



## List of topics to cover

With section titles and brief explanations.

Yiftach Kolb

Berlin, October 14, 2022

Freie Universität



Berlin

# Abstract

punkt. punkt.

# Declaration

punkt. punkt.

# Acknowledgement

punkt.[3] punkt. bip/bop/boop **XxZzYyWw**ℤ so-so-so—so

# List of Tables

# List of Figures

2.1	a neural network graphical description $\sigma(w_1x + w_2y + b)$ . . . . .	14
3.1	VAE graphical model . . . . .	20
4.1	a figure . . . . .	22

# Contents

<b>1</b>	<b>Introduction</b>	<b>8</b>
<b>2</b>	<b>Notations and definitions, preliminary concepts</b>	<b>9</b>
2.1	Basic notations . . . . .	9
2.2	The data . . . . .	10
2.3	Linear algebra preliminary: SVD and PCA . . . . .	10
2.4	Neural networks . . . . .	12
2.4.1	Universal families of parameterized maps . . . . .	12
2.4.2	Neurons . . . . .	12
2.4.3	Loss functions . . . . .	14
2.5	Autoencoders . . . . .	15
2.5.1	Relation between PCA and AE . . . . .	15
<b>3</b>	<b>Variance inference and variational autoencoders</b>	<b>16</b>
3.1	Variational Inference . . . . .	16
3.2	Variational Autoencoder . . . . .	17
3.2.1	Adding parameters . . . . .	17
3.2.2	Rearranging the ELBO . . . . .	18
3.2.3	Thinking of the data as distribution . . . . .	18
3.2.4	Using neural networks for the parametrization . . . . .	19
3.2.5	Graphical representation . . . . .	19
<b>4</b>	<b>Gaussian mixture model VAEs</b>	<b>21</b>
4.0.1	Relation between AE and VAE . . . . .	21

4.0.2	Conditional VAE . . . . .	21
<b>5</b>	<b>Experiments and results</b>	<b>23</b>
5.1	Tests with MNIST and FMNIST . . . . .	23
5.2	Tests with scRNAseq Data . . . . .	23
<b>6</b>	<b>Discussion, some remarks and conclusions</b>	<b>24</b>



# 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. To stress that we talk about matrices rather than vectors we show product ( $\times$ ) in the dimension, i.e  $\mathbf{X} \in \mathbb{R}^{m \times n}$ . Although technically the matrix-space is the tensor product  $\mathbb{R}^m \otimes \mathbb{R}^n$ .

Bold small math letters ( $\mathbf{x}$ ) represent usually row vectors, but in cases where it makes sense may also represent matrices such as a batch of several vectors (each row is a different data point). In few occasions it makes sense to let it represent both a matrix and a vector, for example,  $\boldsymbol{\sigma}$  may represent both the covariance matrix and the variance vector of a diagonal Gaussian distribution. Non-bold math letters ( $x, \sigma, \dots$ ) may represent scalar or vectors in some cases and hopefully it is clear from the context or explicitly stated.

Since we are only dealing with real matrices the transpose and the conjugation operators are the same ( $A^T = A^*$ ) but over  $\mathbb{C}$  conjugation is usually the "natural" operation and we use it to indicate that some property is still valid over  $\mathbb{C}$  with conjugation.

Sometimes matrices are given in row/column/block notations inside brackets where the elements are concatenated in a way that makes positional sense. For example both  $(\mathbf{x}, \mathbf{y})$  and  $(\mathbf{x}|\mathbf{y})$  represent a matrix with 2 **columns**.

As mentioned usually just  $\mathbf{x}$  means a column vector and  $\mathbf{x}^T$  means a row vector but sometimes in matrix notation  $\mathbf{x}$  represents a row when it makes sense. We use **curly** brackets to indicate the **row** representations of a matrix. For example  $\{\mathbf{x}, \mathbf{y}\}$  represents a matrix whose **rows** are  $\mathbf{x}$  and  $\mathbf{y}$  (as row vectors), which alternatively could be represented as  $(\mathbf{x}, \mathbf{y})^T$ .

$(\mathbf{X}, \mathbf{Y})$  represents concatenation of two matrices which implicitly means they have the same number of rows.

Zero-blocks are indicated with 0 or are simply left as voids. For example  $\begin{pmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{0} & \mathbf{D} \end{pmatrix}$  represents block notation of an upper-triangular matrix.

**Definition 2.1.** Let  $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_m\} \in \mathbb{R}^{m \times n}$  be a matrix in **row** notation. Then its *squared Frobenius norm* is

$$\|\mathbf{X}\|_F^2 \triangleq \text{trace}(\mathbf{X}\mathbf{X}^*) = \sum_{i=1}^m \|\mathbf{x}_i\|_2^2 = \sum_{i=1}^m \sum_{j=1}^n x_{ij}^2 \quad (2.1)$$

## 2.2 The data

we assume that the input data unless otherwise stated is real-valued matrix. Rows represent *samples* and columns represent *variables*.

We deal with two type of concrete datasets in this thesis. One of them is Single cell RNAseq data. This data represents gene expression levels in individual cells, where rows represent cells and columns represent genes. So if we see a reading of 0.5 in row 2 column 4 in means that in cell 2 gene 4 has normalized expression of 0.5.

The other type of data is images. For example the MNIST data set contains greyscale  $28 \times 28$  images of hand written digits. We still think of such data set as a matrix. The first axis always represents the samples, so each "row" represent an image. The rest of the axes represent the image. Alternatively we can also flatten the images into one axis and think of an image as a row vector of  $28 * 28$  dimensions.

There could possibly be additional data matrices with information about class or conditions. We use *one-hot encoding* to represent such information. For example in the case of the MNIST dataset every image also comes with a label which indicate which digit it shows. Since there are 10 classes of digits (0 to 9) the class matrix is going to have 10 columns and each row is a one-hot vector indicating the digit of the corresponding image.

**Definition 2.2.** 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 2.3.** A *class matrix*, or also a *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 a *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

In the following state some facts and bring without proof what are the singular value decomposition and the principle components of a matrix. For a full proof see [8].

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  $r = \text{rank}(\mathbf{X}) \leq \min(n, N)$ .

$\mathbf{X}\mathbf{X}^*$  and  $\mathbf{X}^*\mathbf{X}$  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  $s_1^2 \geq s_2^2 \geq \dots \geq s_r^2 > 0$ . 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  $\mathbf{U} = (\mathbf{u}_1 | \dots | \mathbf{u}_N) \in \mathbb{R}^{N \times N}$  be the (column) right eigenvectors of  $\mathbf{X}\mathbf{X}^*$  sorted by their eigenvalues. Then  $\mathbf{U} = (\mathbf{U}_r, \mathbf{U}_k)$  where  $\mathbf{U}_r = (\mathbf{u}_1 | \dots | \mathbf{u}_r) \in \mathbb{R}^{N \times r}$  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 let  $\mathbf{V} = (\mathbf{V}_r, \mathbf{V}_k) \in \mathbb{R}^{n \times n}$  be the (column) right eigenvectors of  $\mathbf{X}^*\mathbf{X}$ , sorted by the eigenvalues, where  $\mathbf{V}_r = (\mathbf{v}_1 | \dots | \mathbf{v}_r) \in \mathbb{R}^{n \times r}$  are the first  $r$  eigenvectors and  $\mathbf{V}_k$  are the  $n - r$  null-eigenvectors.

The critical observations is that  $\mathbf{V}_r = \mathbf{X}^* \mathbf{U}_r \mathbf{S}^{-1}$  and then  $\mathbf{U}_r^* \mathbf{X} \mathbf{V}_r = \mathbf{S}$ .

The *singular value decomposition (SVD)* of  $\mathbf{X}$  is

$$\mathbf{X} = \mathbf{U} \mathbf{D} \mathbf{V}^* \quad (2.2)$$

where  $\mathbf{D} = \begin{pmatrix} \mathbf{S} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix} \in \mathbb{R}^{N \times n}$  is diagonal.

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

$$\mathbf{V}_l = \underset{\mathbf{W}}{\text{argmin}} \{ \|\mathbf{X} - \mathbf{X} \mathbf{W} \mathbf{W}^T\|_F^2 \quad : \quad \mathbf{W} \in \mathbb{R}^{n \times l}, \mathbf{W}^T \mathbf{W} = \mathbf{I}_l \} \quad (2.3)$$

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

If we consider the more general minimization problems:

$$\min_{\mathbf{E}, \mathbf{D}} \{ \|\mathbf{X} - \mathbf{X} \mathbf{E} \mathbf{D}\|_F^2 \quad : \quad \mathbf{E}, \mathbf{D}^T \in \mathbb{R}^{n \times l}, \} \quad (2.4)$$

$$\min_{\mathbf{W}} \{ \|\mathbf{X} - \mathbf{X} \mathbf{W} \mathbf{W}^\dagger\|_F^2 \quad : \quad \mathbf{W} \in \mathbb{R}^{n \times l}, \} \quad (2.5)$$

It can be shown [6] that the last two problems 2.4, 2.5 are equivalent and that for any solution  $\mathbf{E}, \mathbf{D}$  it must hold that  $\mathbf{D} = \mathbf{E}^\dagger$ . ( $\mathbf{D}$  is the Moore–Penrose generalized inverse of  $\mathbf{E}$ ). Moreover,  $\mathbf{V}_l$  still minimizes the general problem 2.4 and for every solution  $\mathbf{W}$ , it must hold that  $\text{span}\{\mathbf{W}\} = \text{span}\{\mathbf{V}_l\}$  (but it isn't necessarily an orthogonal matrix).

## 2.4 Neural networks

We briefly discuss here some of the basics of neural network to provide clarity and motivation. Mostly based on [5].

### 2.4.1 Universal families of parameterized maps

If we take an expression such as  $f_{a,b}(x) = ax + b$ , if we hold  $(a, b)$  fixed on specific values, then we get a linear function on  $x$ . Every assignment of  $(a, b)$  defines a different linear function and in fact every linear function on one dimension can be uniquely described by these  $a$  and  $b$ . So we can say that  $\{f_{a,b}\}_{a,b \in \mathbb{R}}$  is a *parameterization* of the class of all real linear functions on one variable. The distinction between what are the variables and what are the parameters is somewhat arbitrary and in the end,  $f_{a,b}(x)$  is just another way to represent a 3-variable function  $f(a, b, x)$ .

In general we can define one or more multivariate functions  $g : \mathbb{R}^{n+m} \rightarrow \mathbb{R}^k$  (for simplicity of the discussion let's assume it is defined everywhere) and partition the set of its variables into 2.  $g_{\mathbf{w}}(x) \triangleq g(\mathbf{w}, \mathbf{x})$  where  $\mathbf{x} \in \mathbb{R}^n$ ,  $\mathbf{w} \in \mathbb{R}^m$  and let  $g_{\mathbf{w}}(\mathbf{x}) := g(w, x) \in \mathbb{R}^k$ .

We call a class  $\mathcal{F}$  of parameterized functions *universal* if every continuous function can be uniformly approximated by functions of that class. The class of all linear functions is not universal. But taking "any function"  $g$  is too general. What we actually want is a class of parameterized functions that is:

- as simple as possible to construct
- derivable in both the parameters as well as the variables
- can uniformly approximate any continuous function in a bounded domain given sufficiently large set of parameters (i.e. is universal).

But if these requirements are still not enough. For example, the class of multivariate polynomials can uniformly approximate any function. However it may not be a good idea to try to learn very complicated high dimensional data using polynomial representation. For one reason is that the number of terms (monomials) grows very rapidly with the dimension and the degree of the polynomials: for  $n$  dimensions and  $m$  degrees there are something like  $\binom{m+n}{n}$  monomial terms.

We want this class of simpler functions, that are almost as simple as linear, and yet that suits well for statistical learning. For example we want to represent complicated functions with relatively few parameters.

One such class of functions is the feed forward neural networks, which is the class of functions that are comprised from "neurons".

### 2.4.2 Neurons

So first let's see what a neuron is. Basically it is the simplest parametrized scalar function that is not linear. It is simple in the sense that it applies a linear function on the variables and then applies a so called activation function on the output of the linear function.

**Definition 2.4.** An *activation function*  $\sigma : \mathbb{R} \rightarrow \mathbb{R}$  is any one of the following functions:  $x \mapsto 1$  (constant),  $x \mapsto x$  (identity),  $x \mapsto \frac{e^x}{1+e^x}$  (sigmoid), and  $x \mapsto \max(0, x)$  (ReLU).

If  $\mathbf{x} = (x_1, \dots, x_n) \in \mathbb{R}^n$  then  $\sigma(\mathbf{x})$  is the element-wise application  $\sigma(\mathbf{x}) \triangleq (\sigma(x_1), \dots, \sigma(x_n))$ .

**Definition 2.5.** Let  $\sigma : \mathbb{R} \rightarrow \mathbb{R}$  be an activation function and let  $f_{\mathbf{w}} : \mathbb{R}^n \rightarrow \mathbb{R}$  be a *parameterized* linear function. A *neuron*  $\nu$  is the parameterized function  $\nu = \nu_{\mathbf{w}} \triangleq \sigma \circ f_{\mathbf{w}}$ .

The parameters  $\mathbf{w}$  are called the *weights* of the neuron  $\nu$ .

Think of the weights of a neuron as some mutable, tunable property, some sort of memory.

Connecting many neurons together can create pretty powerful parameterized functions which we call neural networks. Connecting means that the output of one neuron is the input to one of the variables of a different neuron. In feed forward networks the information only goes in one direction (no feedback).

**Definition 2.6.** A *feed forward neural network (NN)* is a **parameterized** map  $\phi$  recursively defined follows:

- Activation functions (1, *id*, and  $\sigma$ ) are NNs which are called the *elementary neurons* and they have no parameters ( $\mathbf{w} = \emptyset$ ).
- neurons are NNs
- If  $\psi : \mathbb{R}^n \rightarrow \mathbb{R}$  is NN we say it has a single output (SONN).
- if  $\nu : \mathbb{R}^n \rightarrow \mathbb{R}$  is a NN and  $\psi_1, \dots, \psi_n$  are SONNs, and if the parameter set of  $\nu$  is disjoint from the combined parameters of the  $\psi_i$ 's, then  $\phi = \nu(\psi_1, \dots, \psi_n)$  is a NN.
- if  $f_{\mathbf{w}} : \mathbb{R}^n \rightarrow \mathbb{R}^m$  is a parameterized linear map,  $\psi_1 \dots \psi_n$  are NNs, and the parameter sets of all the maps and NNs involved are pair-wise disjoint, then  $\phi = f_{\mathbf{w}}(g_1, \dots, g_n)$  is a NN.
- if  $\psi$  is a NN, and  $\nu$  is an neuron then  $\phi = \nu(\psi)$  is a NN, which means  $\nu : \mathbb{R} \rightarrow \mathbb{R}$  is applied element-wise on the output of  $\psi$ .

The parameter set  $\mathbf{w}$  is called the *weights* of  $\phi$

Feed forward neural network are depicted as a directed acyclic graph and usually they are arranged in layers. In that description a neuron consists of a node in the graph, the node label (its activation function), incoming edges, and the edge labels.

The label inside the neuron describe its activation function. In the diagrams, we let  $\sigma$  represent the sigmoid function. We represent the identity function either by the name of the variable ( $x_1, y$  etc.) it acts on or simple by *id*. We let the label 1 represent the constant function. We need the constant function because with it we can represent any affine map as a linear map with the first input always clamped to 1. An directed edge between neurons means that the output of the neuron at the tail is multiplied by the edge weight and assigned to the input variable of the neuron it connects to. Input neurons (sources) have no incoming edges and they represent the begining of the computation. Output neurons (sinks) have no outgoing edges and their output is the final result of the computation.

Figure 2.1 shows a single output neural network of four neurons, with three of . input neurons and one output. It describes the function  $f : \mathbb{R}^2 \rightarrow \mathbb{R}, f(x, y) = \sigma(w_1x + w_2y + b)$ .

It turns out [5] that the feed forward neural networks class of functions (will be properly defined later) are "universal", meaning any continuous function  $f$  can be uniformly approximated by a neural network (in fact, a network with just a single hidden layer suffices and just ReLU or Sigmoid as the non-linear activation functions). More precisely, let  $B \subseteq \mathbb{R}^n$  be a bounded domain. Let  $f : B \rightarrow \mathbb{R}^m$  be continuous, and let  $\epsilon \in (0, 1)$ . Then there is a feed forward neural network with a single hidden layer  $\phi = \phi_{\mathbf{w}}$  and there is some value assignment for the parameters  $\mathbf{w}$  such that  $\forall \mathbf{x} \in B \|\phi(\mathbf{x}) - f(\mathbf{x})\| < \epsilon$ . The size of that single hidden layer (the number of parameters) depends on  $f$  and  $\epsilon$ .

### 2.4.3 Loss functions

**Definition 2.7.** 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  $\mathcal{L}(\mathbf{x}) = \sum_{i=1}^m \psi(x_i)$

An example for such a loss function is the square error:  $\mathbf{x} \mapsto \|\mathbf{x}\|_2^2$

**Definition 2.8.** Let  $\mathbf{X} \in \mathbb{R}^{N \times n}$  be a data matrix. A *batch*  $\mathbf{x} \in \mathbb{R}^{b \times n}$  is any subset of  $b$  rows of  $\mathbf{X}$  (Note that in this case  $\mathbf{x}$  represents a matrix).

Batch  $\mathbf{x} = \{\mathbf{x}_1, \dots, \mathbf{x}_b\} \in \mathbb{R}^{b \times n}$  (row notation) represents a subset of  $b$  samples out of the total of  $N$  samples in the dataset. The operations of  $\phi$  and  $\mathcal{L}$  *extend* to batches in natural ways—collect for  $\phi$  and we average for  $\mathcal{L}$ , namely:

**Definition 2.9.** Let  $\phi$  be a neural network as defined in 2.6 and let  $\mathcal{L}$  its associated loss function as defined in 2.7—over vectors. Let  $\mathbf{x} = \{\mathbf{x}_1, \dots, \mathbf{x}_b\} \in \mathbb{R}^{b \times n}$  be a  $b$ -batch (in row notation). Then  $\phi$  and  $\mathcal{L}$  *extended* over batches are:

$$\phi(\mathbf{x}) \triangleq \{\phi(\mathbf{x}_i)\}_{i=1}^b \in \mathbb{R}^{b \times m} \quad (2.6)$$

$$\mathcal{L}(\mathbf{x}) \triangleq \frac{1}{b} \sum_{i=1}^b \mathcal{L}(\mathbf{x}_i) \in \mathbb{R} \quad (2.7)$$

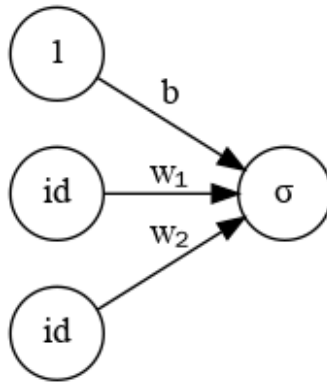


Figure 2.1: a neural network graphical description  $\sigma(w_1x + w_2y + b)$

If  $\mathcal{L}$  is the square error function  $\|\cdot\|_2^2$  on vectors, then its extension 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  $\phi_w$  means finding the weights that minimize the loss function applied on the training set  $X$ , in other words minimizing  $\min_w(\mathcal{L}(\phi_w(X)))$ .

**Definition 2.10.** Let  $\phi_\omega$  be a neural network and  $\mathcal{L}$  its associated loss function. And let the data matrix  $\mathbf{X}$  be our *training set*. Then *Training* of  $\phi_\omega$  with respect to  $\mathcal{L}, \mathbf{X}$  means algorithmically approximating the minimization problem:

$$\min_{\omega} \mathcal{L}(\phi_{\omega}(\mathbf{X})) \quad (2.8)$$

During a *training step* the network is applied on a batch  $x$ . Then the loss function is applied on 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 [7]) and can be regarded as a black box.

## 2.5 Autoencoders

**Definition 2.11.** 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 2.3; Whereas a **linear** autoencoder solves equation 2.5. 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  $\phi$  is linear) is therefore almost equivalent to PCA [6], 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:

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

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

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

In equation 3.2 we found a lower bound  $\mathcal{L}(q, p)$  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  $KL(q(\mathbf{Z})||p(\mathbf{Z}|\mathbf{X}))$ :

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

We can rewrite equation 3.3 as:

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

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.

## 3.2 Variational Autoencoder

### 3.2.1 Adding parameters

Our models will not be fully Bayesian, but rather parametrized. In this case let  $\theta$  represent the set of parameters for  $p$ , and  $\phi$  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  $\theta$  and any  $\phi$ , the equations from the previous chapter hold also in the parametrize 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  $\theta$  from below  $\log p(X) \geq \log p_\theta(X)$ . So together with equation 3.2 we have

$$(\forall \theta, \phi) \log p(X) \geq \log p_\theta(X) \geq \mathcal{L}(q_\phi) = \int \frac{p_\theta(X, Z)}{q_\phi(Z)} dq_\phi(Z) \quad (3.5)$$

So from equation 3.4 we again see that by finding the parameters  $\phi, \theta$  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 Rearranging the ELBO

Equations 3.2 and 3.3 were defined for any distribution  $q(\mathbf{Z})$  and in particular we are allowed to plug in a conditioned distribution  $q(\mathbf{Z}|\mathbf{X})$ . That implies the existence of  $q(\mathbf{Z}, \mathbf{X})$  and  $q(\mathbf{X})$  but we actually don't care about them. We condition everything on  $\mathbf{X}$  but  $\mathbf{X}$  is treated as a given constant from a Bayesian view point and we only want to somehow make  $q(\mathbf{Z}|\mathbf{X})$  to closely approximate  $p(\mathbf{Z}|\mathbf{X})$ .

A second thing we need to achieve is to express the ELBO in terms of  $p(\mathbf{X}|\mathbf{Z})$  and  $q(\mathbf{Z}|\mathbf{X})$  rather than the joint distribution. To that end we need also the prior  $p(\mathbf{Z})$ .

**Definition 3.3.** The *conditioned (on  $\mathbf{X}$ ) ELBO* is

$$\begin{aligned} \mathcal{L}(q, p) &\triangleq \int \log \frac{p(\mathbf{X}, \mathbf{Z})}{q(\mathbf{Z}|\mathbf{X})} dq(\mathbf{Z}|\mathbf{X}) = \int \log \frac{p(\mathbf{X}|\mathbf{Z})p(\mathbf{Z})}{q(\mathbf{Z}|\mathbf{X})} dq(\mathbf{Z}|\mathbf{X}) \\ &= \int \log p(\mathbf{X}|\mathbf{Z}) dq(\mathbf{Z}|\mathbf{X}) - \int \log \frac{q(\mathbf{Z}|\mathbf{X})}{p(\mathbf{Z})} dq(\mathbf{Z}|\mathbf{X}) \\ &= \int \log p(\mathbf{X}|\mathbf{Z}) dq(\mathbf{Z}|\mathbf{X}) - KL(q(\mathbf{Z}|\mathbf{X}) || p(\mathbf{Z})) \end{aligned} \quad (3.6)$$

So to sum it up, if we want to maximize the log evidence  $\log p(\mathbf{X})$  it suffices to maximize  $\mathcal{L}$ , and equation 3.3 shows that this means finding the balance between making the right term (which we call the reconstruction term) large as possible, and making the KL-term small. The KL term is seen as a regularization term.

### 3.2.3 Thinking of the data as distribution

We can think of any data point  $\mathbf{x}_i \in X$ , rather than being a deterministic point, as if is a sample which comes from some distribution  $p_i$ . For example if  $\mathbf{x}_i$  is a vector of binary data, we can think of it as if it was generated by some Bernoulli distribution. If the data is in the range  $[0, 1]$  we can think of every data point as a vector of Bernoulli probabilities rather than a concrete sample. For other type of real data, we think of it as coming from a diagonal Gaussian distribution, where we usually assume the variance is fixed 1 or more generally that all the samples have some common variance  $\sigma$ .

Similarly we can also treat the latent variable as coming from some distribution rather than being deterministic. Actually that is what equation 3.3 is all about.

### 3.2.4 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  $\phi$  and  $\theta$  which define  $p_\theta(\mathbf{X}|\mathbf{Z})$  and  $q_\phi(\mathbf{Z}|\mathbf{X})$ .

Specifically the encoder part of the network is a non-linear map  $f_\theta(\mathbf{Z})$  which is used to define  $P_\theta(X|Z)$ . For example, we can assume that  $P_\theta$  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  $\theta$  in this case are the weights of the encoder neural network. In parametrizing 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  $\phi$  represent the weights of the decoder.

**Definition 3.4.** Let  $\{p_\theta\}$  be a family of distributions over  $\mathbb{R}^n$  and let  $\{q_\phi\}$  be a family of distributions over  $\mathbb{R}^m$  and let  $p$  be some fixed distribution over  $\mathbb{R}^m$ . A *variational autoencoder (VAE)*  $v = (f_\omega, g_\omega)$  consists of neural networks  $f_\omega : \mathbb{R}^n \rightarrow \{\theta\}_\theta$  and  $g_\omega : \mathbb{R}^m \rightarrow \{\phi\}_\phi$  where  $\omega$  are the combined parameters of  $f, g$  and  $\{\theta\}$  means the parameter space (all possible values of  $\theta$  which can define a valid distribution  $p_\theta$  (similarly with  $\phi$ )).

If an autoencoder works on deterministic data, with the encoder mapping input  $\mathbf{x} \mapsto \mathbf{z}$  and the decoder then maps the latent space  $\mathbf{z} \mapsto \hat{\mathbf{x}}$  to the reconstruction, a VAE tries to do basically the same thing but non-deterministically. It maps  $\mathbf{x}$  into a distribution over  $\mathbf{z}$ :  $\mathbf{x} \mapsto q(\mathbf{z}|\mathbf{x})$  and it maps  $\mathbf{z}$  into a distribution over  $\mathbf{x}$ ,  $\mathbf{z} \mapsto p(\mathbf{x}|\mathbf{z})$ .

The loss function associated with a VAE is usually the minus ELBO.

### 3.2.5 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  $\phi, \theta$  subscript but they are still there in reality.

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 ??) 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 or prior).

The squared  $\zeta$  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

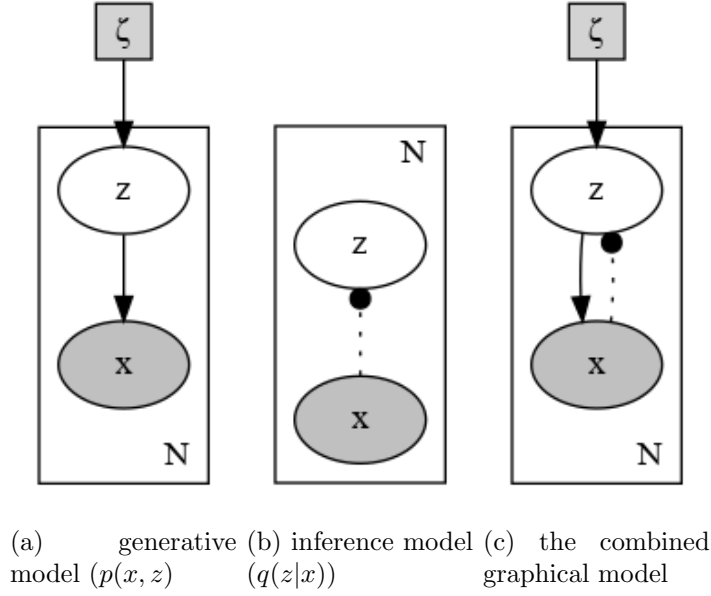


Figure 3.1: VAE graphical model

remind the reader that it can be parametrized 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 it 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 whatever ...) 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 chosen as diagonal Gaussian (whose means and variances we find by training the network).

## Chapter 4

# Gaussian mixture model VAEs

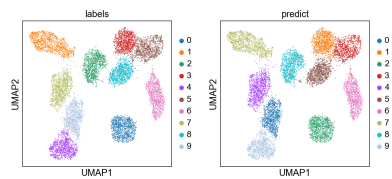
Theoretical background and with some examples from publications and my own tests.

### 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] Michael A Nielsen. *Neural networks and deep learning*. Vol. 25. Determination press San Francisco, CA, USA, 2015.
- [6] Elad Plaut. “From principal subspaces to principal components with linear autoencoders”. In: *arXiv preprint arXiv:1804.10253* (2018).
- [7] Automatic Differentiation In Pytorch. *Pytorch*. 2018.
- [8] Denis Serre. “Matrices: Theory & Applications Additional exercises”. In: *L’Ecole Normale Supérieure de Lyon* (2001).