



VRIJE
UNIVERSITEIT
BRUSSEL



Master thesis submitted in partial fulfilment of the requirements for the degree of
Master of Science In de Ingenieurswetenschappen: Computerwetenschappen

VARIATIONAL GREEDY INFOMAX

Towards indepdenent and interpretable
representations

Fabian Denoodt

2022-2023

Promotor(s): Prof. Dr. Bart de Boer
Science and Bio-Engineering Sciences



VRIJE
UNIVERSITEIT
BRUSSEL



Proefschrift ingediend met het oog op het behalen van de graad van Master of
Science In de Ingenieurswetenschappen: Computerwetenschappen

[DUTCH] VARIATIONAL GREEDY INFOMAX

[Dutch] Towards indepednent and
interpretable representations

Fabian Denoodt

2022-2023

Promotor(s): Prof. Dr. Bart de Boer

Wetenschappen en Bio-ingenieurswetenschappen

Contents

1	Introduction	5
2	Background	7
2.1	Entropy, relative entropy and mutual information	7
2.1.1	Shannon’s Entropy	7
2.1.2	Relative entropy and mutual information	8
2.2	Supervised neural networks	8
2.3	Representation learning through reconstruction error	10
2.3.1	Autoencoders	10
2.3.2	Variational autoencoders	11
2.4	Representation learning through Noise-contrastive estimation	14
2.4.1	Contrastive predictive coding	14
2.4.2	Greedy InfoMax	16
3	Variational InfoNCE	19
3.1	Motivation	19
3.2	Towards decoupled training for probabilistic representations	20
3.3	The learning objective	22
3.3.1	Gradient	22
3.4	Practical benefits	23
3.4.1	Latent representations	23
3.4.2	Training	23
4	Experiments	25
4.1	Greedy Infomax model	25
4.1.1	Dataset	25
4.1.2	Training and validation during training	28
4.2	Decoders for Greedy Infomax	29
4.3	Feature selection for visualisation	29
4.4	Variational contrastive predictive coding	29
4.5	CHATGPT ANSWERS FOR VAE	30
4.6	Results	31
4.6.1	Results CPC Simple v2 w/ 2 modules	31
4.7	GIM: Activations visualisations	32
4.8	Decoder: predictions on test set	34

5	Related work	39
5.1	Representation learning: explainable	39
5.2	Variational learning	39
5.3	Links I should investigate	39
6	Discussion	41

Chapter 1

Introduction

FEEDBACK BART

!!! wat is + waarom het gedaan wordt.

should also explain in my work.

should be around 60 pages.

Discussion: Sectie 4.5: (Bart wil aparte sectie) - Hoe helpt deze techniek om het netwerk interpreteerbaarder te maken. - compare techniques: whether they are better explainable. - Zou na results moeten zijn. en kan zo inleiding zijn naar future work. - bart zou verwachten dat cha 3 zeer groot is.

- UVA: wordt geschreven voor begeleiders, hun hebben meer achtergrond. - Jury aan VUB andere verwachtingen, wil meer uitleg.

- mijn discussie sectie moet langer! en zeker in vertellen waarom het beter is voor visualities.

— Defense: - verduidelijkende vragen - critiek als gaten in argumentatie - hun komen met suggestie die ze uit literatuur kennen, en moet bv mening over geven. - "zou dat ook anders gekund hebben" !! —

- Context: !E GIM for representation learning: generates representations that simplify classification tasks vs when done on raw data
- Problem: If wants to know for what tasks applicable... must know what is present in data
- Solution: Analysis of learned representations on speech data
- My contributions:
 - Decoder ANN for each layer of GIM: Shows what information is maintained through the layers
 - Search correlations between kernels weights and signal features
 - Extension on CPC via VAE

Chapter 2

Background

2.1 Entropy, relative entropy and mutual information

We first discuss specific definitions from information theory. These concepts will be relevant to understand contrastive predictive coding, which we discuss in a following section. The formal definitions are obtained from the book "Elements of Information Theory" [1]. The equations that contain a log function are assumed to be under base two.

2.1.1 Shannon's Entropy

Entropy measures the average amount of information required to describe a random variable [1]. The entropy $H(X)$ of a discrete random variable X , is formally defined in equation 2.1 shown below.

$$H(X) = - \sum_{x \in \mathcal{X}} p(x) \log p(x) \quad (2.1)$$

The alphabet \mathcal{X} represents the set of events the random variable X can take, formerly known as the *sample space*. Additionally, $p : \mathcal{X} \rightarrow [0, 1]$ denotes the probability density function of X . Hence, given an event $x \in \mathcal{X}$, $p(x)$ corresponds to the probability of event x occurring.

Assume a random variable X with possible events x_1, x_2 . Intuitively, when $p(x_1)$ is low, the surprise when the event x_1 occurs will be high. The surprise for one event is denoted in equation 2.2.

$$-p(x) \log p(x) \quad (2.2)$$

Hence, entropy can also be considered as the sum of "surprise" over each event [2]. To understand why equation 2.2 does indeed correspond to a measure of surprise, consider an event $x \in \mathcal{X}$ with $p(x) = 1$. Note that $\log p(x) = 0$, and thus the surprise is zero. Meanwhile, if $p(x)$ approaches 0, $\log p(x)$ goes to $-\infty$. And hence, by the negation sign in formula 2.1 the surprise is large.

Figure 2.1 displays when entropy reaches its maximum for the case of a random variable with 2 outcomes. We can see that the entropy, and thus, the information is largest when the probability of the two outcomes is equal to each other, namely $p(x_1) = p(x_2) = 0.5$. Note that for a random variable X with more than two events, $H(X)$ can be larger than one.

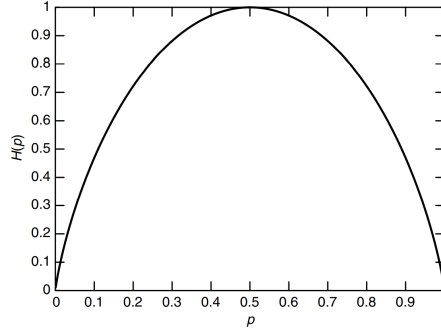


Figure 2.1: $H(p)$ vs p (originates from "Elements of Information Theory", page 16)

2.1.2 Relative entropy and mutual information

Relative entropy, also known as the Kullback Leibler (KL) divergence, is defined in equation 2.3, where p and q denote a probability density function over the same sample space \mathcal{X} [1]. The KL divergence quantifies the "divergence" or "distance" between the two distributions. Note that $D(p||q)$ does not necessarily correspond to $D(q||p)$ and thus the metric is not symmetrical.

$$D_{KL}(q || p) = \sum_{x \in \mathcal{X}} p(x) \log \frac{p(x)}{q(x)} \quad (2.3)$$

The mutual information (MI) between two random variables X and Y can be computed as the KL divergence between their joint probability distribution, $p_{X,Y}(x,y)$, and the product of their marginal probability distributions, $p_X(x)$ and $p_Y(y)$, which is denoted as $p_X(x)p_Y(y)$ [1]. The equation for mutual information then becomes:

$$I(X;Y) = D_{KL}(p_{X,Y}(x,y) || p_X(x)p_Y(y)) \quad (2.4)$$

As described by Cover and Thomas in their book "Elements of Information Theory" [1], $I(X;Y)$ quantifies the amount of information Y describes about X . A alternative definition for $I(X;Y)$ is illustrated in 2.5. The equation provides us with an intuitive meaning for $I(X;Y)$, corresponding to the surprise caused by X , which is reduced by the knowledge of Y . In a following section, we discuss how these concepts from information theory are applied in representation learning, by maximising the mutual information between latent representations.

$$I(X;Y) = H(X) - H(X|Y) \quad (2.5)$$

2.2 Supervised neural networks

We shall now discuss traditional supervised learning approaches, as these will lay the groundwork for the representation learning approaches discussed in the following section.

Typical supervised machine learning problems consider the problem where given a training set of tuples (x_i, y_i) , a mapping function $f : \mathcal{X} \rightarrow \mathcal{Y}$ must be learned, where $x_i \in \mathcal{X}$ corresponds to the input sample and $y_i \in \mathcal{Y}$ the label. A good f will then also generalise well to unseen $x_i \in \mathcal{X}$, which were not part of the training set.

Artificial Neural Networks (ANNs) in particular, tackle this problem by defining f as a fixed set of parameters consisting of a series of layers. During inference, at each layer l a transformation

matrix W^l is applied to the output vector from the previous layer a^{l-1} . This is shown in the equation below.

$$z^l = W^l a^{l-1}$$

Secondly, a non-linear function $\sigma : R^d \rightarrow R^d$ is applied to z^l , as shown in equation 2.6. The resulting vector a^l may then again be the input for a following layer.

$$a^l = \sigma(z^l) \quad (2.6)$$

Hence, during inference, the input vector x is propagated through each layer, resulting in a final output \hat{y}^L . The equation for the forward pass of a neural network with L layers is described in equation 2.7.

$$f_{W^1 \dots W^L}(x) = \hat{y}^L = \sigma(\dots \sigma(\sigma(x^T W^1)^T) W^2 \dots)^T W^L \quad (2.7)$$

During training the ANN's parameters $W^1 \dots W^L$, are optimised according to the learning problem. The viability of the parameters is quantified by the loss function \mathcal{L} over a batch of n training samples, as shown in equation 2.8. y_i corresponds to the ground truth label of the single data sample x_i , $e(y_i, \hat{y}_i)$ the error between a ground truth y_i and the estimation \hat{y}_i of a single sample. The error function e can be replaced depending on the task, for instance by mean squared error for regression or cross entropy for classification problems.

$$\mathcal{L}(W) = \frac{1}{n} \sum_{i=1}^n e(y_i, f_W(x_i)) \quad (2.8)$$

The parameters can be optimised by minimising the loss function \mathcal{L} defined above. This can be done algebraically, however, it becomes difficult when the dimension of the input features is large. Gradient descent is used to find an "appropriate" minimum of the cost function, by iteratively adjusting the parameters $W^1 \dots W^L$ via the following learning rule.

$$W_{ij}^l \leftarrow W_{ij}^l - \alpha \frac{\partial \mathcal{L}}{\partial W_{ij}^l} \quad (2.9)$$

Although gradient descent can be successfully applied for finding local minima, neural networks may still contain many parameters which each must be optimised. This is typically resolved by applying the backpropagation algorithm in combination with gradient descent, such that fewer partial derivatives must be calculated, resulting in more efficient optimisation.

Although supervised learning though ANNs is considered successful, their performance is heavily dependent on the choice of the data representation [3], on which they are applied. For complex representations, a more complex architecture may be required, which requires more data to prevent overfitting. When labelled data is scarce, these models may result in seemingly well performance on the training set, but bad generalisation to data outside the training set. As a consequence, a lot of time in the machine learning pipeline is invested in transforming data into better latent representations. Good latent representations tend to disentangle the data into meaningful features, such that data can be more easily be separable into classes. In the following section we discuss how part of the labour of finding good representations can be relieved with unsupervised learning approaches, which may learn to find good latent representations. These representations could then be used as the input for supervised predictors.

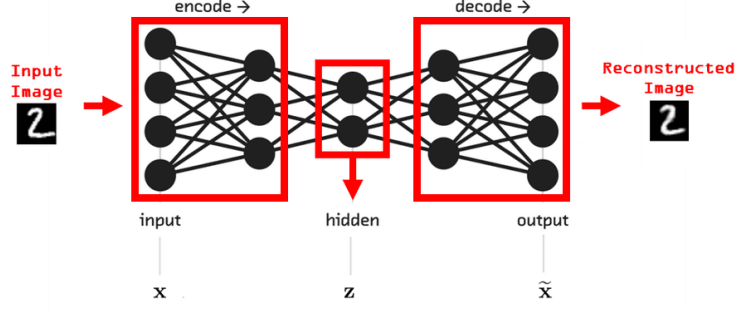


Figure 2.2: Autoencoder neural network architecture adapted from [7].

2.3 Representation learning through reconstruction error

One of the challenges in supervised learning is the constant need of large amounts of labelled data. Hence, when a labelled dataset is small, we would like to leverage a larger unlabelled dataset as basis for learning. By doing so, a mapping can be learned from the raw input data to a representation which makes downstream tasks easier. This process of learning representations from data is commonly referred to as representation learning [4]. Supervised learning algorithms can then learn directly from these disentangled latent representations with fewer labelled data.

In the following two sections we discuss two paradigms of representation learning with ANNs. The first paradigm is learning representations by minimising a reconstruction error. The second learns its representations by contrasting them against noise. These two paradigms will lay the basis for our own contributions in section three.

2.3.1 Autoencoders

Autoencoders were introduced in 1986 by D. Rumelhart et al. [5] as a means to learn compressed representations [6]. This is achieved through a neural network architecture consisting of two blocks. The first block is the encoder E and receives input data which it encodes into a lower dimensional representation. The second block, called the decoder D , receives as input the latent representation and is tasked to reconstruct the original input. The two blocks are simultaneously optimised by *minimising the reconstruction error* shown in the following equation:

$$\mathcal{L} = \sum_{i=1}^N l(\mathbf{x}^{(i)}, D(E(\mathbf{x}^{(i)}))) \quad (2.10)$$

where l refers to the error for a single data point, for instance the l_2 -norm and N the number data points. An example autoencoder architecture is depicted in figure 2.2. $\mathbf{x} \in \mathcal{D}$ corresponds to an input vector which is encoded into latent representation \mathbf{z} . The decoding of \mathbf{z} corresponds to $\tilde{\mathbf{x}} = D(E(\mathbf{x}))$. The dimension of \mathbf{z} is typically bottlenecked to be smaller than the original dimension of \mathbf{x} . This results in the encoder having to define encodings that are as "informative" as possible to reconstruct the original data [6]. \mathbf{z} may then for instance be used for downstream tasks such as classification or clustering.

While capable of learning compressed representations, autoencoders do not pose any restrictions on the latent space they define (the space of latent representations). As a result, the representations may be meaningful to computers, but non-interpretable to humans. For instance, given the left image depicted in figure 2.3 which depicts an autoencoder's two dimensional latent

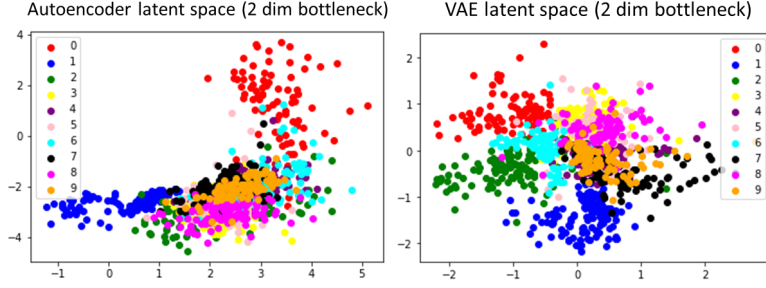


Figure 2.3: Both images represent a the two-dimensional latent space, learned from the MNIST dataset [8]. This is a dataset consisting of images of handwritten numbers between 0 and 9. Each image is associated to a class label, referring to the number. The left image consists of the space learned from a classical autoencoder, the right of a variational autoencoder (VAE). Both autoencoders have not received any explicit information of the labels of the dataset, yet learned to define representations which equal numbers are closer to each other. The VAE’s latent space, depicted in the right image, is optimised to be standard normally distributed. The latent vectors \mathbf{z} are distributed according to the two-dimensional normal distribution $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$, where μ is the two dimensional zero vector. \mathbf{I} is the 2×2 covariance matrix with ones on the diagonal and zeroes elsewhere.

space, it is very difficult to know what the resulting images would be when interpolating between the latent representations of 0 (red) and 1 (blue). Answering this question becomes even more infeasible for higher dimensional latent vectors.

2.3.2 Variational autoencoders

Similar to traditional autoencoders, variational autoencoders (VAE) learn representations that contain the important information that is necessary to reconstruct the data, *however*, an additional constraint is applied to the latent space [9, 10, 11, 12, 13]. The data points from the latent space are defined in such a way that they conform to a certain distribution, typically, the standard normal distribution $\mathcal{N}(\mathbf{0}, \mathbf{I})$ [10]. This behaviour can be observed in the right plot of figure 2.3. Since the latent representations conform a two-dimensional standard normal distribution $\mathcal{N}(\mathbf{0}, \mathbf{I})$, the data samples are more likely to be near the center at (0, 0).

Simulating distributions through neural networks

In variational autoencoders the latent representations of a data point $\mathbf{x}^{(i)}$ does not simply consist of a fixed deterministic vector $\mathbf{z}^{(i)}$, as was the case in the traditional autoencoder. Instead, given $\mathbf{x}^{(i)}$, its latent representation corresponds the following distribution $q_\phi(\mathbf{z} | \mathbf{x}^{(i)})$. A concrete vector $\mathbf{z}^{(i)}$ can be obtained taking a sample $\mathbf{z}^{(i)} \sim q_\phi(\mathbf{z} | \mathbf{x}^{(i)})$. As result a single $\mathbf{x}^{(i)}$ will correspond to multiple latent vectors $\mathbf{z}^{(i)}$. The latent representation $q_\phi(\mathbf{z} | \mathbf{x}^{(i)})$ is modelled as $\mathcal{N}(\mu^{(i)}, \text{diag}(\sigma^{(i)}))$. This means that all the required information to model $q_\phi(\mathbf{z} | \mathbf{x}^{(i)})$ is a mean vector $\mu^{(i)}$ and a covariance matrix $\text{diag}(\sigma^{(i)})$. Thus, given $\mathbf{x}^{(i)}$, a deterministic neural network can simulate $\text{diag}(\sigma^{(i)})$ by predicting two vectors, namely, $\mu^{(i)}$ and $\sigma^{(i)}$. The latent vectors $\mathbf{z}^{(i)}$ can then be obtained by randomly sampling from the distribution. And thus, multiple $\mathbf{z}^{(i)}$ ’s may correspond to a single $\mathbf{x}^{(i)}$. This method is depicted in figure 2.4. Finally, a sample $\mathbf{z}^{(i)}$ can be obtained as follows:

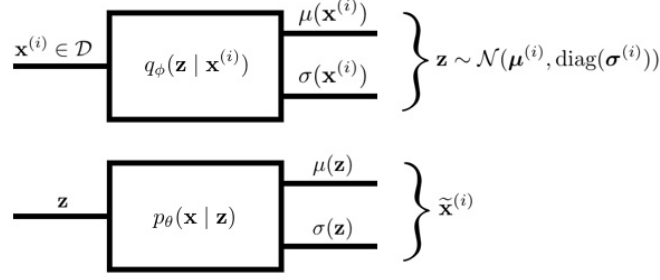


Figure 2.4: High level view of a variational autoencoder, depicting how a data point $\mathbf{x}^{(i)}$ is encoded into a latent distribution and reconstructed as $\tilde{\mathbf{x}}^{(i)}$. Both blocks depict a neural network. The upper block is the encoder and the lower block the decoder. The upper block receives a data points \mathbf{x} and produces the parameters of $q_\phi(\mathbf{z} | \mathbf{x}^{(i)})$. Since we choose to model $q_\phi(\mathbf{z} | \mathbf{x}^{(i)})$ as a Gaussian with independent components, the covariance matrix Σ is zero everywhere except for the diagonal. This way the diagonal values, representing the standard deviations, can be represented via a single vector $\sigma^{(i)}$. The vectors $\mu^{(i)}$ and $\sigma^{(i)}$ are $\mu(\mathbf{x}^{(i)})$ and $\sigma(\mathbf{x}^{(i)})$, respectively. These are the output of the encoder block and form the parameters for $q_\phi(\mathbf{z} | \mathbf{x}^{(i)})$. A single neural network with parameter weights ϕ is used to simulate $q_\phi(\mathbf{z} | \mathbf{x}^{(i)})$ for every $\mathbf{x}^{(i)} \in \mathcal{D}$. This strategy of sharing ϕ across data points is referred to as "amortised variational inference" [12].

$$\mathbf{z}^{(i)} = \mu^{(i)} + \sigma^{(i)} \odot \epsilon^{(i)} \quad (2.11)$$

where $\epsilon^{(i)}$ corresponds to a sampled value $\epsilon^{(i)} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ and \odot is element-wise multiplication. Computing $\mathbf{z}^{(i)}$ through $\epsilon^{(i)}$, rather than directly sampling from $\mathbf{z}^{(i)} \sim \mathcal{N}(\mu^{(i)}, \text{diag}(\sigma^{(i)}))$ is referred to as the parametrisation trick and allows for gradients to freely backpropagate through the layer [10].

The learning objective

So far we have discussed how predictions of a neural network can emulate predicting a distribution $q_\phi(\mathbf{z} | \mathbf{x}^{(i)})$, by predicting the distribution's parameters $\mu^{(i)}$ and $\sigma^{(i)}$. However, no constraints have been set on the quality of the distributions. We will discuss this now. The representations are optimised to minimise two measurements: one, the reconstruction error, and secondly, the distance from the latent distributions to $\mathcal{N}(\mathbf{0}, \mathbf{I})$.

The loss function to be optimised for a single data point $\mathbf{x}^{(i)}$ is shown in the equation below.

$$\mathcal{L}_{\theta, \phi}(\mathbf{x}^{(i)}) = \mathbb{E}_{\mathbf{z} \sim q_\phi(\cdot | \mathbf{x}^{(i)})} \left[-\log p_\theta(\mathbf{x}^{(i)} | \mathbf{z}) \right] + D_{KL} \left(q_\phi(\cdot | \mathbf{x}^{(i)}) || p_\theta(\cdot) \right) \quad (2.12)$$

Although the equation may seem daunting at first, we will decompose its components. The loss function is made up of two terms, the left term corresponds to the reconstruction error, while the second term poses constraints on the latent space. Let us focus on the reconstruction term first.

$$\mathbb{E}_{\mathbf{z} \sim q_\phi(\cdot | \mathbf{x}^{(i)})} \left[-\log p_\theta(\mathbf{x}^{(i)} | \mathbf{z}) \right]$$

The distribution $p_\theta(\mathbf{x}^{(i)} | \mathbf{z})$ is a distribution over $\mathbf{x}^{(i)}$ where \mathbf{z} is instantiated. What comes out are thus probabilities. Notice \mathbf{z} is sampled from the Gaussian distribution $\mathbf{z} \sim q_\phi(\mathbf{z} | \mathbf{x}^{(i)})$.

So the \mathbf{z} 's close to the mean $\mu(\mathbf{x}^{(i)}) = \boldsymbol{\mu}^{(i)}$ are more likely to be sampled than those further away. So for these \mathbf{z} 's, we'd like the probability of corresponding to the actual $\mathbf{x}^{(i)}$ to be high. Finally, by adding a negative sign in front, maximising this probability is equivalent to minimising the negative probability [?]. In practice this term is approximated through mini-batches with the mean squared error.

Optimising the second term of equation 2.12 poses the constraints on the latent distributions $q_\phi(\mathbf{z} | \mathbf{x}^{(i)})$. Again, this metric should be minimised. As we discussed in the chapter on Entropy, the KL divergence can be considered as a (non-symmetric) distance measure between two distributions. Hence, this value is small when the two distributions are similar. As we discussed earlier $p_\theta(\mathbf{Z})$ is often replaced by $\mathcal{N}(\mathbf{0}, \mathbf{I})$ as shown below. Minimising this equation will result in moving each distribution $q_\phi(\mathbf{z} | \mathbf{x}^{(i)})$, corresponding to a value $\mathbf{x}^{(i)}$, close to $\mathcal{N}(\mathbf{0}, \mathbf{I})$.

$$D_{KL} \left(q_\phi(\cdot | \mathbf{x}^{(i)}) \parallel \mathcal{N}(\mathbf{0}, \mathbf{I}) \right)$$

One can algebraically prove, that optimising this equation is equivalent to optimising the following equation [11]:

$$D_{KL} \left(\mathcal{N}(\boldsymbol{\mu}^{(i)}, \text{diag}(\boldsymbol{\sigma}^{(i)})) \parallel \mathcal{N}(\mathbf{0}, \mathbf{I}) \right) = \frac{1}{2} \sum_{k=1}^D \left(-\log(\sigma_k^{(i)})^2 - 1 + (\sigma_k^{(i)})^2 + (\mu_k^{(i)})^2 \right) \quad (2.13)$$

where $\sigma_k^{(i)}$ and $\mu_k^{(i)}$ correspond to the components of the predicted vectors $\boldsymbol{\sigma}^{(i)}$ and $\boldsymbol{\mu}^{(i)}$, respectively.

[TODO] Relation with variational inference / derivation of ELBO loss

$$\begin{aligned} -\mathcal{L}_{ELBO} &= - \mathbb{E}_{\mathbf{z} \sim q_\phi(\cdot | \mathbf{x}^{(i)})} \left[\log p(\mathbf{Z}, \mathbf{x}^{(i)}) - \log q(\mathbf{z} | \mathbf{x}^{(i)}) \right] \\ &= \mathbb{E}_{\mathbf{z} \sim q_\phi(\cdot | \mathbf{x}^{(i)})} \left[-\log p_\theta(\mathbf{x}^{(i)} | \mathbf{z}) \right] + D_{KL} \left(q_\phi(\cdot | \mathbf{x}^{(i)}) \parallel p_\theta(\cdot) \right) \end{aligned} \quad (2.14)$$

Where the first term is reconstruction error and second is regularisation. Hence minimising the KL divergence between the two posteriors is equivalent to minimising the divergence between the approximate posterior $q(\mathbf{z} | \mathbf{x}^{(i)})$ and the **marginal or prior?** $p_\theta(\mathbf{Z})$.

$p(z)$ in the equation is usually chosen to be the standard normal distribution $\mathcal{N}(0, I)$, such that a closed form solution exists. As such, when $p(z)$ corresponds to the standard normal, and $q(\mathbf{z} | \mathbf{x}^{(i)})$ is a multidimensional Gaussian with mean vector $\boldsymbol{\mu}^{(i)}$ and covariance matrix with independent dimensions, such that the diagonal corresponds of a vector standard deviations $\boldsymbol{\sigma}^{(i)}$, then the KL divergence has the following closed form:

$$D_{KL} \left(\mathcal{N}(\boldsymbol{\mu}^{(i)}, \text{diag}(\boldsymbol{\sigma}^{(i)})) \parallel \mathcal{N}(\mathbf{0}, \mathbf{I}) \right) = \frac{1}{2} \sum_{k=1}^D \left(-\log(\sigma_k^{(i)})^2 - 1 + (\sigma_k^{(i)})^2 + (\mu_k^{(i)})^2 \right) \quad (2.15)$$

For the equation above, gradients can easily be back propagated through machine learning libraries such as Tensor Flow. We still require a method for computing the gradient of the first term in \mathcal{L}_{ELBO} .

2.4 Representation learning through Noise-contrastive estimation

2.4.1 Contrastive predictive coding

In what follows next, we discuss Contrastive Predictive Coding (CPC), a representation learning approach that we use as the basis for our own experiment in the following chapter. CPC is an unsupervised learning approach, again with the objective of learning (lower dimensional) representations from high dimensional data [14]. While the objective is thus the same as for the autoencoders discussed in the previous section, CPC achieves its representations entirely differently. An autoencoder’s objective is to define a compressed representation from which the original data can be recovered. However, when working with sequential data, compressing patches of the sequence without considering the relation with nearby patches, will result in lost information as the context between patches is not encoded into the representation.

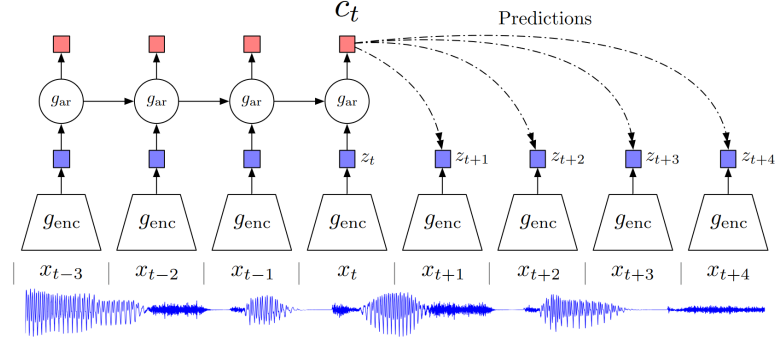


Figure 2.5: Overview of Contrastive Predictive Coding, originates from [14]

CPC deals with these context issues by maximising the shared information between the extracted representations of temporally nearby patches [15]. We will discuss this concept in more detail in a following section. For now, we would like to draw the readers attention to figure 2.5. The figure depicts a high level view of how this idea is achieved by producing two latent representations \mathbf{z}_i and \mathbf{c}_i . An audio sequence of undefined length is split up into patches $\mathbf{x}_1 \dots \mathbf{x}_n$ where each \mathbf{x}_i is a vector of fixed length, containing for instance 10ms of speech audio. Each patch \mathbf{x}_i is encoded into latent representation \mathbf{z}_t , defined as follows:

$$\mathbf{z}_t = g_{enc}(\mathbf{x}_t).$$

$g_{enc}(\cdot)$ is for instance a convolutional fully connected ANN. The latent representations $\mathbf{z}_1 \dots \mathbf{z}_n$ are obtained independently from each other and do not yet contain any information on an historic context. This historic context is achieved through $g_{ar}(\cdot)$, an auto-regressor, which encodes all previous $\mathbf{z}_1 \dots \mathbf{z}_t$ into a single representation \mathbf{c}_t :

$$\mathbf{c}_t = g_{ar}(\mathbf{z}_1 \dots \mathbf{z}_t)$$

Either \mathbf{z}_t or \mathbf{c}_t could be used as latent representation for downstream tasks. Oord et al. suggest to use \mathbf{c}_t for tasks where context about the past is useful, for instance speech recognition, and \mathbf{z}_t when historic context is not useful [14]. As shown in figure 2.5, the encodings from sequential data of undefined length, may correspond a series of latent representations $\mathbf{c}_1, \mathbf{c}_2, \dots$

or $\mathbf{z}_1, \mathbf{z}_2, \dots$. In the case of downstream tasks which require a single representation vector, Oord et al. propose to pool the sequence of vector representations into a single vector.

Slowly varying features

The temporally nearby patches \mathbf{z}_{t+1} and \mathbf{c}_t are optimised to preserve shared information, while discarding differences. Before we discuss how to obtain such representations, we first motivate why defining representations in this fashion makes sense.

Consider we would like to define useful representations for sequential data such as speech signals. Then it not unlikely to believe that the conveyed information at time step t and $t+k$ contains some redundant information, such as pitch, frequency, tone, etc. [16]. Meanwhile, large changes of the signal in a small time window, may be the result of noise. Sequential data which poses these slowly varying features, are commonly referred to as "slow-features" [17]. CPC leverages these slowly varying features, by encoding the underlying shared information between different patches, while at the same time discarding low-level information and noise that is more local [14].

The learning objective

CPC will learn to preserve information between temporally nearby representations, by solving another task. In particular, CPC learns to discriminate subsequent *positive* samples \mathbf{z}_{t+k} from *negative* random samples \mathbf{z}_j . This is achieved through a similarity function $f_k(\cdot)$, which scores the similarity between two latent representations [15]. It is defined as a log bilinear model as follows:

$$f_k(\mathbf{z}_j, \mathbf{c}_t) = \exp(\mathbf{z}_j^T W_k \mathbf{c}_t) \quad (2.16)$$

where W_k is a weight matrix which is learned. $f_k(\mathbf{z}_j, \mathbf{c}_t)$ thus quantifies how likely the context \mathbf{c}_t corresponds to a random vector \mathbf{z}_j . Due to the slowly varying data assumption, a good representation for successive representations \mathbf{z}_{t+k} and \mathbf{c}_t is one where $f_k(\mathbf{z}_{t+1}, \mathbf{c}_{t+1})$ is high and $f_k(\mathbf{z}_j, \mathbf{c}_t)$ is small for random \mathbf{z}_j . Or equivalently, maximising the shared information between temporally nearby patches, while discarding the temporal noise results in large values $f_k(\mathbf{z}_{t+1}, \mathbf{c}_{t+1})$.

The InfoNCE loss, used to optimise g_{enc} , g_{ar} and W_k simultaneously is shown below.

$$\mathcal{L}_{\text{NCE}} = - \sum_k \mathbb{E}_X \left[\log \frac{f_k(\mathbf{z}_{t+k}, \mathbf{c}_t)}{\sum_{\mathbf{z}_j \in X} f_k(\mathbf{z}_j, \mathbf{c}_t)} \right] \quad (2.17)$$

where X corresponds to the set $\{\mathbf{z}_{t+k}, \mathbf{z}_1, \mathbf{z}_2, \dots\}$. Notice that there exists exactly one $\mathbf{z}_{t+k} \in X$, which corresponds to a positive sample and all other $\mathbf{z}_j \in X$ correspond to negative samples (with respect to the context \mathbf{c}_t). Hence good representations for \mathbf{z}, \mathbf{c} and $f_k(\cdot)$, will result in a large similarity score for positive samples and approximate 0 for negative samples, resulting in a minimum fraction equal to 1. This would then be cancelled out by the $\log(\cdot)$ function. Meanwhile, \mathcal{L}_n is large when the denominator is large, indicating in a large high similarity score for negative samples.

Ties with mutual information

Earlier we argued that CPC's encodings will preserve shared information between temporally nearby patches, while discarding the local noise. Oord et al. make this claim even stronger by making ties with mutual information, which we discussed in a previous chapter. In particular, Oord et al. proof that optimising InfoNCE is equivalent to maximising the mutual information between \mathbf{c}_t and \mathbf{z}_{t+1} [14].

$$I(\mathbf{z}_{t+1}; \mathbf{c}_t) = \sum_{\mathbf{z}_{t+1}, \mathbf{c}_t} p(\mathbf{z}_{t+1}, \mathbf{c}_t) \log \frac{p(\mathbf{z}_{t+1} | \mathbf{c}_{t+1})}{p(\mathbf{z}_{t+1})} \quad (2.18)$$

This proof is available in their appendix. Although, we do not repeat the proof here, we give a high level overview.

The first step in proving the relation between the InfoNCE loss and mutual information is to model $f_k(\mathbf{z}_{t+k}, \mathbf{c}_t)$ in a probabilistic manner. The InfoNCE loss is in fact the categorical cross-entropy of classifying the positive sample correctly with $\frac{f_k}{\sum_x f_k}$ as the predicted model [14]. Since this equation may take values between zero and one, it can be considered as a probability. In particular, the optimal probability for the loss can then be written as

$$p(i | X, \mathbf{c}_t)$$

where X corresponds the set of samples $\{\mathbf{z}_{t+k}, \mathbf{z}_1, \mathbf{z}_2, \dots\}$ as discussed in the InfoNCE loss, and i corresponds to indicator that sample \mathbf{z}_i is the "positive" sample. By doing so, one can eventually obtain a proportionality relation to the density distribution presented below.

$$f_k(\mathbf{z}_{t+k}, \mathbf{c}_t) \propto \frac{p(\mathbf{z}_{t+k} | \mathbf{c}_t)}{p(\mathbf{z}_{t+k})} \quad (2.19)$$

Oord et al. utilise this proportionality relation to reformulate $-\mathcal{L}_{\text{NCE}}$ as a lower bound on the mutual information between \mathbf{z}_{t+1} and \mathbf{c}_t as follows [15, 14]:

$$I(\mathbf{z}_{t+1}; \mathbf{c}_t) \geq \log(N) - \mathcal{L}_{\text{NCE}} \quad (2.20)$$

Since the number of samples N is a constant, the mutual information between \mathbf{z}_{t+1} and \mathbf{c}_t becomes greater when \mathcal{L}_{NCE} becomes smaller. Additionally, when the number of samples N increases, the bound becomes tighter.

2.4.2 Greedy InfoMax

So far we discussed how CPC encodes latent representations by maximising the mutual information between temporally nearby patches of data. This method has shown great success in recent years and is considered state-of-the art in self-supervised learning for encoding sequential data [18]. Additionally, CPC has been successfully applied to multiple use cases [18, 19, 20, 21, 22, 23]. This is achieved by minimising the InfoNCE loss discussed earlier in equation 2.17. Through this *global* loss function all parameters are optimised end-to-end via backpropagation.

Although, backpropagation is empirically shown to be highly effective [24, 25], it still suffers from multiple constraints. Firstly, there is a biological perspective to consider, as the human brain lacks a global objective function that can be optimised by backpropagating an error signal [26]. This is especially significant when considering how children can learn to categorise from a few examples, whereas end-to-end backpropagation often requires extensive datasets to achieve good generalisation [15]. Moreover, end-to-end backpropagation also suffers from computational constraints, such as requiring the entire computational graph, including all parameters, activations and gradients to fit in memory [15]. Additionally, during the training process of a neural network, each layer has to wait for the gradients of its subsequent layer, which reduces locality and impedes the efficiency of hardware accelerator design [15].

Towards greedy learning

To overcome these biological and computational constraints, Löwe et al. introduce Greedy InfoMax (GIM) [15], an extension on CPC. Whereas CPC obtains representations \mathbf{z}_t and \mathbf{c}_t through encoder $g_{enc}(\cdot)$ and autoregressor $g_{ar}(\cdot)$, Löwe et al. split up $g_{enc}(\cdot)$'s neural network architecture by depth into M so called "modules":

$$g_{enc}^1(\cdot), g_{enc}^2(\cdot), \dots, g_{enc}^M(\cdot)$$

Each module's output is the input of the successive module. A single module may for instance represent one or more layers.

$$\begin{aligned} g_{enc}^m(\mathbf{z}_t^{m-1}) &= \mathbf{z}_t^m \\ g_{ar}(\mathbf{z}_1^M \dots \mathbf{z}_t^M) &= \mathbf{c}_t \end{aligned}$$

The final representation \mathbf{z}_t^M is obtained by propagating \mathbf{x}_t through each module as follows:

$$g_{enc}^M(\dots g_{enc}^2(g_{enc}^1(\mathbf{x}_t))) = \mathbf{z}_t^M$$

Both \mathbf{z}_t^M or \mathbf{c}_t may serve as representation for following downstream tasks.

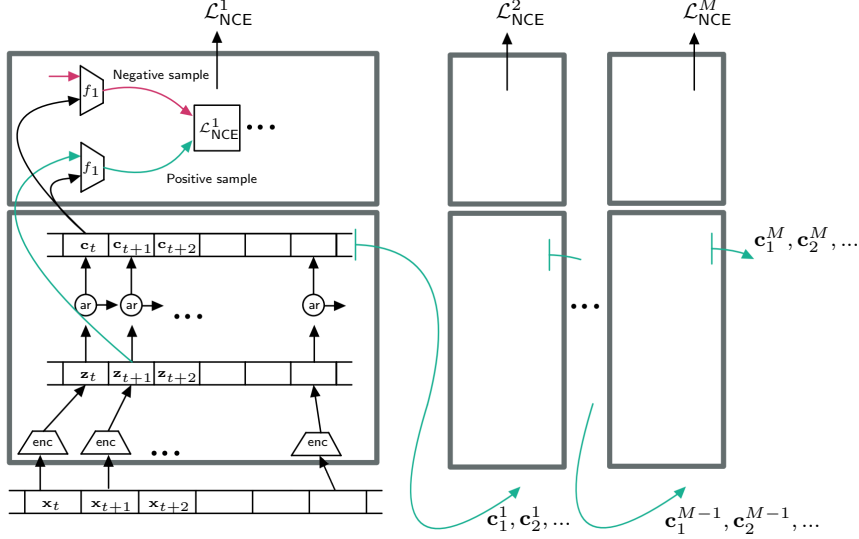


Figure 2.6: Overview of Greedy InfoMax

Each module $g_{enc}^m(\cdot)$ is *greedily* optimised with an adaptation of the InfoNCE Loss. This idea is depicted in figure 2.6, which displays M modules, each trained with their own instance of the InfoNCE loss.

In contrast to CPC where a single loss function \mathcal{L}_{NCE} and scoring function $f_k(\cdot)$ are required, we now require \mathcal{L}_{NCE}^m and $f_k^m(\cdot)$ for each module. The equations to optimise for $g_{enc}^m(\cdot)$ then become:

$$f_k^m(\mathbf{z}_{t+k}^m, \mathbf{z}_t^m) = \exp(\mathbf{z}_{t+k}^{mT} \mathbf{W}_k^m \mathbf{z}_t^m) \quad (2.21)$$

$$\mathcal{L}_{NCE}^m = - \sum_k \mathbb{E}_X \left[\log \frac{f_k^m(\mathbf{z}_{t+k}^m, \mathbf{z}_t^m)}{\sum_{\mathbf{z}_j^m \in X} f_k^m(\mathbf{z}_j^m, \mathbf{z}_t^m)} \right] \quad (2.22)$$

From equation 2.21 we observe, $f_k^m(\cdot)$ no longer receives as input an aggregate \mathbf{c}_t , but instead a second \mathbf{z}_t . This is in contrast to CPC where $f_k^m(\mathbf{z}_{t+k}^m, \mathbf{c}_t^m)$ is maximised instead. Löwe et al. thus omit the autoregressive function within each module. Optionally, for downstream tasks where broad context is helpful, such as classification of phonetic structures in speech recognition, an autoregressive model $g_{ar}(\mathbf{z}_1^M \dots \mathbf{z}_t^M) = \mathbf{c}_t$ can be appended as a final $M+1$ 'th module. The altered scoring function then becomes:

$$f_k^{M+1}(\mathbf{z}_{t+k}^M, \mathbf{c}_t) = \exp\left(\mathbf{z}_{t+k}^{M^T} W_k^{M+1} \mathbf{c}_t\right)$$

As a result of this approach, the mutual information $I(\mathbf{z}_{t+k}^m, \mathbf{z}_t^m)$ is maximised, rather than $I(\mathbf{z}_{t+k}, \mathbf{c}_t)$ as in CPC.

As a result of GIM's greedy approach, modules can be trained in parallel, but also sequentially. By training one module after the other, the memory cost can be decreased during training which is beneficial for memory-constrained scenarios. Furthermore, when a higher level of abstraction is needed, additional modules can be appended to the architecture later on without having to fine-tune previous modules.

Chapter 3

Variational InfoNCE

3.1 Motivation

In the previous section we discussed two categories of representation learning through deep learning. First, we discussed the autoencoder and its variational counterpart, which minimise the reconstruction error. Secondly, we discussed Contrastive Predictive Coding and Greedy InfoMax, both of which optimise the Info NCE objective. This category seeks to maximise the mutual information between the encodings of data patches that are temporally nearby. The latent representations obtained from all four methods can then be utilised for downstream tasks [3, 28, 14, 15]

The autoencoder’s sole objective is to define representations to reconstruct the original data. As a result, the representations may serve well for data compression, however, no additional constraints are enforced, such as feature disentanglement and thus the latent space may still be hard to work with for downstream tasks [29]. Meanwhile, VAEs’ additional regularisation term, results in representations which break down or disentangle each feature into a narrowly defined variable and encodes them as separate dimensions [28]. This additional constraint may result in better suited representations for downstream tasks.

Both autoencoders and VAEs merely learn to reconstruct the data. Hence, all the “information” that is important to reconstruct the data will be maintained in the latent representation, whether the information is useful for the downstream task or not. Meanwhile, optimising latent representations for the InfoNCE objective will maintain shared information between temporally nearby patches, while discarding local noise. Reconstruction is thus not needed for training. This strategy has the tremendous benefit that a decoder block is not required, resulting in a significantly simplified architecture, meanwhile maintaining state-of-the-art performance [18]. A second benefit of these mutual information maximisation models is that they are directly compatible with sequential data.

Both categories (reconstruction and information maximisation algorithms) possess the ability to obtain useful representations for various downstream tasks. However, the content of these representations may not always be intuitive to humans and their structure may be difficult to comprehend. While CPC and GIM are considered state-of-the-art, their performance comes at a cost of having the least interpretable representations. Autoencoders maintain interpretability by using a decoder to reveal the information contained in the latent representation. The same transparency can also be achieved with VAEs. Additionally, by using a standard Gaussian as a prior and constraining the latent distributions to be similar to this prior, we can interpolate between representations and observe the effects through the decoder. As such, we can observe

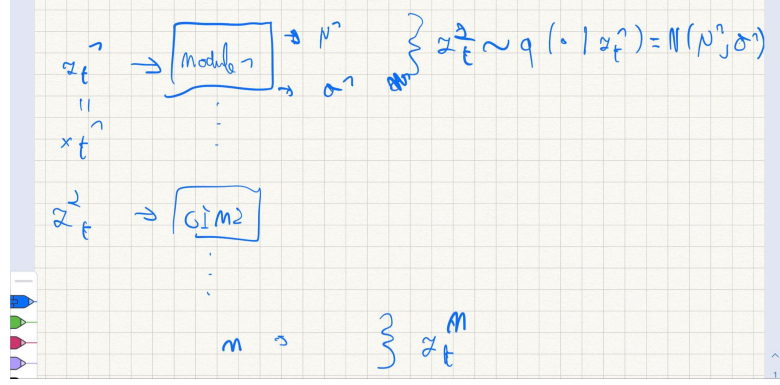


Figure 3.1:

the specific information that is contained in each of the representation’s features. VAEs can also result in disentangled features, further enhancing interpretability [30]. In contrast, CPC and GIM do not contain a built in decoder mechanism, nor pose constraints on the latent space, significantly reducing interpretability.

3.2 Towards decoupled training for probabilistic representations

In what follows next we introduce ”Variational Greedy Infomax”, maintaining the state-of-the-art performance obtained from optimising InfoNCE, while leveraging the interpretable and disentangled benefits from VAEs. This is achieved by optimising a novel loss function, *Variational-InfoNCE*, a combination of InfoNCE and the regularisation term from VAEs. Additionally, by splitting up the neural network into modules, as introduced in [15], we greedily optimise each module with its own instance of this loss function. As a result, the interpretability benefits from VAEs will also be applicable in-between modules. This is in contrast to VAEs where solely the final output representations are interpretable.

As discussed in the section on Contrastive Predictive Coding (CPC), a patch of sequential data \mathbf{x}_t is encoded through $g_{enc}(\mathbf{x}_t) = \mathbf{z}_t$ and aggregated over previous encodings through auto-regressor $g_{ar}(\mathbf{z}_1 \dots \mathbf{z}_t) = \mathbf{c}_t$, where both \mathbf{z}_t or \mathbf{c}_t may serve as representations for downstream tasks. The encoder function $g_{enc}(\cdot)$ is represented as neural network, eg via a CNN, and $g_{ar}(\cdot)$ for instance as a GRU. Finally, the encoding functions $g_{enc}(\cdot)$ and $g_{ar}(\cdot)$ are obtained by optimising a global loss function, the InfoNCE loss, end-to-end via backpropagation.

Instead, in this study, we split up $g_{enc}(\cdot)$ ’s network architecture by depth into M modules

$$g_{enc}^1(\cdot), g_{enc}^2(\cdot), \dots, g_{enc}^M(\cdot)$$

and prevent gradients from flowing between modules, as introduced in [15]. An additional optional $M + 1$ ’th module $g_{ar}(\cdot)$ can be appended to the architecture. Each module is greedily optimised via a novel loss function, \mathcal{L}_{V-NCE} , which we will define in a following subsection. Each module’s output serves as input for the successive module, as presented in the following equations, and depicted in figure 3.1.

$$\begin{aligned}
g_{enc}^1(\mathbf{x}_t) &= \mathbf{z}_t^1 \\
g_{enc}^m(\mathbf{z}_t^{m-1}) &= \mathbf{z}_t^m \\
g_{ar}(\mathbf{z}_1^M \dots \mathbf{z}_t^M) &= \mathbf{c}_t
\end{aligned}$$

The final representation \mathbf{c}_t is obtained by propagating \mathbf{x}_t through each modules as follows:

$$g_{ar}(g_{enc}^M(\dots g_{enc}^2(g_{enc}^1(\mathbf{x}_t))))$$

Additionally, taking inspiration from VAEs, the outputs from $g_{enc}^m(\cdot)$ and $g_{ar}(\cdot)$ are in fact samples from a distribution denoted by $q(\mathbf{z}_t^m | \mathbf{z}_t^{m-1})$, defined as a multivariate Gaussian with diagonal covariance matrix, as follows:

$$q(\cdot | \mathbf{z}_t^{m-1}) = \mathcal{N}(\boldsymbol{\mu}, \text{diag}(\boldsymbol{\sigma}^2))$$

with $\boldsymbol{\mu}$ and $\boldsymbol{\sigma}$ dependent on \mathbf{z}_t^{m-1} , specified in more detail in a following subsection. The outputs for $g_{enc}^m(\cdot)$ and $g_{ar}(\cdot)$ are obtained by sampling from this distribution, denoted respectively, as follows:

$$\mathbf{z}_t^m \sim q(\cdot | \mathbf{z}_t^{m-1}) \quad (3.1)$$

$$\mathbf{c}_t \sim q(\cdot | \mathbf{z}_t^M) \quad (3.2)$$

Modules are thus stochastic and computing $g_{enc}^m(\mathbf{z}_t^{m-1})$ twice will likely result in two different representations of \mathbf{z}_t^m . This is in contrast to CPC and GIM's encodings which remain fixed depending to the input [14, 15].

We achieve these stochastic modules by defining each module $g_{enc}^m(\cdot)$ consisting of two blocks. The first block receives as input \mathbf{z}_t^{m-1} and predicts the parameters $\boldsymbol{\mu}$ and $\boldsymbol{\sigma}$. These two parameters describe the distribution $q(\cdot | \mathbf{z}_t^{m-1})$. Since we defined q as Gaussian with a diagonal covariance matrix, the distribution can be fully described by those two vectors. The second block samples $\mathbf{z}_t^m \sim q(\cdot | \mathbf{z}_t^{m-1})$ from this distribution and produces an output representation. This is depicted in figure 3.2.

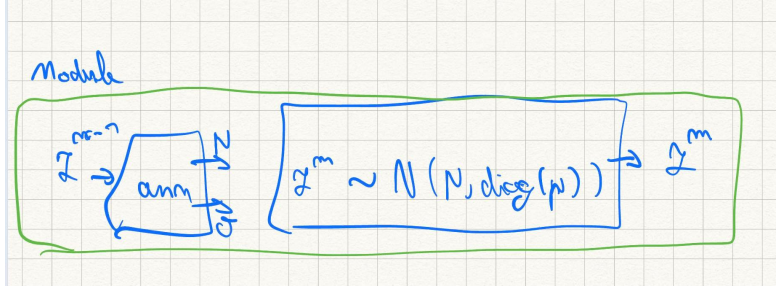


Figure 3.2:

In practice, sampling from q is achieved through a reparametrisation trick, as introduced in [11]. The equation to compute \mathbf{z}_t^m then becomes:

$$\mathbf{z}_t = \boldsymbol{\mu} + \boldsymbol{\sigma} \odot \boldsymbol{\epsilon}$$

where ϵ corresponds to a sampled value $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ and \odot is element-wise multiplication. The procedure to obtain \mathbf{c}_t is analogous to \mathbf{z}_t^m which we described above.

Because of this probabilistic approach, a single patch of data \mathbf{x}_t will have multiple representations \mathbf{z}_t^M , providing increased variance in the representations. This can potentially benefit downstream tasks, particularly when labelled data is scarce [28], leading to improved performance.

3.3 The learning objective

Instead of training the neural network end-to-end with a global loss function, the network is split up into modules, which each are optimised greedily with their own personal loss function. Through the introduction of the novel *Variational-InfoNCE* loss, mutual information between temporally nearby representations is maximised, while regularising the latent space to be approximate to the standard Gaussian $\mathcal{N}(\mathbf{0}, \mathbf{I})$. The Variational-InfoNCE loss is defined as follows:

$$\mathcal{L}_{\text{V-NCE}}^m = \underbrace{\sum_k \mathbb{E}_{\substack{\mathbf{z}_{t+k}^m \sim q(\cdot | \mathbf{z}_{t+k}^{m-1}) \\ \mathbf{z}_t^m \sim q(q(\cdot | \mathbf{z}_t^{m-1}))}} \left[\log \frac{f_k^m(\mathbf{z}_{t+k}^m, \mathbf{z}_t^m)}{\sum_{\mathbf{z}_j^m \in X} f_k^m(\mathbf{z}_j^m, \mathbf{z}_t^m)} \right]}_{\text{Maximise } I(\mathbf{z}_{t+k}^m, \mathbf{z}_t^m)} + \underbrace{\beta D_{KL}(q(\cdot | \mathbf{z}_t^{m-1}) || \mathcal{N}(\mathbf{0}, \mathbf{I}))}_{\text{Regularisation}} \quad (3.3)$$

$m \in \mathbb{N}$ refers to the m 'th module. $k \in \mathbb{N}$ corresponds to the number of patches in the future the similarity score $f_k^m(\mathbf{z}_{t+k}^m, \mathbf{z}_t^m)$ must rate. \mathbf{z}_{t+k}^m and \mathbf{z}_t^m are encoded samples produced by $g_{enc}^m(\mathbf{z}_{t+k}^{m-1})$ and $g_{enc}^m(\mathbf{z}_t^{m-1})$, respectively. X is a set of samples $\{\mathbf{z}_{t+k}^m, \mathbf{z}_1^m, \mathbf{z}_2^m, \dots\}$ where $\mathbf{z}_j^m \neq \mathbf{z}_{t+k}^m$ are random samples.

The similarity score $f_k^m(\cdot)$'s definition is identical to [15]:

$$f_k^m(\mathbf{z}_{t+k}^m, \mathbf{z}_t^m) = \exp(\mathbf{z}_{t+k}^{mT} W_k^m \mathbf{z}_t^m)$$

3.3.1 Gradient

The gradient of the first term in $\mathcal{L}_{\text{V-NCE}}^m$ can be approximated through mini-batches, and optimised directly in PyTorch. With regards to the second term, since $q(\cdot | \mathbf{z}_t^{m-1})$ is a Gaussian, a closed form solution exists [11], the term can be differentiated without approximated method.

$$\frac{1}{2} \sum_{j=1}^J \left(1 + \log((\sigma_j^{(i)})^2) - (\mu_j^{(i)})^2 - (\sigma_j^{(i)})^2 \right) \quad (3.4)$$

$$D_{KL}(\mathcal{N}(\boldsymbol{\mu}^{(i)}, \text{diag}(\boldsymbol{\sigma}^{(i)})) || \mathcal{N}(\mathbf{0}, \mathbf{I})) = \frac{1}{2} \sum_{k=1}^D \left(-\log(\sigma_k^{(i)})^2 - 1 + (\sigma_k^{(i)})^2 + (\mu_k^{(i)})^2 \right) \quad (3.5)$$

where $z^{(i,l)} = \sigma^{(i)} \odot \epsilon^{(l)}$ and $\epsilon^{(l)} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$

todo: variables should maybe be bold.

3.4 Practical benefits

3.4.1 Latent representations

Representations: Interpretability of latent representations: - Interpolation – ζ adds decoder

Better generalisation for downstream tasks - Overfitting: reduction of required labelled data needed. Similar data is similar region, the kl divergence makes regions bigger.

Overfitting during inference: - The same datapoint has multiple (similar) representations, such that learning techniques for downstream tasks will not be able to "memorise" the latent space as easily.

- Holes: more predictable inference, such that unseen data is more likely to be near clusters. And thus downstream tasks receive latents that are more similar to what is seen before. = better generalisation

3.4.2 Training

- Batch normalisation: is useful for following modules, BUT ALSO DOWNSTREAM TASKS!
 - GIM advantages remain: maintain the benefits such as smaller networks that can learn in-dep, ALSO NO EXPLICIT NEED FOR A DECODER. THIS REMAINS, SO ARCH CAN BE SIMPLER.

- Independent latent dimensions - During training similar behaviour to batch normalization in-between layers

Other sources:

Chapter 4

Experiments

4.1 Greedy Infomax model

- train CNN
- architecture: 4 layers ...
- dataset

4.1.1 Dataset

The Greedy Infomax model is trained on speech data. The model takes as input a raw speech signal of a fixed length and outputs a latent representation for that signal. The dataset is split up into 729 training files and 122 test files. In each file consists of a single spoken sound consisting of three consonants and three vowels, where the consonants and vowels alternate each other. Some Examples are the sounds "gi-ga-bu" and "ba-bi-gu". All the sounds are spoken by the same person, at a constant and peaceful **todo: describe emotional aspects of speech audio**.

The following transformations are applied to the audio files. Although the original contains a sample rate of 441 Khz, the audio files are downsampled to 16 Khz, matching the sample rate used by Löwe [15]. This significantly reduces the size of the latent representations, and thus the required amount of VRAM during training. Additionally, two types of noise are added to the data. We apply Gaussian white noise, at different decibels ranging between zero and fifteen **TODO**

- also background noise from dataset. is a way to enlargen our dataset.
- Each audio file is cut to have length **10240**. Additionally,

Since the recordings are very consistent in loudness and are noise free, we can split up the files per syllable, obtaining three files per original sound (one for each syllable).

A sliding window is used of size 0.02 seconds. With a sample rate of 22050, this corresponds to roughly 500 samples per window. The maximum is computed for each window. Speech signals can then be split up when a severe dip happens in the signal. Regions where the amplitude is greater than 0.2 are considered **klinkers**, the regions with with lower values are considered **medeklinkers**. Apart from a few edge cases, this technique worked well enough for this purpose. In those cases, the splitting points closest to the one-third and two-third splitting points were considered.

ERR: OUDE AFBEELDINGEN ZIJN WEG

note: we do need a hard threshold which is based on the signal's intensity level. One could consider the alternative approach of looking at the gradient at each point and selecting the points with largest negative gradient. This will work in many cases, however, not for temporal envelopes which gradually move towards zero, s.t: 4.1. Instead we use a dynamic threshold. This threshold is computed by creating transforming the signal into bins of 90'th percentile, creating a histogram of the single signal and applying otsu's image segmentation algorithm to obtain the threshold of that single audio sample. We also tried directly applying otsu to the moving average and maximum of the bins. This either gave a threshold that was too small or too large. the 90th percent resulted in an acceptable compromise.

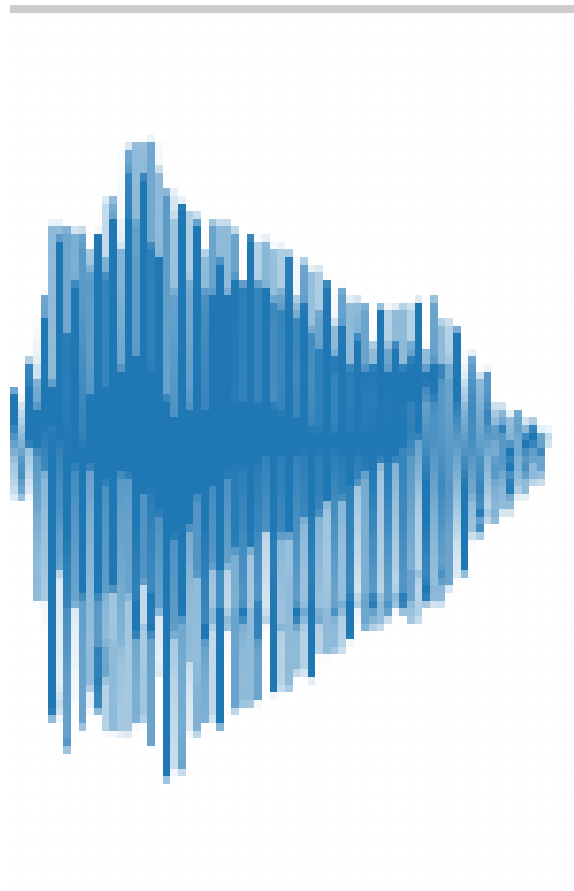


Figure 4.1:

Example where the explained strategy does not work:

Reference images for in the text:

audio padded to maximum length. (added zeros in front and back)

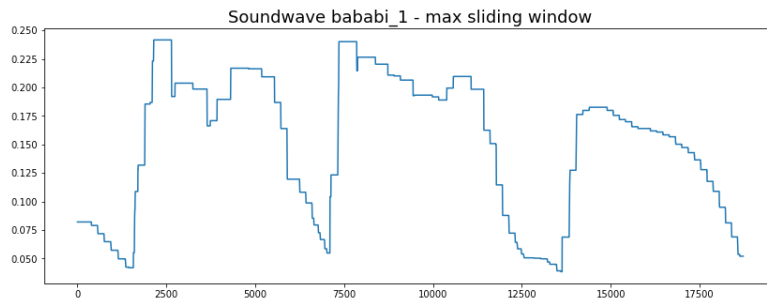


Figure 4.2:

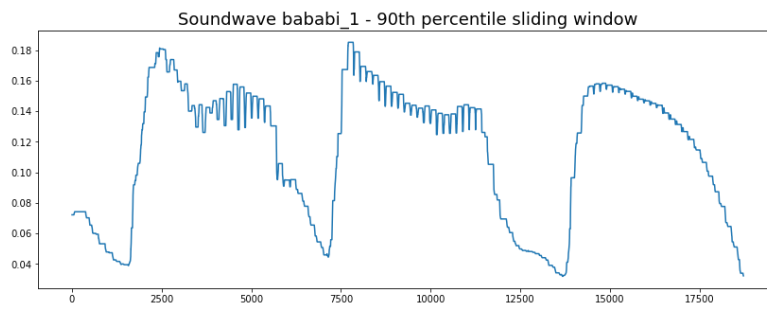


Figure 4.3:

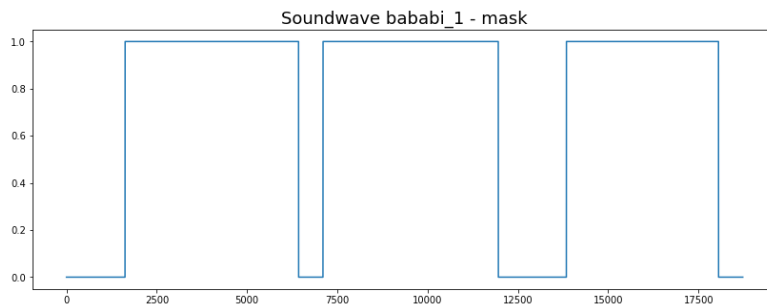


Figure 4.4:

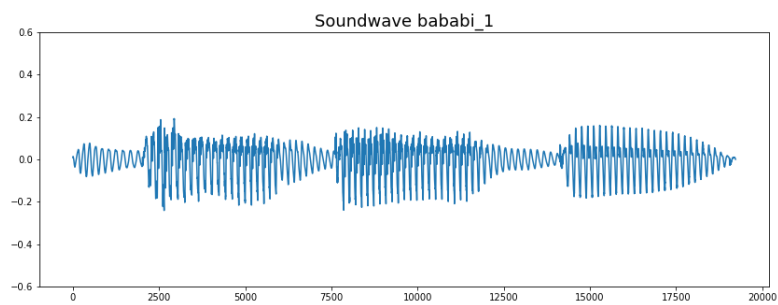


Figure 4.5:

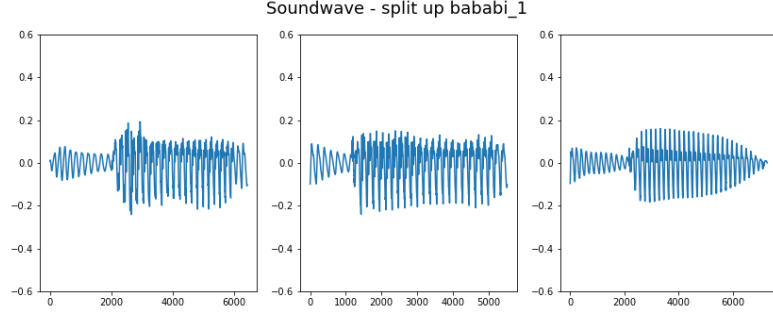


Figure 4.6:

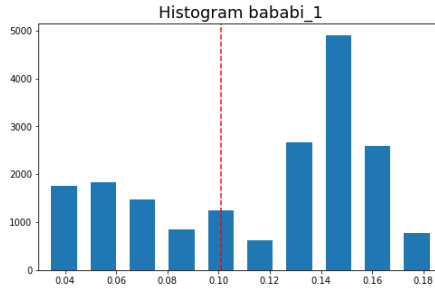


Figure 4.7:

4.1.2 Training and validation during training

Although loss is used as evaluation metric, it only shows part of the picture. The main objective for the representations is to obtain some form of "decoupled/dis-entangled" features that are more easily separable. This evaluation is done by projecting the latent representations to a 2D plane, via t-SNE. Then datapoints are coloured in depending on their the syllable that was pronounced, eg: "gi" or "ga".

Training is done speech signals of fixed length, eg 8,800 samples. Notice however, that the neural network only makes use of convolutional neural networks layers and GRU's, no fully connected layers. The input dimensions can therefore be variable during inference. Only the number of channels in the latent representations should be constant, but the length can change.

During inference, (in this context obtaining the latent representations for our input signals), depending on the length of the input signal, the length of the output latent representation will differ. If we wish to look at how separable latent representations are for syllables, the length can be variable. Some input sounds could be 6,600 samples, while others 8,800 samples. We therefore pad the syllables with zeroes in front and end of the signal, to obtain fixed length of equal to that of the longest syllable; 8,800 samples.

Training happens on longer data samples, and every \mathbf{X} epochs t-SNE visualisations are made to observe evolutional of dis-entanglement.

4.2 Decoders for Greedy Infomax

Objective: understand GIM learned representations. To do so we first discuss the model's architecture. To verify whether GIM does in fact contain the necessary information, we reconstruct the original speech signal from a GIM representations. We define train a decoder model for each of the four layers of the GIM model.

- Layer 1 ... 4
- Dataset from Bart
- Reduce audio to 16khz, reducing the size of the dataset and speeding up training time
- Add background noise to the data.
- Use this dataset: [openslr.org](https://open.slr.org) instead of white noise as is more representative.

4.3 Feature selection for visualisation

The GIM architecture consists of 4 layers, each layer 512 kernels Due to the large number of learned kernels, we filter based on standard deviation. This will result in looking at weights which are sensitive to particular features.

4.4 Variational contrastive predictive coding

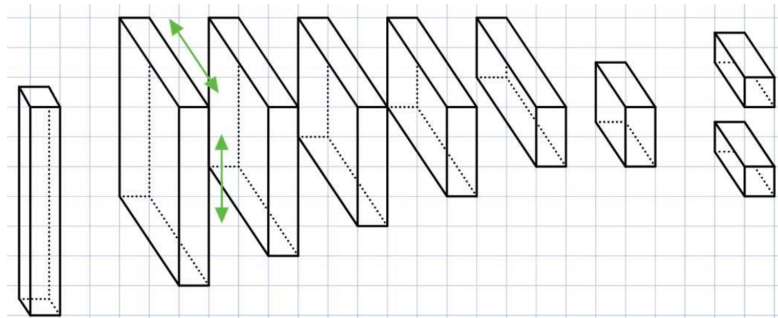


Figure 4.8:

Figure 4.8 displays how via contrastive predictive coding an input speech signal is transformed into two latent vectors. The two vectors combined describe a Gaussian distribution for each feature of the latent representation. One vector corresponds to the means and the other to the standard deviations of the latent distributions. We variational autoencoders assume independent latent features, such that the covariance matrix is non negative on the diagonal and zero off the diagonal. This allows the covariance matrix to be discribed using a single vector. We again, make use of this (plausible incorrect?) assumption of having independent latent features.

4.5 CHATGPT ANSWERS FOR VAE

Yes, you can incorporate the regularization component from a VAE into the contrastive predictive coding (CPC) loss function to encourage the learned latent representations to be Gaussian distributed. This can be achieved by adding a KL divergence term between the learned latent distribution and a Gaussian prior distribution to the CPC loss function.

The modified loss function would be:

$$\mathcal{L}N + KL = - \sum k \mathbb{E}_X \left[\log \frac{f_k(zt + k, c_t)}{\sum_{z_j \in X} f_k(z_j, c_t)} \right] + \lambda \text{KL}(q(z|X) || \mathcal{N}(0, I)) \quad (4.1)$$

where $q(z|X)$ is the learned distribution of latent representations, and $\mathcal{N}(0, I)$ is the standard Gaussian prior distribution. The hyperparameter λ controls the weight of the KL divergence term in the overall loss function.

To compute the KL divergence term, we need to calculate the mean and covariance of the learned latent distribution $q(z|X)$. Let μ and Σ be the mean and covariance of $q(z|X)$, respectively. Then the KL divergence can be computed as:

$$\text{KL}(q(z|X) || \mathcal{N}(0, I)) = \frac{1}{2} (\text{tr}(\Sigma) + \mu^T \mu - k - \log \det(\Sigma)) \quad (4.2)$$

where k is the dimensionality of the latent space.

The function $f_k(z_j, c_t)$ can be defined using a linear layer as follows:

code

Here, ‘hidden size’ is the size of the hidden layer, ‘latent size’ is the size of the latent space, and c is the context vector for a given time step. The function $W_k(z)$ applies a linear transformation to the input z , resulting in a latent vector of size ‘latent size’.

To incorporate the KL divergence term into the CPC loss function, you can compute the mean and covariance of the learned latent distribution using the `torch.mean` and `torch.var` functions, respectively. The modified loss function can then be implemented as follows:

code

Here, z is the learned latent representation, c is the context vector, and `lambd` is the weight of the KL divergence term. The ‘compute logits’ function computes the logits for the CPC loss as described in the original CPC paper.

Note that this implementation assumes that the ...

To incorporate the regularization component from VAE into the CPC loss function, we can add the Kullback-Leibler (KL) divergence term to the loss function, which will encourage the latent representations to be Gaussian distributed. The KL divergence term measures the difference between the learned latent distribution and the target Gaussian distribution.

The KL divergence term can be defined as:

$$\mathcal{L}KL = \frac{1}{2} \sum_i = 1^d (\mu_i^2 + \sigma_i^2 - \log(\sigma_i^2) - 1)$$

where μ_i and σ_i are the mean and standard deviation of the i -th dimension of the learned latent distribution, and d is the dimension of the latent space.

To combine the CPC loss and the KL divergence term, we can simply add them together, with a weighting factor β to balance their contributions:

$$\mathcal{L} = \mathcal{L}N + \beta \mathcal{L}KL$$

where \mathcal{L}_N is the CPC loss function defined in equation (2.17), and β is a hyperparameter that controls the strength of the regularization term.

The final loss function will encourage the learned latent representations to be both informative for the CPC task and Gaussian distributed, which can be useful for downstream tasks such as clustering or classification.

Sure! In the context of contrastive predictive coding (CPC), the ‘compute logits’ function is used to compute the similarity scores (logits) between the predicted future and the actual future.

In more detail, ‘compute logits’ takes in two inputs, z and c , which represent the latent representation and the context for a given time step, respectively. The function computes the similarity score between the predicted future representations and the actual future representations by performing a dot product between the predicted future representation and a set of parameters (i.e., weights) that are learned during training. Specifically, ‘compute logits’ applies a linear transformation to the context c and uses it to score each element in the predicted future z .

The output of ‘compute logits’ is a vector of logits that corresponds to the similarity scores between each predicted future representation and the actual future representation. These logits are then used to compute the contrastive loss, which is a measure of the similarity between the predicted future and the actual future.

The ‘compute logits’ function is a crucial component of the CPC algorithm, as it determines how well the predicted future representations match the actual future representations, which in turn affects the quality of the learned latent representation.

4.6 Results

4.6.1 Results CPC Simple v2 w/ 2 modules

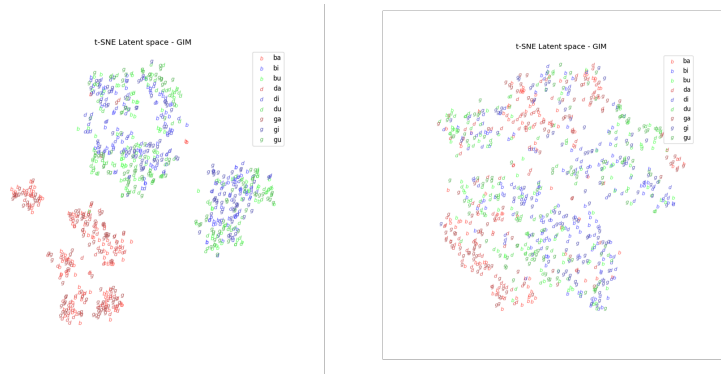


Figure 4.9:

4.9: can see that second model harms performance. We believe this can be explained via the learning rate. 4.10, there we see that second module performs better separation, indicating that the intermediate KL convergence constraint (causing the normal Gaussian distributions) also serves as a batch normalisation term, and thus resulting in faster convergence.

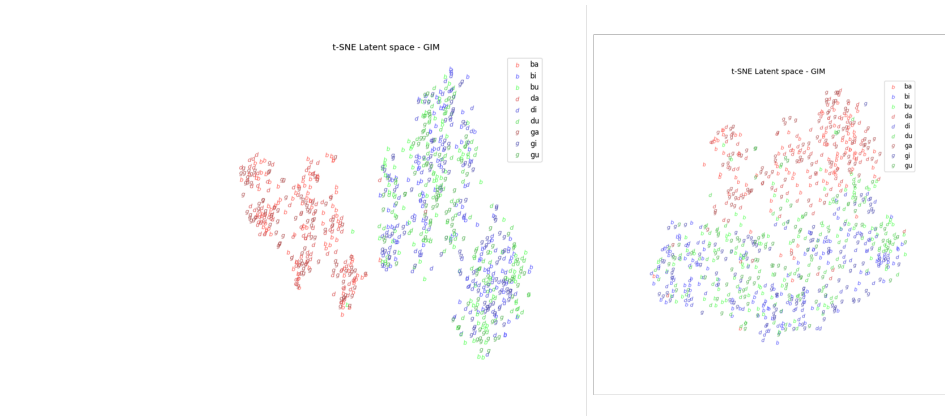


Figure 4.10:

4.7 GIM: Activations visualisations

thought for later: its actually weird i was able to play enc as audio as enc is 512 x something so huh? that means that a lot of info is already in first channel? what do other 511 channels then contain? """ Observations: First layer decoded still contains the same sound, but with some added noise (could be because decoder hasn't trained very). However, the encoded first layer, still contains the exact sound as the original sound. It is however downsampled a lot - from 16khz to 3khz """ thought for later: its actually weird i was able to play enc as audio as enc is 512 x something so huh? that means that a lot of info is already in first channel? what do other 511 channels then contain?

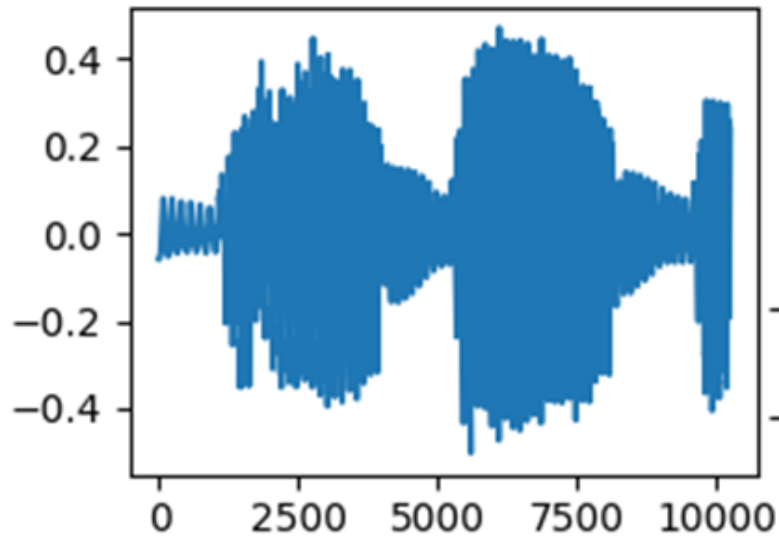


Figure 4.11: "BA-BA-BA" time domain

No batch normalisation, so although channels appear to have larger activations than other



Figure 4.12: Activations of the sound "BA-BA-BA" through GIM

channels, size of activation does not really say anything about information. eg activations 0.01 could still contain more information than 3.0 activation.

Since the activations from convolutional neural networks, the order is still maintained. Hence, can align activations with original signal.

Observations in latent representations:

Layer 1: The activations of the first decoder still contain a lot of similarity with the original signal, in terms of structure. There is a lot of redundant data within the representation. Eg: the one channel could be replayed

Layer 2

Layer 3:

Layer 4: Still notices multiple channels which have high activations when signal is has high amplitudes and small activations when amplitude is low.

Also activations which are high when volume is low. $-i$ indicates that certain kernel weights are sensitive for "**klinkers**" and other kernels for **medeklinkers**. see 4.13.

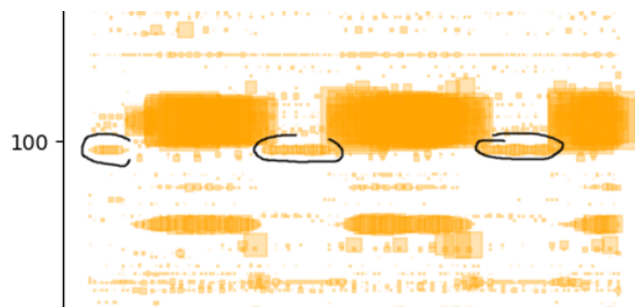


Figure 4.13: zoomed in

Observe that activations happen in clusters/sequences. So it is usually a patch of signal samples that cause high activations. This could for instance indicate that both kernels are sensitive for the **medeklinker** "b", but sensitive for different features. eg the letter B has spoken sound "buh". so maybe one is sensitive for "b" and other for "uh".

Figure 4.14 also nicely shows how different channels have clusters of activations at slightly different times.

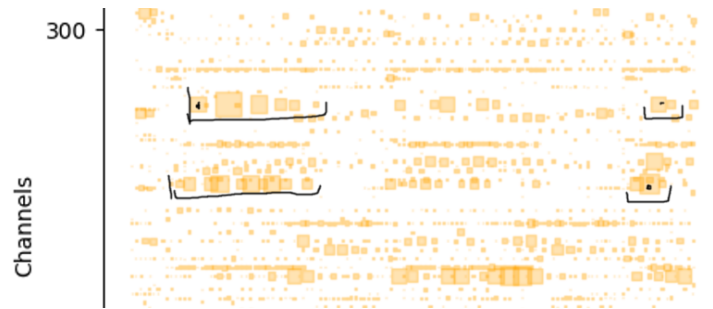


Figure 4.14: Zoomed in

T-sne visualisations: Multiple models trained, trained GIM, but single module (exact CPC architecture), one with autoregressor and once without autoregressor layer, so only CNN layers. 200 epochs, trained on split up data samples.

- 1) Only CNN: We observe better linear separability
- 2) CNN + 1 autoregressor layer:
- 3) This is the pure data visualised (no latent representation). Were at least doing a bit better than the original data, so that's good! our work is not for nothing.
- 4) old GIM with all modules each one layer. l1 .. 5 - cnn, l6 = gru. img shows l5 = cnn: model can more easily distinguish A's from other **klinkers**. Partly, it kinda makes sense for GIM to learn to separate **klinkers**. Since they last longer (longer duration), the loss function will more likely randomly sample a subwindow from the "aa's" than from the **medeklinker** part.

4.8 Decoder: predictions on test set

Fig ?? displays the reconstructed signal from the vocal sound "ba-gi-di". The two images on the left displays the original signal, while the right two images contain the reconstructed signal. The upper images displays the signals in time domain, the bottom images spectral domain. The reconstructed signal is an audio sample, for instance which is encoded via Greedy Infomax (up to the fourth (and final) convolution layer), this output is then given to a decoder to reconstruct the original signal.

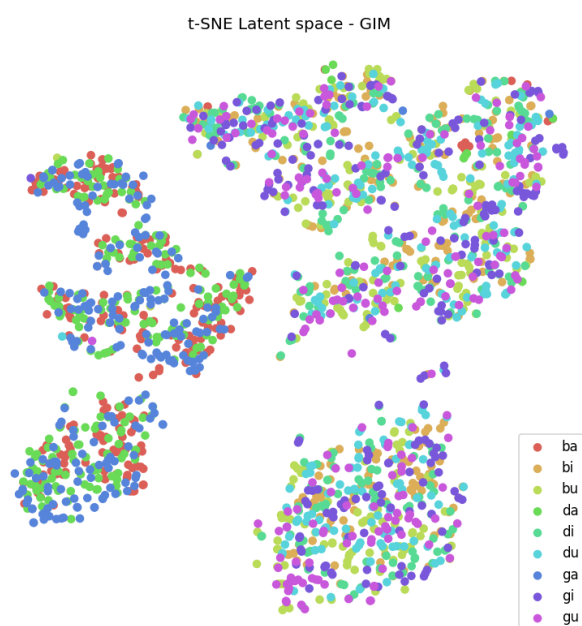


Figure 4.15:

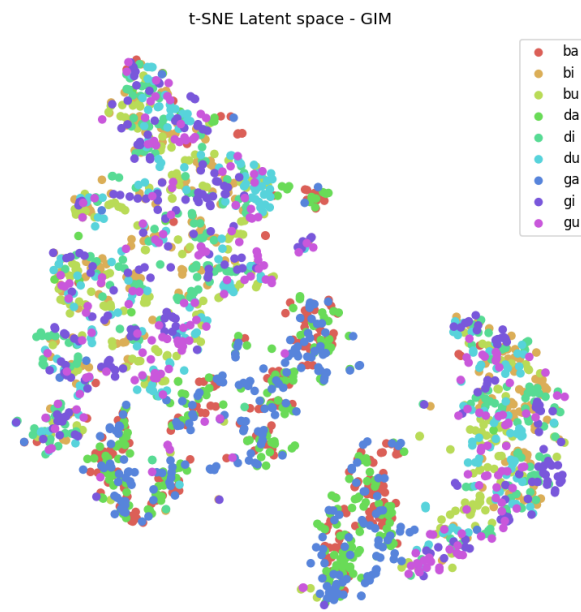


Figure 4.16:

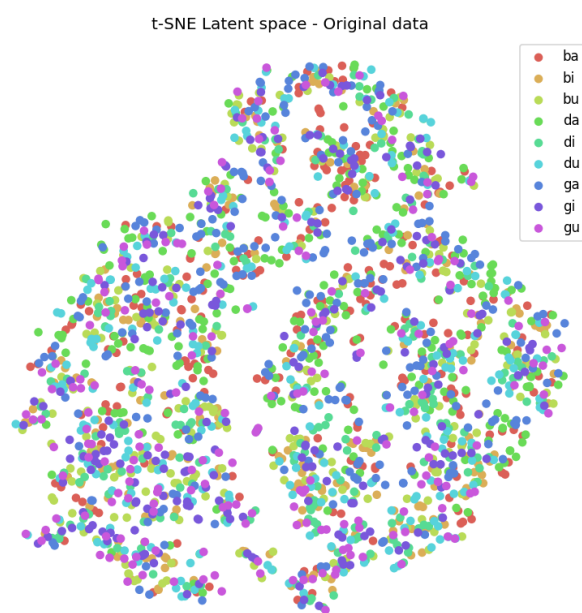


Figure 4.17:

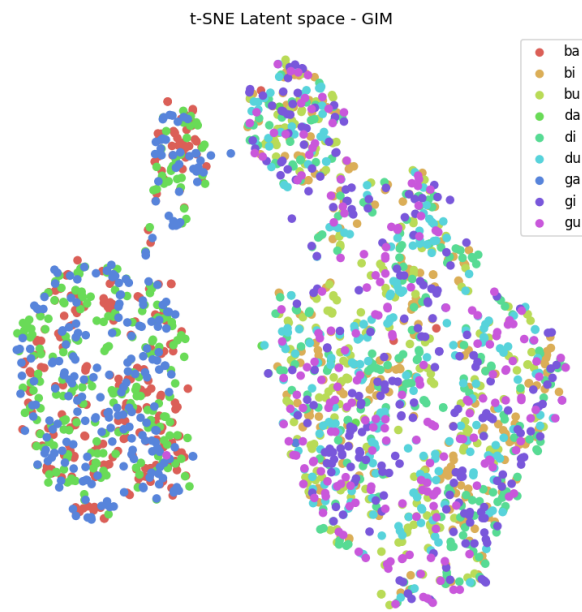


Figure 4.18:

Chapter 5

Related work

5.1 Representation learning: explainable

5.2 Variational learning

There already were a few papers with variational contrastive predictive coding

5.3 Links I should investigate

- Adversarial and Contrastive Variational Autoencoder for Sequential Recommendation <https://arxiv.org/pdf/2103.10693>.
- S3VAE: Self-Supervised Sequential VAE for Representation Disentanglement and Data Generation <https://arxiv.org/abs/2005.11437>
- Implementation of Sequential VAE <https://github.com/ermongroup/Sequential-Variational-Autoencoder>
- Contrastively Disentangled Sequential Variational Autoencoder <https://proceedings.neurips.cc/paper/2021/file/53c55b1e1e1e1e1e1e1e1e1e1e1e1e1e-Paper.pdf>
- Sequential Variational Autoencoders for Collaborative Filtering <https://arxiv.org/pdf/1811.09975.pdf>
- !! Variational noise contrastive estimation: <https://arxiv.org/abs/1810.08010>

Chapter 6

Discussion

while GIM argues to resolve memory constraints, not entirely true. In fact we even countered the opposite as containing multiple neural networks, each with their own personal loss function (the loss function is based on `fk` which contains parameters that must be learned), and thus for early layers where the sequence is still long, a lot of memory is required. We went for a compromise on GIM by splitting up the architecture in merely two modules, significantly reducing the memory constraints.

— The second module in GIM clearly doesn't have as much effect. This can be explained because there may not be as much common information anymore between the patches. There may be a source that says that `cpc` learns low level features, but the second module is supposed to learn more high level features, which `cpc` may have trouble with? —

Future work: - Related work in VAE shows that gradually increasing regularisation term, results in better disentanglement, while avoiding posterior collapse. could have a `kldweight` scheduler.

- not constrained solely to InfoNCE loss, the GIM architecture could work for other losses too that allow for greedy optimisation.

- I didn't add an autoregressor as i didn't find a performance benefit. Potentially, with larger architecture could further improve performance.

- Explainability of latents is dependent on the performance of the decoder.
- Intermediate loss function with `kld` resulted in similar behaviour as batch normalisation. Resulting in faster convergence than without `kld`.
- We observed no quality loss in the learned representations. Data was equally easily separable.

Bibliography

- [1] T. M. Cover and J. A. Thomas, *ELEMENTS OF INFORMATION THEORY*.
- [2] “Ali Ghodsi, Lec : Deep Learning, Variational Autoencoder, Oct 12 2017 [Lect 6.2].” [Online]. Available: <https://www.youtube.com/watch?v=uaaqyVS9-rM>
- [3] Y. Bengio, A. Courville, and P. Vincent, “Representation Learning: A Review and New Perspectives,” vol. 35, pp. 1798–1828.
- [4] P. H. Le-Khac, G. Healy, and A. F. Smeaton, “Contrastive Representation Learning: A Framework and Review,” vol. 8, pp. 193 907–193 934.
- [5] D. Rumelhart, G. Hinton, and R. Williams, “Learning Internal Representations by Error Propagation,” in *Readings in Cognitive Science*. Elsevier, pp. 399–421. [Online]. Available: <https://linkinghub.elsevier.com/retrieve/pii/B9781483214467500352>
- [6] D. Bank, N. Koenigstein, and R. Giryes. Autoencoders. [Online]. Available: <http://arxiv.org/abs/2003.05991>
- [7] S. Karagiannakos. How to Generate Images using Autoencoders. AI Summer. [Online]. Available: <https://theaisummer.com/Autoencoder/>
- [8] Papers with Code - MNIST Dataset. [Online]. Available: <https://paperswithcode.com/dataset/mnist>
- [9] C. Doersch. Tutorial on Variational Autoencoders. [Online]. Available: <http://arxiv.org/abs/1606.05908>
- [10] David Foster, “3. Variational Autoencoders,” in *Generative Deep Learning, 2nd Edition*, 2nd ed. [Online]. Available: <https://www.oreilly.com/library/view/generative-deep-learning/9781098134174/>
- [11] D. P. Kingma and M. Welling. Auto-Encoding Variational Bayes. [Online]. Available: <http://arxiv.org/abs/1312.6114>
- [12] —, “An Introduction to Variational Autoencoders,” vol. 12, no. 4, pp. 307–392. [Online]. Available: <http://arxiv.org/abs/1906.02691>
- [13] L. P. Cinelli, M. A. Marins, E. A. Barros da Silva, and S. L. Netto, *Variational Methods for Machine Learning with Applications to Deep Networks*. Springer International Publishing. [Online]. Available: <https://link.springer.com/10.1007/978-3-030-70679-1>
- [14] p. d. u. family=Oord, given=Aaron, Y. Li, and O. Vinyals. Representation Learning with Contrastive Predictive Coding. [Online]. Available: <http://arxiv.org/abs/1807.03748>

- [15] S. Löwe, P. O'Connor, and B. S. Veeling. Putting An End to End-to-End: Gradient-Isolated Learning of Representations. [Online]. Available: <http://arxiv.org/abs/1905.11786>
- [16] S. S. Rao. Understanding the Gradient-Isolated Learning of Representations and intuition to the Greedy.... Analytics Vidhya. [Online]. Available: <https://medium.com/analytics-vidhya/understanding-the-gradient-isolated-learning-of-representations-and-intuition-to-the-greedy-cb6c3598e317>
- [17] Z. Zhang and D. Tao, "Slow Feature Analysis for Human Action Recognition," vol. 34, no. 3, pp. 436–450.
- [18] K. Stacke, C. Lundström, J. Unger, and G. Eilertsen, "Evaluation of Contrastive Predictive Coding for Histopathology Applications," in *Proceedings of the Machine Learning for Health NeurIPS Workshop*. PMLR, pp. 328–340. [Online]. Available: <https://proceedings.mlr.press/v136/stacke20a.html>
- [19] p. u. family=Haan, given=Puck and S. Löwe. Contrastive Predictive Coding for Anomaly Detection. [Online]. Available: <http://arxiv.org/abs/2107.07820>
- [20] M. Y. Lu, R. J. Chen, J. Wang, D. Dillon, and F. Mahmood. Semi-Supervised Histology Classification using Deep Multiple Instance Learning and Contrastive Predictive Coding. [Online]. Available: <http://arxiv.org/abs/1910.10825>
- [21] S. Bhati, J. Villalba, P. Želasko, L. Moro-Velazquez, and N. Dehak. Segmental Contrastive Predictive Coding for Unsupervised Word Segmentation. [Online]. Available: <http://arxiv.org/abs/2106.02170>
- [22] S. Deldari, D. V. Smith, H. Xue, and F. D. Salim, "Time Series Change Point Detection with Self-Supervised Contrastive Predictive Coding," in *Proceedings of the Web Conference 2021*, ser. WWW '21. Association for Computing Machinery, pp. 3124–3135. [Online]. Available: <https://dl.acm.org/doi/10.1145/3442381.3449903>
- [23] O. Henaff, "Data-Efficient Image Recognition with Contrastive Predictive Coding," in *Proceedings of the 37th International Conference on Machine Learning*. PMLR, pp. 4182–4192. [Online]. Available: <https://proceedings.mlr.press/v119/henaff20a.html>
- [24] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "ImageNet classification with deep convolutional neural networks," in *Advances in Neural Information Processing Systems*, F. Pereira, C. Burges, L. Bottou, and K. Weinberger, Eds., vol. 25. Curran Associates, Inc. [Online]. Available: https://proceedings.neurips.cc/paper_files/paper/2012/file/c399862d3b9d6b76c8436e924a68c45b-Paper.pdf
- [25] S. Ioffe and C. Szegedy. Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift. [Online]. Available: <http://arxiv.org/abs/1502.03167>
- [26] A. H. Marblestone, G. Wayne, and K. P. Kording, "Toward an Integration of Deep Learning and Neuroscience," vol. 10. [Online]. Available: <https://www.frontiersin.org/articles/10.3389/fncom.2016.00094>
- [27] N. Caporale and Y. Dan, "Spike timing-dependent plasticity: A Hebbian learning rule," vol. 31, pp. 25–46.

- [28] R. Wei and A. Mahmood, “Recent Advances in Variational Autoencoders With Representation Learning for Biomedical Informatics: A Survey,” vol. 9, pp. 4939–4956.
- [29] M. Tschannen, O. Bachem, and M. Lucic. Recent Advances in Autoencoder-Based Representation Learning. [Online]. Available: <http://arxiv.org/abs/1812.05069>
- [30] M. Grossutti, J. D’Amico, J. Quintal, H. MacFarlane, A. Quirk, and J. R. Dutcher, “Deep Learning and Infrared Spectroscopy: Representation Learning with a Beta-Variational Autoencoder,” vol. 13, no. 25, pp. 5787–5793. [Online]. Available: <https://doi.org/10.1021/acs.jpcclett.2c01328>