To be completed individually or in groups of two people. **Please be sure matriculation numbers are clearly included at the top of your submission.** Submissions can be handwritten or in LaTeX formatting, but hard-to-read handwritten submissions will not be graded.

Please submit via Ilias. Submissions should be a single PDF document (note that Jupyter notebooks can and should also be downloaded as PDFs, and not submitted as .ipynb files).

Each question will be graded "pass" (full points) or "fail" (no points). We award 0.5 bonus points for the exam for each theory and practical question solved. You must complete 50% of all exercises to enter the final exam.

1. **EXAMple Question** *(Affine transformations)*

   Consider an $m$-dimensional normally-distributed random vector $\mathbf{u} \sim \mathcal{N}(\mu, \boldsymbol{\Sigma})$. Let $\mathbf{x} = f(\mathbf{u}|\mathbf{A}, \mathbf{b}) := \mathbf{A}\mathbf{u} + \mathbf{b}$ be an affine transformation, where $\mathbf{A} \in \mathbb{R}^{n \times m}, \mathbf{b} \in \mathbb{R}^n$.

   (a) Using the formula for the transformation of random variables from the lecture, write the conditions for $f(\mathbf{u}|\mathbf{A}, \mathbf{b})$ to be a bijection. What is $f^{-1}(\mathbf{x})$?

   (b) Using $f^{-1}(\mathbf{x})$ from the previous part, derive $p_{\mathbf{x}}(\mathbf{x})$ in terms of $\mathbf{A}, \mathbf{b}, \mu, \boldsymbol{\Sigma}$.

   (c) For $m = n$, find a parameterization of $\mathbf{A}, \mathbf{b}$ such that $f(\mathbf{u}|\mathbf{A}, \mathbf{b})$ is always invertible. (*Hint*: enforce invertibility of $\mathbf{A}$ using Cholesky factorization or other methods.)

2. **Theory Question** *(Density estimation)* You have familiarized yourselves with the **Kullback-Liebler (KL) divergence** in the previous tutorial. This question deals with certain properties of the KL divergence, and its use in density estimation using normalizing flows. The KL divergence, which measures the discrepancy between two distributions $p(x)$ and $q(x)$, is defined as

   $$D_{\mathrm{KL}}\left[p(x)\|q(x)\right] = \mathbb{E}_{p(x)}\left[\log \frac{p(x)}{q(x)}\right].$$

   (a) Is $D_{\mathrm{KL}}[p(x)\|q(x)] = D_{\mathrm{KL}}[q(x)\|p(x)] \ \forall p, q$? If not, provide a counterexample.

   (b) You are given an i.i.d. dataset $\{x_n\}_{n=1}^N$. $x_n$'s are samples from the true data distribution $p_{\mathrm{true}}(x)$. Given an approximating data distribution $p(x|\theta)$, the goal is to estimate the maximum likelihood estimate (MLE) of $\theta$ over the dataset.
   Prove that as the dataset size $N \to \infty$,

   $$\theta_{\mathrm{MLE}} = \underset{\theta}{\mathrm{argmax}} \ p(\{x_n\}_{n=1}^N|\theta) \to \underset{\theta}{\mathrm{argmin}} \ D_{\mathrm{KL}}[p_{\mathrm{true}}(x)\|p(x|\theta)].$$

   In other words, for large number of datapoints, the MLE estimate converges to the parameter that minimizes the KL divergence between the true data distribution and the approximating data distribution. If needed, use the *uniform law of large numbers*[1].

   (c) Let $x_n \in \mathbb{R}^d$. Our goal is to approximate the density $p_{\mathrm{true}}(x)$ by $p(x|\theta)$, such that we can generate new datapoints $x' \sim p(x|\theta)$. This can be done using a normalizing flow, which transforms a random vector which can be sampled *easily* into data from the *desired distribution*.
   Consider a general bijective transformation $x = f(u|\theta) : \mathbb{R}^d \to \mathbb{R}^d$. Assuming $u \sim \mathcal{N}(0, I)$, density estimation over $x$ is performed by learning the parameters of the transformation $f$. The distribution over $x$ induced by $f$ can be learnt by minimizing $D_{\mathrm{KL}}[p_{\mathrm{true}}(x)\|p(x|\theta)]$. Using the result of the previous part, write a loss function for finding optimal parameters of the normalizing flow.

---

[1] Law of large numbers Section 4.4 – Wikipedia, the Free Encyclopedia.

(d) (*Optional*) $D_{\mathrm{KL}}[p(x|\theta)\|p_{\mathrm{true}}(x)]$ is usually referred to as the reverse KL divergence, as opposed to the previously mentioned $D_{\mathrm{KL}}[p_{\mathrm{true}}(x)\|p(x|\theta)]$, which is called the forward KL divergence. Minimizing both of these takes the approximating distribution closer to the true distribution. In the context of learning normalizing flows, how should you choose between the two objectives? (*Hint*: note that the true distribution may be either hard to sample from, or hard to evaluate.)

3. **Practical Question** In this exercise you will train normalizing flows for density estimation using `Pyro`, a probabilistic programming framework based on `Pytorch`. See `Exercise_11.ipynb`.