

VIR 2022

Name: Chytrý Student

Exam test

Variant: A

Points 20

1. Let us consider gradient learning of the linear regressor  $y = \mathbf{w}^\top \mathbf{x}$ . Given the single training example  $(\mathbf{x} = [\sqrt{3}, 1]^\top, y = 0)$ , the least squares learning reduces to the minimization of the following criterion

$$f(\mathbf{x}, \mathbf{w}) = \frac{1}{2} \|\mathbf{w}^\top \mathbf{x}\|_2^2 = \frac{1}{2} ((w_1 x_1)^2 + (w_2 x_2)^2)$$

$$\left. \frac{\partial f}{\partial w_1} \right|_{x_1} = w_1 x_1^2 \Big|_{x_1} = 3 w_1$$

$$\left. \frac{\partial f}{\partial w_2} \right|_{x_2} = w_2 x_2^2 \Big|_{x_2} = w_2$$

**TASK 1.1** Derive the recurrent formula for values of weights in the  $k$ -th iteration

$$w_1^k = \rho_1(\alpha)^k w_1^0 = (1 - 3\alpha)^k w_1^0$$

$$w_2^k = \rho_2(\alpha)^k w_2^0 = (1 - \alpha)^k w_2^0$$

$$w_1^k = w_1^{k-1} - \alpha \cdot 3 w_1^{k-1}$$

$$w_2^k = w_2^{k-1} - \alpha \cdot w_2^{k-1}$$

**TASK 1.2** For which learning rate  $\alpha$  the gradient descent converges (at least slowly) in both dimensions?

**Hint:** The smaller the  $|\rho_i(\alpha)|$ , the faster the convergence. Find  $\alpha$  for which both formulas converge to zero.

$$\alpha^{\text{convergent}} \in (0, 2/3)$$

$$|1 - 3\alpha| < 1 \quad \wedge \quad |1 - \alpha| < 1$$

**TASK 1.3** What is the best learning rate  $\alpha^*$ , which guarantees the fastest convergence rate for arbitrary weight initialization  $\mathbf{w}^0$  and this particular training example.

**Hint:** Choose alpha, which minimizes the maximum of both convergence rates:

$$\alpha^* = \arg \min_{\alpha} \max\{|\rho_1(\alpha)|, |\rho_2(\alpha)|\} = \frac{1}{2}$$

$$1 - \alpha^* = 3\alpha^* - 1$$

$$\alpha^* = \frac{1}{2}$$

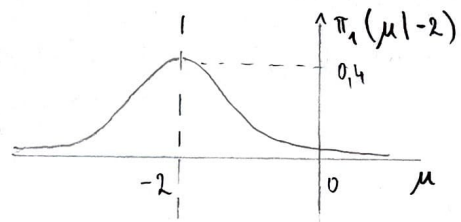
2. Consider stochastic continuous policy, that selects the action  $\mathbf{u} \in \mathbb{R}$  in the state  $\mathbf{x} \in \mathbb{R}$  according to the following probability distribution:

$$\pi_{\theta}(\mathbf{u}|\mathbf{x}) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}(\theta\mathbf{x} - \mathbf{u})^2\right)$$

with scalar parameter  $\theta = 1$ . This policy maps one-dimensional state  $\mathbf{x}$  on the Gaussian probability distribution (with the unit variance) of possible actions  $\mathbf{u}$ .

**TASK 2.1** Let us assume that the robot/agent is in state  $\mathbf{x}_1 = -2$ . Sketch the shape of probability distribution  $\pi_{\theta}(\mathbf{u}|\mathbf{x}_1 = -2)$  from which the actions are drawn.

$$\pi_{\theta}(\mu|-2) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}(\mu + 2)^2\right) \rightarrow \mathcal{N}(-2, 1)$$



**TASK 2.2** The policy performs the action  $\mathbf{u}_1 = 1$  (that has been randomly generated from the probability distribution), and the robot ends up in the state  $\mathbf{x}_2 = +3$ . The reward function for the resulting training trajectory  $\tau = [\mathbf{x}_1, \mathbf{u}_1, \mathbf{x}_2]$  is  $r(\tau) = 2$ . Estimate REINFORCE policy gradient:

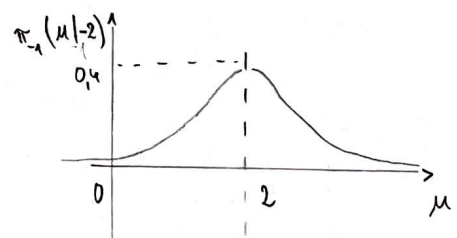
$$\frac{\partial \log \pi_{\theta}(\mathbf{u}|\mathbf{x})}{\partial \theta} \Bigg|_{\substack{\mathbf{x} = \mathbf{x}_1 \\ \mathbf{u} = \mathbf{u}_1}} \cdot r(\tau) = \frac{\partial}{\partial \theta} \log \left[ \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}(\theta\mathbf{x} - \mu)^2\right) \right] \Bigg|_{\substack{\mathbf{x}_1 \\ \mu_1}} \cdot r(\tau) =$$

$$= -(\theta\mathbf{x} - \mu) \times \Bigg|_{\substack{\mathbf{x}_1 \\ \mu_1}} \cdot r(\tau) = -12$$

**TASK 2.3** Update policy parameters by the gradient ascent method with  $\alpha = 1/6$  and sketch the shape of the updated distribution  $\pi_{\theta^{\text{updated}}}(\mathbf{u}|\mathbf{x}_1 = -2)$

$$\theta^{\text{updated}} = \theta + \alpha(-12) = -1$$

$$\pi_{\theta}(\mu|-2) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}(\mu - 2)^2\right) \rightarrow \mathcal{N}(2, 1)$$



3. You are given an input feature map (image)  $x$ , a convolution layer  $\text{Conv2d}(\text{in\_channels}=3, \text{out\_channels}=6, \text{kernel\_size}=5, \text{stride}=1, \text{padding}=0, \text{dilation}=1)$ , an activation function  $\text{ReLU}$ , a batch normalization layer  $\text{BatchNorm2d}(6)$ , a max pooling layer  $\text{MaxPool2d}(2, 2)$  and an output  $y$ .

**TASK 3.1:** Consider the following architecture / nemaji vliv na RF

$$x \rightarrow \text{Conv2d} \rightarrow \text{ReLU} \rightarrow \text{BatchNorm2d} \rightarrow \text{MaxPool2d} \rightarrow y$$

and compute the receptive field (RF) of the output, i.e., the size of the region in the input  $x$  that produces the feature  $y_{i,i}$ :

Conv 2D: zajimá nás kernel\_size a dilation => 5x5

$k_s = 5$   $d = 1$

Max Pool 2D: + 1 k RF (2x2)

$$\text{RF} = 6 \times 6$$

**TASK 3.2:** Tick the correct answer (multiple choice).

- A receptive field depends on the size of the input image.
- A batch normalization procedure consists of feature-wise operations which do not alter the receptive field of the network.
- Some linear layers increase the size of the receptive field.
- The larger the convolutional stride, the larger the receptive field.
- By adding more convolutional layers, an arbitrarily large receptive field can be achieved.
- A large receptive field usually negatively impacts the ability of the neural network to understand the context of the input image.

# Nebude v testu

4. Consider the composite normalizing flow  $f : \mathbb{R}^3 \rightarrow \mathbb{R}^3$ ,  $f = f_1 \circ f_2$  of length 2

$$P_Z \sim z \begin{array}{c} \xrightarrow{g_1} \\ \xleftarrow{f_1} \end{array} \mathbf{y} \begin{array}{c} \xrightarrow{g_2} \\ \xleftarrow{f_2} \end{array} \mathbf{x} \sim P_X$$

$Z \sim U([-1, 1]^3)$ , i.e.,  $Z$  is a real random vector in  $\mathbb{R}^3$  with uniform distribution over the cube of edge length 2.

**TASK 4.1:**  $g_1$  is specified as a linear transformation

$$g_1 : \mathbf{y} = A\mathbf{z} + \mathbf{b},$$

where  $A$  is a  $3 \times 3$  square matrix and  $\mathbf{b}$  is a  $3 \times 1$  column vector

$$A = \begin{bmatrix} 1 & 0 & 0 \\ 3 & 2 & 1 \\ 1 & 2 & 3 \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} 0 \\ 1 \\ 1 \end{bmatrix}.$$

We have  $f_1 = g_1^{-1}$ . Calculate the determinant of Jacobian of  $f_1$ , i.e., calculate  $\det(\mathbf{J}_{f_1}) = \det(\mathbf{J}_{g_1^{-1}})$ . Note that you do not need to know the inverse matrix  $A^{-1}$  to complete this task.

**TASK 4.2:**  $f_2$  is a simple coupling flow  $\mathbb{R}^3 \rightarrow \mathbb{R}^3$  that is specified as follows,  $\mathbf{y} = f_2(\mathbf{x})$ :

$$\begin{aligned} y_1 &= x_1, \\ y_2 &= x_2 \cdot \exp(+2x_1) + x_1, \\ y_3 &= x_3 \cdot \exp(-2x_1) + x_1. \end{aligned}$$

Calculate the determinant of the Jacobian  $f_2$ .

# Nebude v testu

**TASK 4.3:** Consider the real data point  $\mathbf{x}^* = (0, 1, 1)^T$ . Assume that  $\mathbf{x}^*$  was generated from distribution  $P_X$  which is further normalized by the flow transformation  $f$  to the distribution  $P_Z \sim U([-1, 1]^3)$ .

Calculate the latent representation  $\mathbf{z}^*$  of  $\mathbf{x}^*$  under  $f$  (Hint:  $A^{-1}$  is not required).

$$\mathbf{z}^* = f(\mathbf{x}^*) = f_1 \circ f_2(\mathbf{x}^*) =$$

**TASK 4.4:** What is the value of density  $p_X$  at this point? Use the change of variable formula

$$p_X(\mathbf{x}) = p_Z(f(\mathbf{x})) \cdot |\det(\mathbf{J}_f)|$$

and results from TASK 4.1, 4.2, 4.3 to complete the task.

$$p_X(\mathbf{x}) =$$