



List of topics to cover

With section titles and brief explanations.

Yiftach Kolb

Berlin, October 6, 2022

Freie Universität



Berlin

Abstract

punkt. punkt.

Declaration

punkt. punkt.

Acknowledgement

punkt.[3] punkt.

List of Tables

List of Figures

3.1	VAE graphical model	16
4.1	a figure	18

Contents

1	Introduction	8
2	Notations and definitions, preliminary concepts	9
2.1	Basic notations	9
2.2	The data	9
2.3	Linear algebra preliminary: SVD and PCA	10
2.4	Neural networks	11
2.5	Autoencoders	12
2.5.1	Relation between PCA and AE	12
3	Variance inference and variational autoencoders	13
3.1	Variational Inference	13
3.2	Variational Autoencoder	15
3.2.1	Adding parameters	15
3.2.2	Using neural networks for the parametrization	15
3.2.3	Graphical representation	15
4	Gaussian mixture model VAEs	17
4.0.1	Relation between AE and VAE	17
4.0.2	Conditional VAE	17
5	Experiments and results	19
5.1	Tests with MNIST and FMNIST	19
5.2	Tests with scRNAseq Data	19

6	Discussion, some remarks and conclusions	20
----------	---	-----------

Chapter 1

Introduction

punkt. *punkt*, *punkt*.

Chapter 2

Notations and definitions, preliminary concepts

2.1 Basic notations

Throughout this paper (modulo typing errors) we use capital bold math Latin or Greek letters ($\mathbf{X}, \mathbf{\Sigma}$) to represent matrices, bold small math letters (\mathbf{x}) represent (usually row) vectors, but in cases where it makes sense may also represent a matrices such as a batch of several vectors (each row is a different data point). Or another example, $\boldsymbol{\sigma}$ may represent both the covariance matrix and the variance vector of a diagonal Gaussian distribution. Non-bold math letters (x, σ, \dots) may represent scalar or vectors in some cases and hopefully it is clear from the context or explicitly stated.

2.2 The data

we are assuming that the input data unless otherwise stated is real and 2-dimensional. Its rows represent *samples*, for example—cells in the case of scRNAseq dataset, while its columns represent *variables*, for example—genes in scRNAseq dataset.

In addition to \mathbf{X} there may be additional data with information about class or conditions. We use *one-hot encoding* to represent such information.

Definition 2.1. A *data matrix* is simply a real valued matrix $\mathbf{X} \in \mathbb{R}^{N \times n}$ which represent a set of N n -dimensional data points. The N rows are also called *observations* and the n columns are *variables*.

Definition 2.2. A *class matrix*, or also *condition matrix* $\mathbf{C} \in \mathbb{R}^{N \times c}$ is simply a real matrix which represents one-hot encoding of c classes or conditions over N samples. For example if sample i has class j , then $(\forall k \in 1, \dots, c) \mathbf{C}[i, k] = \delta_{jk}$.

We say that that \mathbf{C} is a *class probability matrix* or *relaxed class matrix* (same with condition) if instead of being one-hot it is a distribution matrix, namely each row is non-negative and sums up to 1.

Usually if the input data includes class/condition information, it comes as a class matrix

(pure one-hot) but the output (the prediction) is naturally probabilistic and hence is relaxed.

2.3 Linear algebra preliminary: SVD and PCA

Let $\mathbf{X} \in \mathbb{R}^{N \times n}$ be a real-valued matrix representing N samples of some n -dimensional data points and let $r = \text{rank}(\mathbf{X}) \leq \min(n, N)$.

$\mathbf{X}\mathbf{X}^T$ and $\mathbf{X}^T\mathbf{X}$ are both symmetric and positive semi-definite. Their eigenvalues are non-negative, and they both have the same positive eigenvalues, exactly r such, which we mark $s_1^2 \geq s_2^2 \geq \dots s_r^2 > 0$. The values $s_1 \dots s_r$ are called the *singular values* of \mathbf{X} .

$$\text{Let } \mathbf{S} = \begin{pmatrix} s_1 & & & \\ & s_2 & & \\ & & \ddots & \\ & & & s_r \end{pmatrix} \in \mathbb{R}^{r \times r}$$

Let $\mathbf{U} \in \mathbb{R}^{N \times N}$ be the (column) eigenvectors of $\mathbf{X}\mathbf{X}^T$ sorted by their eigenvalues. Then $\mathbf{U} = (\mathbf{U}_r, \mathbf{U}_k)$ where $\mathbf{U}_r \in \mathbb{R}^{N \times r}$ are the first r eigenvectors corresponding to the non-zero eigenvalues, and \mathbf{U}_k are the eigenvectors corresponding to the $N - r$ 0-eigenvalues. Similarly, and let $\mathbf{V} = (\mathbf{V}_r, \mathbf{V}_k) \in \mathbb{R}^{n \times n}$ be the (column) eigenvectors of $\mathbf{X}^T\mathbf{X}$, sorted by the eigenvalues, where $\mathbf{V}_r \in \mathbb{R}^{n \times r}$ are the first r eigenvectors and \mathbf{V}_k are the $n - r$ null-eigenvectors.

Then the *singular value decomposition (SVD)* of \mathbf{X} is

$$\mathbf{X} = \mathbf{U}\mathbf{D}\mathbf{V}^T \quad (2.1)$$

where $\mathbf{D} = \left(\begin{array}{c|c} \mathbf{S} & \mathbf{0} \\ \hline \mathbf{0} & \mathbf{0} \end{array} \right) \in \mathbb{R}^{N \times n}$ is diagonal.

\mathbf{V}_r are called the (*right*) *principle components* of \mathbf{X} . Note that $\mathbf{V}_r^T \mathbf{V}_r = \mathbf{I}_r$ and that $\mathbf{X} = \mathbf{X}\mathbf{V}_r\mathbf{V}_r^T$. If one looks at the second expression, it means that the each row of \mathbf{X} is spanned by the orthogonal basis \mathbf{V}_r (because the other vectors of \mathbf{V} are in $\ker(\mathbf{X})$).

More generally For every $l \leq r$, let $\mathbf{V}_l \in \mathbb{R}^{n \times l}$ be the first l components, Then $\mathbf{X}\mathbf{V}_l\mathbf{V}_l^T$ is as close as we can get to \mathbf{X} within an l -dimensional subspace of \mathbb{R}^n , and \mathbf{V}_l minimizes

$$\mathbf{V}_l = \text{argmin}\{\|\mathbf{X} - \mathbf{X}\mathbf{W}\mathbf{W}^T\|_F^2 : \mathbf{W} \in \mathbb{R}^{n \times l}, \mathbf{W}^T\mathbf{W} = \mathbf{I}_l\} \quad (2.2)$$

Where $\|\cdot\|_F^2$ is simply the sum of squares of the matrix' entries.

If we consider the more general minimization problems:

$$\min\{\|\mathbf{X} - \mathbf{X}\mathbf{E}\mathbf{D}\|_F^2 : \mathbf{E}, \mathbf{D}^T \in \mathbb{R}^{n \times l},\} \quad (2.3)$$

$$\min\{\|\mathbf{X} - \mathbf{X}\mathbf{W}\mathbf{W}^\dagger\|_F^2 : \mathbf{W} \in \mathbb{R}^{n \times l},\} \quad (2.4)$$

It can be shown [5] that the last two problems 2.3, 2.4 are equivalent and that for any solution E, D it must hold that $D = E^\dagger$. (D is the Moore-Penrose generalized inverse of E). Moreover, V_l still minimizes the general problem 2.3 and for every solution W , it must hold that $\text{span}\{W\} = \text{span}\{V_l\}$ (but it isn't necessarily an orthogonal matrix).

2.4 Neural networks

Definition 2.3. A *feed forward neural network* is simply a parameterized differentiable map $\phi_w : \mathbb{R}^n \rightarrow \mathbb{R}^m$. $\phi_w(x)$ is differentiable both in its input variable x as well as in its parameters w , which is also called its *weights*.

For example each affine function has the form $f : x \rightarrow a \cdot x + b$. a and b are its trainable parameters (weights). The input is treated as fixed (the data we are trying to explain).

Normally ϕ is a sequence of compositions of more simple functions. We call such more basic unit in a composition sequence a *layer*. Each layer is itself a composition, with exactly a single affine map, followed by 0 or more dimension preserving functions such as normalization functions or activation functions. An *activation function* is a real values non-linear function which is applied element-wise over the input later. For example the sigmoid function and the ReLU (rectified linear unit) are well-known and often used activation functions.

Definition 2.4. Usually together with a neural network comes an associated differentiable function which is the *loss function* $\mathcal{L} : \mathbb{R}^m \rightarrow \mathbb{R}$.

Typically the loss function is additive on the dimension, meaning it has the form $\mathcal{L}(x) = \sum_{i=1}^m \psi(x_i)$

An example for such a loss function is the square error $x \rightarrow \frac{1}{2} \|x\|_2^2$

Definition 2.5. Let $X \in \mathbb{R}^{N \times n}$ be a data matrix. A *batch* $x \in \mathbb{R}^{b \times n}$ is any subset of b rows of X .

Batch $x \in \mathbb{R}^{b \times n}$ represents a subset of b samples out of the total of N samples in the dataset. The operations of ϕ and \mathcal{L} *extend naturally* to batches (note that in this case x is a matrix and x_i is a vector representing a single sample of the batch)—we collect for ϕ and we average for \mathcal{L} , namely: $\phi(x) = (\phi(x_i))_{i=1}^b \in \mathbb{R}^{b \times m}$, and $\mathcal{L}(x) = \frac{1}{b} \sum_{i=1}^b \mathcal{L}(x_i)$.

If \mathcal{L} is the square error function $\|\cdot\|_2^2$ on vectors, then its expansion to batches is $\frac{1}{b} \|\cdot\|_F^2$. The reason why we sum and don't average over the dimensions will be cleared later when we get into variational inference.

Training the neural network ϕ_w means finding the weights that minimize the loss function applied on the training set X , in other words minimizing $\min_w (\mathcal{L}(\phi_w(X)))$.

During a training step the network is applied on a batches x . Then the loss function is applied at the output and a gradient (with relation to the weights) is taken. This gradient is used for the weight update rule, which varies depending on the specific training algorithm. Typical training algorithms are SGD (stochastic gradient decent) and Adam, which is the one used throughout this work.

We only need to define the network, the loss function and the specific training algorithm.

The rest (derivation, weight update etc.) is taken care for us by the backend of the software (Pytorch) and can be regarded as a black box.

2.5 Autoencoders

Definition 2.6. An *Autoencoder* (AE) is a neural network $\phi = D \circ E$ with a "bottleneck" layer and which approximates the identity function on the training input.

We call E which projects into the bottleneck, the *encoder*, and D which expands back into the input dimensions, the *decoder*.

2.5.1 Relation between PCA and AE

For **centered** data, meaning every variable (column of \mathbf{X}) has 0 sample mean, the first $k \leq \text{rank}(\mathbf{X})$ principle components \mathbf{P} are the solution for equation 2.2; Whereas a **linear** autoencoder solves equation 2.4. As mentioned, it must hold that $E = D^\dagger$ (the encoder must be the Moore-Penrose inverse of the decoder).

A linear autoencoder (an AE where ϕ is linear) is therefore almost equivalent to PCA [5], in that in the optimum, a bottleneck space of dimension k is spanned by the first k principle components of the input \mathbf{X} . In general, an AE can be seen a PCA-like, but non-linear method for dimensionality reduction.

Chapter 3

Variance inference and variational autoencoders

3.1 Variational Inference

Here we briefly explain the idea behind variational inference and introduce the ELBO which is the loss function we'll use throughout this text. For more details see Christopher M Bishop and Nasser M Nasrabadi. *Pattern recognition and machine learning*. Vol. 4. 4. Springer, 2006.

We treat the data matrix as a set of independent observations (its rows) $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ which we try to explain by a probabilistic model. We assume that the \mathbf{x}_i 's are i.i.d with some distribution function $p(\mathbf{x})$ and therefore for the entire dataset it holds that $p(\mathbf{X}) = \prod p(\mathbf{x}_i)$.

Definition 3.1. Let $\mathbf{X} \in \mathbb{R}^{Nn}$ be a data matrix and let $\{\mathbf{x}_i\}_1^n$ be its rows, which we assume to be i.i.d with some (unknown) distribution $p(\mathbf{x})$. Then $\log p(\mathbf{X}) = \sum_1^N \log p(\mathbf{x}_i)$ is called the *log evidence* of our data.

The r.vs \mathbf{X} are high dimensional however we have some reason to believe that behind the scenes there are some hidden (latent), smaller dimensional, r.vs $\mathbf{Z} = \{\mathbf{z}_1 \dots \mathbf{z}_N\}$ that generate the observations \mathbf{X} . In other words we think that \mathbf{X} is conditioned on \mathbf{Z} and we can speak of the joint distribution $p(\mathbf{X}, \mathbf{Z}) = p(\mathbf{X}|\mathbf{Z})p(\mathbf{Z})$. Because we assume i.i.d for both \mathbf{X} and \mathbf{Z} all the distributions factor over the individual samples multiplicatively, e.g. $p(\mathbf{X}|\mathbf{Z}) = \prod p(\mathbf{x}_i|\mathbf{z}_i)$.

Suppose that we have a fully Bayesian model. In this case there are no parameters because the parameters are themselves stochastic variables with some suitable priors. We can therefore pack all the latent variables and stochastic parameters into one latent "meta variable" $\mathbf{Z} = (\mathbf{z}_1, \mathbf{z}_2, \dots)$, where each \mathbf{z}_i is some multidimensional r.v and possibly composed of several simpler r.vs (for example a categorical and a normal r.vs). We similarly pack all the observed variables into one meta variable \mathbf{X} . Together we have a distribution $p(\mathbf{X}, \mathbf{Z})$ and the working assumption is that it is easy to factorize $p(\mathbf{X}, \mathbf{Z}) = p(\mathbf{X}|\mathbf{Z})p(\mathbf{Z})$, however $p(\mathbf{Z}|\mathbf{X})$ is intractable and $p(\mathbf{X})$ is unknown.

We are being Bayesian here so we consider $\mathbf{X} = (\mathbf{x}_1, \mathbf{x}_2, \dots)$ to be a constant a set of observations and we want to best explain $p(\mathbf{X})$ by finding as high as possible lower

bound for it (or rather to $\log p(\mathbf{X})$, the *log evidence*). A second goal is to approximate the intractable $p(\mathbf{Z}|\mathbf{X})$ by some simpler distribution $q(\mathbf{Z})$ taken from some family of distributions.

Definition 3.2. Let \mathbf{x}, \mathbf{z} be random variables with joint distribution $p(\mathbf{x}, \mathbf{z})$ and let $q(\mathbf{z})$ be any distribution. The *evidence lower bound (ELBO)* with respect to p, q is:

$$\mathcal{L}(q) \triangleq \int \log \frac{p(\mathbf{x}, \mathbf{z})}{q(\mathbf{z})} dq(\mathbf{z}) \quad (3.1)$$

The following equation shows that the *ELBO* is a lower bound for the *log evidence*. (using Jensen's inequality)

$$\begin{aligned} \log p(\mathbf{X}) &= \log \int p(\mathbf{X}, \mathbf{Z}) d\mathbf{Z} = \log \int \frac{p(\mathbf{X}, \mathbf{Z})}{q(\mathbf{Z})} q(\mathbf{Z}) d\mathbf{Z} \\ &= \log \int \frac{p(\mathbf{X}, \mathbf{Z})}{q(\mathbf{Z})} dq(\mathbf{Z}) \geq \int \log \frac{p(\mathbf{X}, \mathbf{Z})}{q(\mathbf{Z})} dq(\mathbf{Z}) \triangleq \mathcal{L}(q) \end{aligned} \quad (3.2)$$

In equation 3.2 we found a lower bound $\mathcal{L}(q)$ for the log evidence $\log p(\mathbf{X})$, the *ELBO*. Whatever distribution q we put in ELBO will not be greater than the real log evidence so we are looking for the q which **maximizes** it.

Now we show that maximizing the ELBO actually obtains the log evidence and it is equivalent to minimizing $KL(q(\mathbf{Z})||p(\mathbf{Z}|\mathbf{X}))$:

$$\begin{aligned} \mathcal{L}(q) &\triangleq \int \log \frac{p(\mathbf{X}, \mathbf{Z})}{q(\mathbf{Z})} dq(\mathbf{Z}) = \int \log \frac{p(\mathbf{Z}|\mathbf{X})p(\mathbf{X})}{q(\mathbf{Z})} dq(\mathbf{Z}) \\ &= \int \log p(\mathbf{X}) dq(\mathbf{Z}) - \int \log \frac{q(\mathbf{Z})}{p(\mathbf{Z}|\mathbf{X})} dq(\mathbf{Z}) = \log p(\mathbf{X}) - KL(q(\mathbf{Z})||p(\mathbf{Z}|\mathbf{X})) \end{aligned} \quad (3.3)$$

We can rewrite equation 3.3 as:

$$\log p(\mathbf{X}) = \mathcal{L}(q) - KL(q(\mathbf{Z})||p(\mathbf{Z}|\mathbf{X})) \quad (3.4)$$

Equation 3.4 shows that the ELBO minus the kl-divergence are constant and equal the log evidence. Therefore minimizing the kl-divergence (which is always non-negative) simultaneously maximizes the ELBO and vicer-versa.

In reality we can't search in the set of **all** possible distributions $q(\mathbf{Z})$. Instead we limit our search to some parameterized family of simple distributions $\{q_\alpha(\mathbf{Z})\}_\alpha$ where α is the parameter set. For a simple one dimensional concrete example, we can consider the set of all Gaussian distributions $\{\mathcal{N}(\mathbf{z}; \mu, \sigma)\}_{\mu, \sigma}$.

The task therefore is to use the data \mathbf{X} to find the best parameter $\hat{\alpha}$ that maximizes $\mathcal{L}(q_\alpha)$:

$$\hat{\alpha} \triangleq \arg \max_\alpha (\mathcal{L}(q_\alpha)) \quad (3.5)$$

3.2 Variational Autoencoder

3.2.1 Adding parameters

Our models will not be fully Bayesian, but rather parametrized. In this case let θ represent the set of parameters for p , and ϕ the parameters for q . Meaning we are dealing with a family of distributions $p_\theta(x, z)$ and another family $q_\phi(z)$.

For any θ and any ϕ , the equations from the previous chapter hold also in the parametrized form, i.e $\log p_\theta(x) = \mathcal{L}(q_\phi) - KL(q_\phi(Z)||p_\theta(Z|X))$.

We assume that we can only approach the "real" distribution using θ from below $\log p(X) \geq \log p_\theta(X)$. So together with equation 3.2 we have

$$(\forall \theta, \phi) \log p(X) \geq \log p_\theta(X) \geq \mathcal{L}(q_\phi) = \int \frac{p_\theta(Z|X)}{q_\phi(Z)} dq_\phi(Z) \quad (3.6)$$

So from equation 3.4 we again see that by finding the parameters ϕ, θ that maximize the elbo we approach the real log evidence as much as we can within the limits of the parametrized family of distributions we use.

3.2.2 Using neural networks for the parametrization

In this text we deal with variational autoencoders (VAE). A VAE is a neural network which is used to define and optimize the parameters ϕ and θ .

Specifically the encoder part of the network is a non-linear map $f_\theta(Z)$ which is used to define $P_\theta(X|Z)$. For example, we can assume that P_θ is a family of multivariate Gaussians and in this case $f_\theta(Z) = (\mu(Z), \Sigma(Z))$. Meaning the encoder maps Z to the location vector and covariance matrix. The parameter θ in this case are the weights of the encoder neural network. In parametrized the prior $p(Z)$ however in this case its parameter is not a function of X . In practice there is no reason to do this for most VAEs and we choose some simple fixed prior distribution for $p(Z)$.

The decoder network is similarly defined as a non-linear function $g_\phi(X)$ which maps X into the parameters defining the family $q_\phi(Z)$. Here too ϕ represent the weights of the decoder.

Remark 3.1. In many papers about VAEs (including this text) the decoder is described as $q_\phi(X|Z)$. This is an abuse of notation, because there is neither joint distribution $q_\phi(X, Z)$ nor marginal $q_\phi(X)$ to speak of. However $q_\phi(Z)$ is meant to approximate the intractable $p(Z|X)$ and that's presumably the reason for the abused notation.

3.2.3 Graphical representation

It is both convenient as well as informative to include a graphical description of our probabilistic models by way of plate diagrams.

Please note that we drop the ϕ, θ subscript but they are still there in reality.

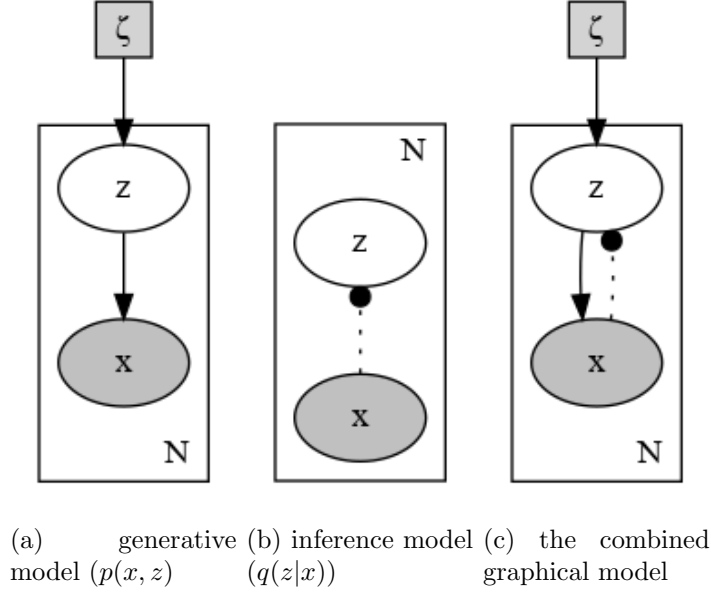


Figure 3.1: VAE graphical model

In a plate diagram nodes represent random variables and arrows represent dependency. Figure 3.1 is a plate diagram of the VAE model with slight adaptation. We use dotted arrows to represent the arrows of the inference model, and regular arrows for the generative model. Regular (triangular) arrowhead represents real probabilistic dependency whereas rounded arrows are reminding us that this is not a real probabilistic dependency (recall 3.1) which maybe we can call 'parametric dependency'.

Plate represents packing of N i.i.ds since we have N observations $X = (x_i)_1^N$ and correspondingly N latent variables $Z = (z_i)$.

Shaded node represent known values (either observation of prior).

The squared ζ node represent some fixed parameters which describes the prior distribution of $P(z)$. Usually it is not shown in the papers about VAE but we just wanted to remind the reader that it can be parametrize in general.

The generative model therefore factors as: $p(x, z) = p(x|z)p(z|\zeta) = p(x|z)p(z)$

The inference model in this case is just $q(z)$ but we might denote is as $q(z|x)$ because it tries to approximate $p(z|x)$.

Note that the graphical model has no assumption about the specific types of distributions involved (Gaussian, Dirichlet or whateve ...) and that is left for the actual implementation.

In the case of a "vanilla" VAE, We chose the prior to be diagonal standard Gaussian $p(z) \sim N(0, 1)$. And $p(z|x)$ is then chosed as diagonal Gaussian (whose means and variances we find by training the network).

Chapter 4

Gaussian mixture model VAEs

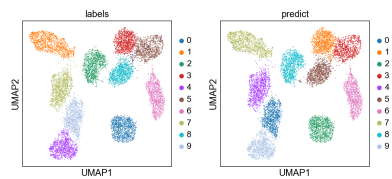
Theoretical background and with some examples from publications and my own tests.

4.0.1 Relation between AE and VAE

4.0.2 Conditional VAE



(a)



(b)

Figure 4.1: a figure

Chapter 5

Experiments and results

5.1 Tests with MNIST and FMNIST

5.2 Tests with scRNAseq Data

some words about (sc)RNAseq and published papers where AE and VAE models have been applied. What we were hoping to achieve and compare with.

Chapter 6

Discussion, some remarks and conclusions

punkt. punkt. punkt.

Bibliography

- [1] Christopher M Bishop and Nasser M Nasrabadi. *Pattern recognition and machine learning*. Vol. 4. 4. Springer, 2006.
- [2] Xifeng Guo et al. “Improved deep embedded clustering with local structure preservation.” In: *Ijcai*. 2017, pp. 1753–1759.
- [3] Diederik P Kingma and Max Welling. “Auto-encoding variational bayes”. In: *arXiv preprint arXiv:1312.6114* (2013).
- [4] Mohammad Lotfollahi, F Alexander Wolf, and Fabian J Theis. “Generative modeling and latent space arithmetics predict single-cell perturbation response across cell types, studies and species”. In: *bioRxiv* (2018), p. 478503.
- [5] Elad Plaut. “From principal subspaces to principal components with linear autoencoders”. In: *arXiv preprint arXiv:1804.10253* (2018).
- [6] Denis Serre. “Matrices: Theory & Applications Additional exercises”. In: *L’Ecole Normale Supérieure de Lyon* (2001).