Most rule.

# Lecture 11. Kohonen maps, autoencoders, transfer learning, generative adversarial networks

## └ Kohonen Map (Self Organizing Map, SOM)

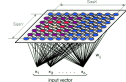
We introduce a rectangular grid of clusters

$\{1, \dots, SimX\} \times \{1, \dots, SimY\}$

Each node  $(x, y)$  is assigned a Kohonen neuron  $w_{xy} \in \mathbb{R}^n$

Along with the metric  $\rho(x, w_{xy})$ , a metric on the grid is introduced:

$$r((x, y), (a, b)) = \sqrt{(x - a)^2 + (y - b)^2}$$



Finally, now we can formulate a new unusual approach to the clustering problem, which is called Kohonen Map (or Self Organizing Map). So we introduce a rectangular grid of clusters, where each node  $(x, y)$  is assigned to a Kohonen neuron  $w_{xy} \in \mathbb{R}^n$ .

# Lecture 11. Kohonen maps, autoencoders, transfer learning, generative adversarial networks

## └ Kohonen Map Training

Input: sample  $X^I$ , learning rate  $\eta$

Output:  $w_{ij} \in \mathbb{R}^D$  are weight vectors

1 initialize weights:  $w_{ij} = \text{random}(-\frac{1}{\sqrt{D}}, \frac{1}{\sqrt{D}})$

2 repeat

• choose a random object  $x_i$  from  $X^I$

• WTA: compute cluster coordinates:

$$(a, b) = \arg \min_{(a,b)} p(x_i, w_{ab})$$

• for all  $(a, b) \in \text{Neighborhood}(a, b)$

WTA: do gradient descent step:

$$w_{ab} := w_{ab} + \eta(x_i - w_{ab})K(r((a, b), (a, b)))$$

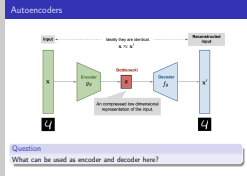
3 yet clustering does not stabilize

How do we train it?



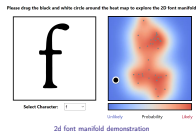
# Lecture 11. Kohonen maps, autoencoders, transfer learning, generative adversarial networks

## └ Autoencoders



Ok, great. Now let's move on to another interesting approach in machine learning. It is autoencoders and it's main idea is quite simple. So you have some input vector, for example an image. Then you have two parts of the model, which are the encoder part, where you get some representation of the input, and the decoder part, where you try to get back to the initial input vector.

# Lecture 11. Kohonen maps, autoencoders, transfer learning, generative adversarial networks



Here let's enjoy a little demo where different fonts are nested in a 2D plane.

# Lecture 11. Kohonen maps, autoencoders, transfer learning, generative adversarial networks

## └ Denoising Auto Encoder

Stability of code vectors  $z_i$  with respect to noise in  $x_i$ :

$$\mathcal{L}_{DAE}(z, \beta) = \sum_{i=1}^I \mathbb{E}_{\tilde{x} \sim q(\tilde{x}|x_i)} \mathcal{L}(g(f(\tilde{x}, \alpha), \beta), x_i) \rightarrow \min_{z, \beta}$$

Instead of calculating the expectation  $\mathbb{E}_z$  in the stochastic gradient method,  $x_i$  objects are sampled and noisy one at a time:  $\tilde{x} \sim q(\tilde{x}|x_i)$

- Gaussian noise:  $\tilde{x} \sim N(x, \sigma^2 I)$
- zeroing components of vector  $x_i$  with probability  $p_0$



P. Vincent, H. Larochelle, Y. Bengio, P.-A. Manzagol. Extracting and composing robust features with denoising autoencoders. ICML-2008.

Ok, now let's discuss briefly a few more autoencoders. For a example, the denoising one.

# Lecture 11. Kohonen maps, autoencoders, transfer learning, generative adversarial networks

## └ Variational Auto Encoder

A generative model is constructed capable of generating new objects  $x$

similar to the objects of the sample  $X^t = \{x_1, \dots, x_t\}$

$q_\alpha(x|x)$  – probabilistic encoder with  $\alpha$  parameter

$p_\beta(\tilde{x}|z)$  – probabilistic decoder with  $\beta$  parameter

$$\begin{aligned}\mathcal{L}_{VAE}(\alpha, \beta) &= \sum_{i=1}^t \log p(x_i) = \sum_{i=1}^t \log \int q_\alpha(x_i|x) \frac{p_\beta(x_i|x) p(x)}{q_\alpha(x_i|x)} dx \geq \\ &\geq \sum_{i=1}^t \int q_\alpha(x_i|x_i) \log p_\beta(x_i|x) dx - KL(q_\alpha(x_i|x_i) || p(x)) \rightarrow \max_{\alpha, \beta}\end{aligned}$$

D.P.Kingma, M.Welling. Auto-encoding Variational Bayes. 2013.

C.Doruch. Tutorial on variational autoencoders. 2016.

The next one is Variational Auto Encoder and it may be difficult to understand from the first view, but it's main feature is to use lower variational bound of log likelihood of your model.

# Lecture 11. Kohonen maps, autoencoders, transfer learning, generative adversarial networks

$$\sum_{i=1}^I \underbrace{E_{x \sim q_i(x)} \log p_\theta(x|x_i)}_{\text{quality reconstruction}} - \underbrace{KL(q_i(x|x_i) \| p(x))}_{\text{regularizer by } \alpha} \rightarrow \max_{\alpha, \theta}$$

where  $p(x)$  is the prior distribution, usually  $N(0, \sigma^2 I)$

Reparametrization  $q_i(x|x_i) : z = f(x_i, \alpha, \epsilon), \epsilon \sim N(0, I)$

Stochastic gradient method:

- sample  $x_i \sim \mathcal{X}^i, \epsilon \sim N(0, I), z = f(x_i, \alpha, \epsilon)$
- gradient step  $\alpha = \alpha + \eta \nabla_{\alpha} [\log p_\theta(x_i | f(x_i, \alpha, \epsilon)) - KL(q_i(x|x_i) \| p(x))]$   
 $\beta = \beta + \eta \nabla_{\beta} [\log p_\theta(x_i | z)]$

Generation of similar objects:

$$x \sim p_\theta(x | f(x_i, \alpha, \epsilon)), \epsilon \sim N(0, I)$$

So we two parts of functional, and first one if responsible for quality reconstruction, second one is actually regularizer by  $\alpha$ .

To evaluate the expactation of some function with z distribution you need to use reparametrization trick. What is it? On the first step you sample z-value as value of new function f with new parameter  $\epsilon$ , which comes from some random distribution, for example normal. And then on the second step you get the expectation with your z-value.

# Lecture 11. Kohonen maps, autoencoders, transfer learning, generative adversarial networks

## └ Multi-task learning

### Multi-task learning

- ♦  $f(x, \alpha)$  – universal part of the model (vectorization)
  - ♦  $g(x, \beta)$  – specific parts of the model for problems  $t \in T$
- Simultaneous training of the model  $f$  on tasks  $X_t, t \in T$ :

$$\sum_{t \in T} \sum_{x \in X_t} \mathcal{L}_t(f(x_t, \alpha), g(x_t, \beta_t)) \rightarrow \min_{\alpha, \{\beta_t\}}$$

*Learnability*: the quality of solving a particular problem  $(X_t, \mathcal{F}_t, g_t)$  improves with increasing sample size  $l_t = |X_t|$ .

*Learning to learn*: the quality of the solution of each of the problems  $t \in T$  improves with the growth of  $l_t$  and the total number of problems  $|T|$ .

*Few-shot learning*: a small number of examples, sometimes even one, is enough to solve the  $t$  problem.

M. Cranenban. Multi-task learning with deep neural networks: a survey. 2020  
Y. Wang et al. Generalizing from a few examples: a survey on few-shot learning. 2020

The hot topic of the last couple of years is exactly multitask and multi-domain training. So we want to have one AI model, which is good for all tasks => Something like another one step on the way to general AI.

## Temporary page!

$\text{\LaTeX}$  was unable to guess the total number of pages correctly. A was some unprocessed data that should have been added to the this extra page has been added to receive it.

If you rerun the document (without altering it) this surplus page away, because  $\text{\LaTeX}$  now knows how many pages to expect for t document.