



UNIVERSIDAD
DE GRANADA

MUTUAL INFORMATION IN SELF-SUPERVISED LEARNING

FRANCISCO JAVIER SÁEZ MALDONADO

Bachelor's Thesis

Computer Science and Mathematics

Tutor

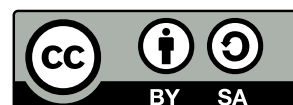
Nicolás Pérez de la Blanca Capilla

FACULTY OF SCIENCE

H.T.S. OF COMPUTER ENGINEER AND TELECOMMUNICATIONS

Granada, Thursday 12th August, 2021

This work is licensed under a [Creative Commons](#) “[Attribution-ShareAlike 4.0 International](#)” license.



The source code of this text and developed programs are available in the Github repository [fjsaezm/Mutual-Information-in-Unsupervised-Machine-Learning](#)

ABSTRACT

Abstract goes here

RESUMEN

El resumen en español va aquí

INTRODUCTION

MAIN GOALS AND RESULTS ACHIEVED

CONTENTS

I BASIC NOTIONS

1	PROBABILITY THEORY	11
1.1	Probability spaces	11
1.2	Expectation of a random variable	13
1.3	Random vectors	15
2	DISTRIBUTIONS AND KULLBACK-LEIBLER DIVERGENCE	17
2.1	Examples of distributions	18
2.1.1	Bernoulli	18
2.1.2	Gaussian Distribution	19
3	STATISTICAL INFERENCE	20
3.1	Parametric Modeling	20
3.2	Minimal sufficient statistics	21
4	NOISE CONTRASTIVE ESTIMATION	23
4.1	Logistic regression	23
4.1.1	Fitting a Logistic Regression model	24
4.2	Motivation	25
4.3	Formalization of the NCE problem	26

II INFORMATION THEORY

5	ENTROPY	31
5.1	Definition	31
5.2	Properties of the entropy. Conditional entropy	32
6	MUTUAL INFORMATION	35
6.1	Lower bounds on Mutual Information	36
6.1.1	Variational Lower Bound	37
6.1.2	Donsker-Varadhan Representation lower bound	37
6.1.3	Contrastive Lower Bound	38

III REPRESENTATION LEARNING

7	CONTEXT	42
8	GENERATIVE MODELS	45
8.1	Autoregressive Models	46
9	THE INFO NCE LOSS	48
9.1	Contrastive Predictive Coding	49
9.2	Good views for Contrastive Learning	51
10	TRIPLET LOSSES	54
10.1	From deep metric learning to triplet losses and its generalization	54
10.2	Generalization of triplet losses and N -pairs loss	56
10.2.1	InfoNCE Bound as a triplet loss	57

IV NEW FRAMEWORKS FOR REPRESENTATION LEARNING

11	INTRODUCTION TO DEEP LEARNING	60
11.1	Neural Networks	61
11.2	Convolutional Neural Networks	62
11.2.1	ResNet	63
11.3	Data augmentation	64
12	SIMCLR	68
12.1	Introduction	68
12.2	The framework	68
12.3	Findings of SimCLR	71
13	BOOSTRAP YOUR OWN LATENT	73
13.1	Motivation	73
13.2	BYOL Algorithm	73
13.3	Results obtained by this framework	76

V EXPERIMENTS

14	INTRODUCTION	78
14.1	Hardware and basic libraries	78
14.2	The dataset: CIFAR10	79
14.3	Tensorflow	80
14.3.1	Tensorboard	81
14.4	Metrics	82
15	EXPERIMENTATION	83
15.1	SimCLR exploration	83
15.1.1	First approach	83

VI APPENDIX

A	APPENDIX A	88
B	APPENDIX B	89
B.1	SimCLR Experiments	89

LIST OF FIGURES

Figure 1	Bimodal Distribution	17
Figure 2	Plants' height	17
Figure 3	Examples of bimodal and multimodal distributions. . .	17
Figure 4	Representation of $H(p)$ in the example 4.	32
Figure 5	Image from (Oord <i>et al.</i> , 2019). Overview of Contrastive Predictive Coding framework using audio signal as input.	50
Figure 6	Example of an anchor x , a positive instance x^+ and a negative instance x^- . Images obtained from Google. . .	55
Figure 7	Sigmoid. Image from this Medium article	61
Figure 8	Hyperbolic Tangent. Image from this Medium article	61
Figure 9	ReLU. Image from this Medium article	61
Figure 10	Residual Block. Image from He <i>et al.</i> (2015)	64
Figure 11	The image on the left is the original image, and the image on the right is a new example generated by first a rotation of a random angle and then a flip. . . .	65
Figure 12	The image on the left is the original image, and the image on the right is a new example generated by performing a random crop of the image and then re-sizing the crop to the original image size.	66
Figure 13	Histogram for both the original image on blue and the random crop on orange.	66
Figure 14	Color jitter applied to tulips image.	67
Figure 15	Histograms of the color jitter image in orange and cropped image in blue.	67
Figure 16	Figure obtained from Chen <i>et al.</i> (2020b). A simple framework for contrastive learning of visual representations.	69
Figure 17	Algorithm that summarizes the learning process that SimCLR follows.	70
Figure 18	Image from (Grill <i>et al.</i> , 2020). Overview of Bootstrap Your Own Latent Framework.	75
Figure 19	Ten examples of each class in the CIFAR10 dataset. . .	79
Figure 20	Tensorflow logo.	80
Figure 21	Results of the batch-size experiment.	85
Figure 22	Accuracy score following the color jitter parameter and both batch sizes considered.	86
Figure 23	Accuracy score following the temperature and both batch sizes considered.	86

LIST OF TABLES

Table 1	Resnet 18 architecture.	64
Table 2	Comparison between SimCLR and BYOL Top1 and Top5 accuracies using the same architectures on the ImageNet dataset.	76
Table 3	Best results for the grid search experiment with Sim- CLR.	84
Table 4	All results for first experiment using SimCLR.	90

Part I

BASIC NOTIONS

Underneath each experiment involving any grade of uncertainty there is a *random variable*. This is no more than a *measurable* function between two *measurable spaces*. A probability space is composed by three elements: $(\Omega, \mathcal{A}, \mathcal{P})$. We will define those concepts one by one.

1.1 PROBABILITY SPACES

Definition 1.1.1. Let Ω be a non empty sample space. \mathcal{A} is a σ -algebra over Ω if it is a family of subsets of Ω that verify that the emptyset is in \mathcal{A} , and it is closed under complementation and countable unions. That is:

- $\emptyset \in \mathcal{A}$.
- If $A \in \mathcal{A}$, then $\Omega \setminus A \in \mathcal{A}$.
- If $\{A_i\}_{i \in \mathbb{N}} \in \mathcal{A}$ is a numerable family of \mathcal{A} subsets, then $\bigcup_{i \in \mathbb{N}} A_i \in \mathcal{A}$.

The pair (Ω, \mathcal{A}) is called a *measurable space*. To get to our probability space, we need to define a *measure* on the *measurable space*.

Definition 1.1.2. Given (Ω, \mathcal{A}) a measurable space, a *measure* \mathcal{P} is a countable additive, non-negative set function on this space. That is: $\mathcal{P} : \mathcal{A} \rightarrow \mathbb{R}_0^+$ satisfying:

- $\mathcal{P}(A) \geq \mathcal{P}(\emptyset) = 0$ for all $A \in \mathcal{A}$,
- $\mathcal{P}(\bigcup_n A_n) = \sum_n \mathcal{P}(A_n)$ for any countable collection of disjoint sets $A_n \in \mathcal{A}$.

If $\mathcal{P}(\Omega) = 1$, \mathcal{P} is a *probability measure* or simply a *probability*. With the concepts that have just been explained, we get to the following definition:

Definition 1.1.3. A *measure space* is the tuple $(\Omega, \mathcal{A}, \mathcal{P})$ where \mathcal{P} is a *measure* on (Ω, \mathcal{A}) . If \mathcal{P} is a *probability measure* $(\Omega, \mathcal{A}, \mathcal{P})$ will be called a *probability space*.

Throughout this work, we will be always in the case where \mathcal{P} is a probability measure, so we will always be talking about probability spaces and we will note \mathcal{P} simply as P . Some notation for these measures must be introduced. Let A and B be two events. The notation $P(A, B)$ refers to the probability of the intersection of the events A and B , that is: $P(A, B) := P(A \cap B)$. It is clear that since $A \cap B = B \cap A$, then $P(A, B) = P(B, A)$. We remark the next definition since it will be important.

Definition 1.1.4. Let A, B be two events in Ω . The *conditional probability* of B given A is defined as:

$$P(B|A) = \frac{P(A, B)}{P(A)}.$$

There is an alternative way to state the definition that we have just made.

Theorem 1.1.1 (Bayes' Theorem). Let A, B be two events in Ω , given that $P(B) \neq 0$. Then

$$P(B|A) = \frac{P(A|B)P(B)}{P(A)}.$$

Proof. Straight from the definition of the conditional probability we obtain that:

$$P(A, B) = P(A|B)P(B).$$

We also see from the definition that

$$P(B, A) = P(B|A)P(A).$$

Hence, since $P(A, B) = P(B, A)$,

$$P(A|B)P(B) = P(B|A)P(A) \implies P(A|B) = \frac{P(B|A)P(A)}{P(B)}.$$

□

However, events might not give any information about another event occurring. When this happens, we call those events to be *independent*. Mathematically, if A, B are independent events:

$$P(A, B) = P(A)P(B)$$

and as a consequence of this, the conditional probability of those events is $P(A|B) = P(A)$. For a finite set of events $\{A_i\}_{i=1}^n$, we say that they are mutually independent if, and only if, every event is independent of any intersection of the other events. That is, if $\{B_i\} \subset \{A_i\}$, then

$$P\left(\bigcap_{i=1}^k B_i\right) = \prod_{i=1}^k P(B_i) \quad \text{for all } k \leq n.$$

Random variables (RV) can now be introduced. Their first property is that they are measurable functions. This kind of functions are defined as it follows:

Definition 1.1.5. Let $(\Omega_1, \mathcal{A}), (\Omega_2, \mathcal{B})$ be measurable spaces. A function $f : \Omega_1 \rightarrow \Omega_2$ is said to be *measurable* if, $f^{-1}(B) \in \mathcal{A}$ for every $B \in \mathcal{B}$.

As a quick note, we can affirm that if f, g are real-valued measurable functions, and $k \in \mathbb{R}$, it is true that kf , $f + g$, fg and f/g (if g is not the identically zero function) are also *measurable functions*.

We are now ready to define one of the concepts that will lead us to the main objective of this thesis.

Definition 1.1.6. Let $(\Omega, \mathcal{A}, \mathcal{P})$ be a probability space, and (E, \mathcal{B}) be a measurable space. A *random variable* is a measurable function $X : \Omega \rightarrow E$, from the probability space to the measurable space. This means: for every subset $B \in (E, \mathcal{B})$, its pre-image

$$X^{-1}(B) = \{\omega : X(\omega) \in B\} \in \mathcal{A}.$$

Using that sums, products and quotients of measurable functions are measurable functions, we obtain that *sums, products and quotients of random variables are random variables*.

Let now X be a R.V. The *probability* of X taking a concrete value on a measurable set contained in E , say, $S \in E$, is written as:

$$P_X(S) = P(X \in S) = P(\{a \in \Omega : X(a) \in S\}).$$

We use the notation P_X to remark that P is related to the distribution of X . Most of the times, when it is clear from the context, the subindex will be omitted and we will simply refer to it $p(x)$. A very simple example of random variable is the following:

Example 1. Consider tossing a coin. The possible outcomes of this experiment are *Heads* or *Tails*. Those are our random events. We can give our random events a possible value. For instance, let *Heads* be 1 and *Tails* be 0. Then, our random variable looks like this:

$$X = \begin{cases} 1, & \text{if we obtain heads,} \\ 0, & \text{if we obtain tails.} \end{cases}$$

In the last example, our random variable is *discrete*, since the set $\{X(\omega) : \omega \in \Omega\}$ is finite. A random variable can also be *continuous*, if it can take any value within an interval.

1.2 EXPECTATION OF A RANDOM VARIABLE

Definition 1.2.1. The *cumulative distribution function* F_X of a real-valued random variable X is its probability of taking value below or equal to x . That is:

$$F_X(x) = P(X \leq x) = P(\{\omega : X(\omega) \leq x\}) = P_X((-\infty, x]) \quad \text{for all } x \in \mathbb{R}.$$

We can difference between certain types of random variables. If the image, \mathcal{X} , of X is countable, we call it a *discrete* random variable. Its *probability mass function* p gives the probability of the R.V. being equal to a certain value:

$$p(x) = P(X = x).$$

If the cumulative distribution function of our random variable X is continuous everywhere, then X is a *continuous* random variable. In this case there might exist a non-negative Lebesgue-integrable function f such that:

$$F_X(x) = \int_{-\infty}^x f(t)dt,$$

called the *probability density function* of X .

During this document, distributions $p(x)$ will be addressed many times. We will use $p(x)$ to denote the probability mass function of a discrete random variable or the probability density function of a continuous random variable.

We are now ready to introduce the *expectation* of a random variable. Imagine observing a wide number of outcomes from our random variable, and taking the average of these random values. The expectation is the value of this average when we take *infinite* outcomes of our random variable.

Definition 1.2.2. Let X be a non negative random variable on a probability space $(\Omega, \mathcal{A}, \mathcal{P})$. The *expectation* $E[X]$ of X is defined as:

$$E[X] = \int_{\Omega} X(\omega) dP(\omega).$$

Sometimes we might be referring to multiple random variables. In these cases, in order to make reference to the variable (or distribution function, that will be presented later) for which we calculate the expectation, we will denote it as E_X (or E_P , in the case that we are addressing a distribution).

The expectation of a random variable will be also denoted as μ . Now, if X is generic R.V, the expectation is defined as:

$$E[X] = E[X^+] - E[X^-],$$

where X^+, X^- are defined as it follows:

$$X^+(\omega) = \max(X(\omega), 0), \quad X^-(\omega) = \min(X(\omega), 0).$$

The expectation $E[X]$ of a random variable is a linear operation. That is, if Y is another random variable, and $\alpha, \beta \in \mathbb{R}$, then

$$E[\alpha X + \beta Y] = \alpha E[X] + \beta E[Y].$$

This is a trivial consequence of the linearity of the *Lebesgue integral*.

As a note, if X is a *discrete* random variable and \mathcal{X} is its image, its expectation can be computed as:

$$E[X] = \sum_{x \in \mathcal{X}} x P_X(x),$$

where x is each possible outcome of the experiment, and $P_X(x)$ the probability under the distribution of X of the outcome x . The expression given in Def. 1.2.2 generalizes this particular case.

Using the definition of the expectation of a random variable, we can approach to the concept of the *moments* of a random variable.

Definition 1.2.3. If $k \in \mathbb{N}$, then $E[X^k]$ is called the k -th moment of X .

If we take $k = 1$, we have the definition of the *expectation*. It is sometimes written as $m_X = E[X]$, and called the *mean*. We use the *mean* in the definition of the variance:

Definition 1.2.4. Let X be a random variable. If $E[X^2] < \infty$, then the *variance* of X is defined to be

$$\text{Var}(X) = E[(X - m_X)^2] = E[X^2] - m_X^2.$$

Thanks to the linearity of the *expectation* of a random variable, it is easy to see that, if $a, b \in \mathbb{R}$, then

$$\text{Var}(aX + b) = E[(aX + b) - E[aX + b]]^2 = a^2 E[(X - m_X)^2] = a^2 \text{Var}(X).$$

1.3 RANDOM VECTORS

Usually, when it comes to applying these concepts to a real problem, we will be observing multiple features that a phenomenon in nature presents. We would like to have a collection of random variables each one representing one of this features. In order to set the notation for these kinds of situations, we will introduce *random vectors*.

Definition 1.3.1. A random vector is a row vector $\mathbf{X} = (X_1, \dots, X_n)$ whose components are real-valued random variables on the same probability space (Ω, \mathcal{A}, P) .

The probability distribution of a random variable can be extended in to the *joint probability distribution* of a random vector.

Definition 1.3.2. Let $\mathbf{X} = (X_1, \dots, X_n)$ be a random vector. The *cumulative distribution function* (or simply, the *distribution function*) $F_{\mathbf{X}} : \mathbb{R}^n \rightarrow [0, 1]$ of \mathbf{X} is defined as:

$$F_{\mathbf{X}}(x) = P(X_1 \leq x_1, \dots, X_n \leq x_n).$$

We also name it *multivariate distribution*. Before, we presented the concept of independence between a pair of events. Using the cumulative distribution function, we can now define the independence between random variables.

Definition 1.3.3. A finite set of n random variables $\{X_1, \dots, X_n\}$ is mutually independent if, and only if, for any sequence $\{x_1, \dots, x_n\}$, the events $\{X_1 \leq x_1\}, \dots, \{X_n \leq x_n\}$ are mutually independent. Equivalently, this finite set is mutually independent if, and only if,:

$$F_{X_1, \dots, X_n}(x_1, \dots, x_n) = F_{X_1}(x_1) \dots F_{X_n}(x_n), \quad \text{for all } x_1, \dots, x_n.$$

We can also extend the notion of expectation to a random vector. Let $\mathbf{X} = (X_1, \dots, X_n)$ be a random vector and assume that $E[X_i]$ exists for all $i \in \{1, \dots, n\}$. The expectation of \mathbf{X} is defined as the vector containing the expectations of each individual random vector, that is:

$$E[\mathbf{X}] = \begin{bmatrix} E[X_1] \\ \vdots \\ E[X_n] \end{bmatrix}.$$

To generalize the variance of a random variable, we have to build the following matrix.

Definition 1.3.4. Let $\mathbf{X} = (X_1, \dots, X_n)$ be a random vector. Then, the *covariance matrix* of \mathbf{X} is defined as:

$$\Sigma = \text{Cov}(\mathbf{X}) = E[(\mathbf{X} - \mu_{\mathbf{X}})(\mathbf{X} - \mu_{\mathbf{X}})^T] = \begin{pmatrix} \sigma_{11} & \cdots & \sigma_{1n} \\ \vdots & \ddots & \vdots \\ \sigma_{n1} & \cdots & \sigma_{nn} \end{pmatrix},$$

where $\sigma_{ij} = \text{Cov}(X_i, X_j) = E[(X_i - \mu_i)(X_j - \mu_j)] = \sigma_{ji}$.

It can also happen that, given a random vector, we would like to know the probability distribution of some of its components. That is called the *marginal distribution*.

Definition 1.3.5. Let $\mathbf{X} = (X_1, \dots, X_n)$ be a random vector. The *marginal distribution* of a subset of \mathbf{X} is the probability distribution of the variables contained in the subset.

In the simple case of having two random variables, e.g. $\mathbf{X} = (X_1, X_2)$, then the marginal distribution of X_1 is:

$$P(x) = \int_{x_2} P(x_1, x_2) dx_2.$$

2

DISTRIBUTIONS AND KULLBACK-LEIBLER DIVERGENCE

We have introduced the concepts of *random variable*, *random vector* and its *probability distribution*. We will explain some concepts related to the latter ones.

Definition 2.1. The *mode* of a distribution is the value at which the probability mass function takes its maximum value. That is, the value that is most likely to be sampled.

Distributions can be *unimodal*, when their distribution has a single peak, *bimodal* when their distribution has two peaks, and *multimodal* when the number of peaks is equal or greater to 2.

Example 2. We can simulate the two following distributions:

1. The distribution of the marks obtained in a test by the students of certain class.
2. The distribution of the height of the plants from three different species.

The result is the following:

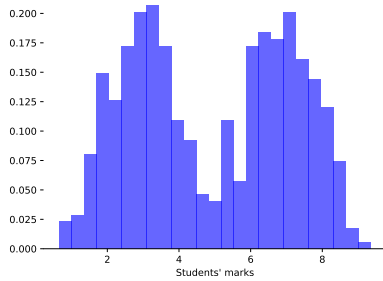


Figure 1: Bimodal Distribution

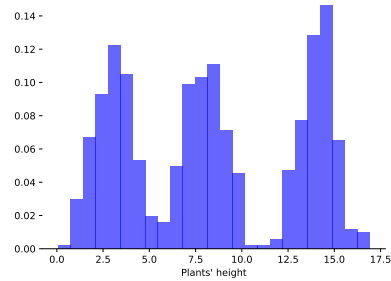


Figure 2: Plants' height

Figure 3: Examples of bimodal and multimodal distributions.

Now, given two distributions, we would like to determine how different they are from each other. In order to compare them, we enunciate the definition of the Kullback-Leibler divergence.

Definition 2.2. Let P and Q be probability distributions over the same probability space Ω . Then, the Kullback-Leibler divergence is defined as:

$$D_{KL}(P \parallel Q) = E_P \left[\log \frac{P(x)}{Q(x)} \right].$$

It is defined if, and only if, P is *absolutely continuous with respect to* Q , that is, if $P(A) = 0$ for any A subset of Ω where $Q(A) = 0$. There are some properties of this definition that must be stated.

Proposition 1. *If P, Q are two probability distributions over the same probability space, then $D_{KL}(P|Q) \geq 0$.*

Proof. Firstly, note that if $a \in \mathbb{R}^+$, then $\log a \leq a - 1$. Then:

$$\begin{aligned} -D_{KL}(P \parallel Q) &= -E_P \left[\log \frac{P(x)}{Q(x)} \right] \\ &= E_P \left[\log \frac{Q(x)}{P(x)} \right] \\ &\leq E_P \left[\left(\frac{Q(x)}{P(x)} - 1 \right) \right] \\ &= \int P(x) \frac{Q(x)}{P(x)} dx - 1 \\ &= 0. \end{aligned}$$

So we have obtained that $-D_{KL}(P \parallel Q) \leq 0$, which implies that $D_{KL}(P \parallel Q) \geq 0$. \square

As a corollary of this proposition, we can affirm that $D_{KL}(P \parallel Q)$ equals zero if and only if $P = Q$ almost everywhere. We will also remark the discrete case, as it will be used later. Let P, Q be discrete probability distributions defined on the same probability space Ω . Then,

$$D_{KL}(P \parallel Q) = \sum_{x \in \Omega} P(x) \log \left(\frac{P(x)}{Q(x)} \right).$$

2.1 EXAMPLES OF DISTRIBUTIONS

Let us present some examples of common distributions. They will be used further in this document.

2.1.1 Bernoulli

Think for a moment that you want to model the possible outcomes of an experiment with two possibilities: success or failure. Imagine also that you already know that in your experiment there is a probability p of achieving success. That is the intuitive idea of a Bernoulli distribution. We can define it more formally as follows:

The *Bernoulli distribution* is a discrete probability distribution of a random variable that takes two values, $\{0, 1\}$, with probabilities p and $q = 1 - p$, respectively. We will say that our distribution is a *Bern*(p).

If k is a possible outcome, we can define the probability mass function f of a Bernoulli distribution as:

$$f(k, p) = \begin{cases} p, & \text{if } k = 1, \\ 1 - p, & \text{if } k = 0. \end{cases}$$

Using the expression of the mean for discrete random variables, we obtain that $E[X] = p$ and

$$\text{Var}[X] = E[X^2] - E[X]^2 = E[X] - E[X]^2 = p - p^2 = p(1 - p) = pq.$$

As a note, this is just a particular case of the *Binomial distribution* with $n = 1$.

2.1.2 Gaussian Distribution

The Gaussian (or normal) distribution is used to represent real-valued random variables whose distributions are not known. Its importance relies in the fact that, using the *central limit theorem*, we can assume that the average of many samples of a random variable with finite mean and variance is a random variable whose distribution converges to a normal distribution as the number of samples increases.

Definition 2.1.1. We say that the real valued random variable X follows a *normal distribution* of parameters $\mu, \sigma \in \mathbb{R}$ if, and only if, its probability density function exists and it is determined by

$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}, \quad (1)$$

where μ is the mean and σ is its standard deviation. We denote this normal distribution as $X \sim \mathcal{N}(\mu, \sigma)$.

The particular case where $\mu = 0$ and $\sigma = 1$ is widely used in statistics. In this case, the density function is simpler:

$$f(x) = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}x^2}.$$

A remarkable property of these distributions is that, if $f : \mathbb{R} \rightarrow \mathbb{R}$ is a real-valued function defined as $f(x) = ax + b$, then $f(X) \sim \mathcal{N}(a\mu + b, |a|\sigma)$.

In the same way that we extended random variables to random vectors, we can extend the normal distribution to a multivariate random distribution.

Definition 2.1.2. We say that a random vector $\mathbf{X} = (X_1, \dots, X_n)$ follows a multivariate normal distributions of parameters $\mu \in \mathbb{R}^n$, $\Sigma \in \mathcal{M}_N(\mathbb{R})$ if, and only if, its probability density function is:

$$f(x) = \frac{1}{\sqrt{\det(2\pi\Sigma)}} e^{-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)}.$$

It is denoted $X \sim \mathcal{N}(\mu, \Sigma)$. In this case, μ is the mean vector of the distribution and Σ denotes the covariance matrix.

3

STATISTICAL INFERENCE

Statistical inference is the process of deducing properties of an underlying distribution by analyzing the data that it is available. With this purpose, techniques like deriving estimates and testing hypotheses are used.

Inferential statistics are usually contrasted with descriptive statistics, which are only concerned with properties of the observed data. The difference between these two is that in inferential statistics, we assume that the data comes from a larger population that we would like to know.

In *machine learning*, subject that concerns us the most, the term inference is sometimes used to mean *make a prediction by evaluating an already trained model*, and in this context, inferring properties of the model is referred as *training or learning*.

3.1 PARAMETRIC MODELING

In the following chapters, we will be trying to estimate density functions in a dataset. To do this we will be using *parametric models*. We say that a *parametric model*, $P_\theta(x)$, is a family of density functions that can be described using a finite numbers of parameters θ . We can get to the concept of *log-likelihood* now.

Definition 3.1.1. The *likelihood* $\mathcal{L}(\theta|x)$ of a parameter set θ is a function that measures how plausible is θ , given an observed point x in the dataset \mathcal{D} . It is defined as the value of the density function parametrized by θ at x . That is:

$$\mathcal{L}(\theta|x) = P_\theta(x).$$

In a finite dataset \mathcal{D} consisting of independent observations, we can write:

$$\mathcal{L}(\theta|X) = \prod_{x \in \mathcal{D}} P_\theta(x).$$

In practice, it is often convenient to work with the natural logarithm of the likelihood function.

Definition 3.1.2. Let \mathcal{D} be a dataset of independent observations and θ a set of parameters. Then, we define the *log-likelihood* ℓ as the sum of the logarithms of the evaluations of p_θ in each x in the dataset. That is:

$$\ell(\theta|X) = \sum_{x \in \mathcal{D}} \log P_\theta(x).$$

Since the logarithm is a monotonic function, the maximum of \mathcal{L} and the maximum of its logarithm will occur at the same θ . Our goal would be

to find the optimal value $\hat{\theta}$ that maximizes the likelihood of observing the dataset \mathcal{D} . We get to the following definition:

Definition 3.1.3. We say that $\hat{\theta} = \hat{\theta}(\mathcal{D})$ is a *maximum likelihood estimator* (MLE) for θ if

$$\hat{\theta} \in \arg \max_{\theta} \mathcal{L}(\theta|\mathcal{D})$$

for every observation \mathcal{D} .

Usually, we seek for likelihood functions that are differentiable, so the derivative test for determining maxima can be applied. Sometimes, the first-order conditions of the likelihood function can be solved explicitly, like in the case of the ordinary least squares estimator which maximizes the likelihood of a linear regression model. However, most of the times, we have to make use of numerical methods to be able to find the maximum of the likelihood function.

There is another concept related to the probability and the set of parameters that the distribution takes:

Definition 3.1.4. The *prior probability* is the probability distribution that it is believed to exist before evidence is taken into account.

The *posterior probability* is the probability of the parameters θ given the sampled data X , that is, $P(\theta, X)$.

The relation with the likelihood function $P(X|\theta)$ is that, given a prior belief that a p.d.f. is $P(\theta)$ and observations x have a likelihood $P(x|\theta)$, the posterior probability is defined using the prior probability as follows:

$$P(\theta|x) = \frac{P(x|\theta)}{P(x)} P(\theta),$$

where we have simply used Bayes' theorem.

3.2 MINIMAL SUFFICIENT STATISTICS

In parametric modeling, the goal was to determine the density function under a distribution. Another interesting task can be determining specific parameters or quantities related to a distribution, given a sample $X = (x_1, \dots, x_n)$.

Definition 3.2.1. Let (Ω, \mathcal{A}) be a measurable space where \mathcal{A} contains all singletons. A statistic is a measurable function of the data, that is: $T : X \rightarrow \Omega$ where T is measurable.

Remark 1. A statistic is also a random variable.

However, not all statistics will provide useful information for the statistical inference problem, since almost anything can be a statistic. We would like to find statistics that provide relevant information.

Definition 3.2.2. Let $X \sim P_{\theta}$. Then, the statistic $T(X) = T : (\Omega, \mathcal{A}) \rightarrow (\mathbb{T}, \mathcal{B})$, is sufficient for a family of parameters $\{P_{\theta} : \theta \in \Theta\}$ if the

conditional distribution of X , given $T = t$, is independent of θ .

Example 3. The simplest example of a sufficient statistic is the mean μ of a gaussian distribution with known variance. Oppositely, the *median* of an arbitrary distribution is not sufficient for the mean since, even if the median of the sample is known, more information about the mean of the population can be obtained from the mean of the sample itself.

Although it will not be shown in this document, sufficient statistics are not unique. In fact, if T is sufficient, $\psi(T)$ is sufficient for any bijective mapping ψ . It would be interesting to find a sufficient statistic T that is *the smallest* of them.

Definition 3.2.3. A sufficient statistic T is minimal if, for every sufficient statistic U , there exists a mapping f such that $T(x) = f(U(x))$ for any $x \in \Omega$.

4

NOISE CONTRASTIVE ESTIMATION

In this chapter, we will present the mathematical fundamentals that inspired the methods that nowadays obtain the best results in the field of unsupervised representation learning, which is a very important part of the new advances in machine learning. All these concepts will be clarified later in this document.

4.1 LOGISTIC REGRESSION

The logistic model is used to model the probability of a certain *discriminative* class or event existing, such as pass/fail or dead/alive. It can be extended to a multiple class problem, but we will focus on the binary problem.

Definition 4.1.1. The *odds* is a likelihood measure of a particular outcome. It is calculated as the ratio of the number of events that produce that outcome to the numbers of events that produce a different outcome.

The *log-odds* is the logarithm of the odds.

Usually, logistic models make use of an indicator variable (which is a random variable that takes values 0 or 1). The log-odds of the value labeled as 1 is a combination of one or more independent variables. Let us see how this idea is formalized and how we use this.

Let X be the input variables to a model and Y a label with values $\{0, 1\}$ as we said before. The labels are distributed following a $Bern(p_i)$ distribution. That is:

$$y_i | x_{1,i}, \dots, x_{n,i} \sim Bern(p_i)$$

We would like to have the conditional distribution $P(Y|X)$, that is, the distribution that models the probability of our label Y given the input X .

Since Y is an indicator variable, we have that $P(Y = 1) = E[Y]$ and, similarly, $P(Y = 1|X = x) = E[Y|X = x]$. Knowing this, there exists methods for estimating the regression function for the indicator variable using a smoother (Cosma, 2021). However this could be a bad idea because we can not assure that any smoother would return a value between 0 and 1. However, we can find a way to go around this problem.

Let us now assume $P(Y = 1|X = x) = p(x; \theta)$ for some function p parametrized by θ and that the observations x_i are independent from each other. Then, the conditional likelihood function is:

$$\prod_{i=1}^n P(Y = y_i | X = x_i) = \prod_{i=1}^n p(x_i, \theta)^{y_i} (1 - p(x_i; \theta))^{(1-y_i)}. \quad (2)$$

We would like to find the $p(x; \theta)$, constraining it to be a probability, that is, to be between 0 and 1. To achieve this, we make use of the logistic transformation, that is, we consider:

$$\log \frac{p}{1-p}.$$

Let us set in the case of a single explanatory variable x . An idea could be to make this model a linear regressor. That is, solving

$$\log \frac{p(x; \theta)}{1-p(x; \theta)} = \beta_0 + \beta x,$$

which gives us the solution:

$$p(x; \beta) = \frac{\exp(\beta_0 + \beta x)}{1 + \exp(\beta_0 + \beta x)} = \frac{1}{1 + \exp(-\beta_0 + \beta x)},$$

which is called the logistic function. We have constructed the following definition:

Definition 4.1.2. The *standard logistic function* $\sigma : \mathbb{R} \rightarrow (0, 1)$ is defined as follows:

$$\sigma(t) = \frac{1}{1 + e^{-t}}$$

4.1.1 Fitting a Logistic Regression model

Consider now a generalized linear model, which we assume parametrized by a set of parameters θ :

$$p(X, \theta) = \frac{1}{1 + e^{-\theta^T X}} = P(Y = 1|X; \theta).$$

WhereTherefore,

$$P(Y = 0|X; \theta) = 1 - p(X, \theta).$$

Since Y is an indicator variable, we have that

$$P(y|X; \theta) = p(X, \theta)^y (1 - p(X, \theta))^{(1-y)}.$$

As we have seen in Equation (2), we have that the likelihood function in:

$$\begin{aligned} L(\theta|y; x) &= P(Y|X; \theta) \\ &= \prod_i P(y_i|x_i; \theta) \\ &= \prod_i p(x_i; \theta)^{y_i} (1 - p(x_i, \theta))^{(1-y_i)}. \end{aligned}$$

If we consider now the log-likelihood function, the logarithm turns products into sums, so we obtain:

$$N^{-1} \log L(\theta|y; x) = N^{-1} \sum_{i=1} \log P(y_i|x_i; \theta),$$

which can be maximized using optimization techniques such as *gradient descent*, a very common technique in machine learning.

Remark 2. If we assume that the pairs (x_i, y_i) are drawn uniformly from the underlying distribution, if we consider the limit in N , we obtain:

$$\begin{aligned}
& \lim_{N \rightarrow +\infty} N^{-1} \sum_{i=1}^N \log P(y_i | x_i; \theta) \\
&= \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} P(X = x, Y = y) \log P(Y = y | X = x; \theta) \\
&= \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} P(X = x, Y = y) \left(-\log \frac{P(Y = y | X = x)}{P(Y = y | X = x; \theta)} + \log P(Y = y | X = x) \right) \\
&= -D_{KL}(Y || Y_\theta) - H(Y | X),
\end{aligned}$$

where $H(Y | X)$ is the conditional entropy¹, and in the first equality we have use the *Law of the Large Numbers*. This means that, intuitively, if we maximize the log-likelihood of a model, we are minimizing the KL-divergence of our model from the maximal entropy distribution, which is sort of *searching for the model that makes the fewest assumptions in its parameters*.

4.2 MOTIVATION

Our problem now is to estimate the probability density function (p.d.f.) of some observed data. However, the data that we have available has been sampled in a specific way, let us present see how it is done.

Consider that a sample $X = \{x_1, \dots, x_{T_d}\}$ has been observed, where each element has been sampled from a random variable. This data follows an unknown p.d.f. P_d , that we assume to belong to a parametrized family of functions, that is

$$P_d \in \{P_m(\cdot; \theta)\}_{\theta},$$

where θ is a vector of parameters. In other words

$$P_d(\cdot) = P_m(\cdot; \theta^*) \quad \text{for some } \theta^*.$$

Our problem now will be to find the θ^* that matches the distribution.

Any estimate $\hat{\theta}$ must meet the constraints that a normalized p.d.f. should satisfy.

Definition 4.2.1. Let $\hat{\theta}$ be a set of parameters and $P_m(\cdot; \hat{\theta})$ a probability density function with parameters $\hat{\theta}$. We say that $P_m(\cdot; \hat{\theta})$ is normalized if it satisfies:

$$\int P_m(u; \hat{\theta}) du = 1, \quad \text{and} \quad P_m(\cdot; \hat{\theta}) \geq 0.$$

If the constraints are satisfied for any θ in the set of parameters, we say that the model is normalized, and then we can use the maximum likelihood principle to estimate θ .

Consider now some noisy data Y . Let us assume that the noisy data Y is an i.i.d. sample $\{y_1, \dots, y_{T_n}\}$ of a random variable with p.d.f. P_n . The ratio

¹ A definition and explanation of what the entropy is will be given in the next Part.

P_d/P_n of the density functions that generate X and Y respectively, can give us a relative description of the data X . That means, if we know the noisy distribution P_n , then we can estimate the data distribution P_d using the ratio that we have just mentioned. Shortly: if we know the differences between X and Y and we know the properties of Y , we can deduce the properties of X .

Our goal in *Noise Contrastive Estimation* (NCE) is to be able to discriminate between elements that have been sampled from the original data distribution P_d and elements that have been sampled from the noise distribution. In order to discriminate between elements of X and Y , it is needed to compare their properties. We will show that we can provide a relative description of X in the form of an estimate of the ratio P_d/P_n .

4.3 FORMALIZATION OF THE NCE PROBLEM

Let $U = \{u_1, \dots, u_{T_d+T_n}\}$ be the union of the sets X and Y that we mentioned before. We assign to each u_t a binary class label, depending if it belongs to the original data X or the noise data Y :

$$C_t(u_t) = \begin{cases} 1 & \text{if } u_t \in X \\ 0 & \text{if } u_t \in Y \end{cases}$$

With this labels, will now make use of logistic regression, where the posterior probabilities of the classes given the data are estimated. We know that the distribution of the data P_d is unknown, we want to model the class conditional probability $P(.|C = 1)$ with $P_m(.,\theta)$. Note that θ may include a parameter for the normalization of the model, if it is not normalized. Hence, we have:

$$P(u|C = 1, \theta) = P_m(u; \theta), \quad P(u|C = 0) = P_n(u).$$

Furthermore, since we know the numbers of examples of each of the X and Y sets, we also know that the prior probabilities are:

$$P(C = 1) = \frac{T_d}{T_d + T_n}, \quad P(C = 0) = \frac{T_n}{T_d + T_n}.$$

We also have to consider that $P(u)$ can be decomposed as:

$$P(u) = P(C = 1)P_m(u; \theta) + P(C = 0)P_n(u) = \frac{T_d}{T_d + T_n}P_m(u; \theta) + \frac{T_n}{T_d + T_n}$$

Hence, if we rename the quotient of the probability of the negative class by the probability of the positive class $v = P(C = 0)/P(C = 1) = T_n/T_d$, the posterior probabilities for the positive class $C = 1$ is:

$$\begin{aligned} P(C = 1|u; \theta) &= \frac{P(u|C = 1; \theta)P(C = 1)}{P(u)} \\ &= \frac{P_m(u; \theta)P(C = 1)}{P(C = 1)P_m(u; \theta) + P(C = 0)P_n(u)} \\ &= \frac{P_m(u; \theta)}{P_m(u; \theta) + vP_n(u)} \end{aligned}$$

Where in the first equality we have used Bayes' rule. Similarly, the posterior probability for the negative class $C = 0$ is:

$$\begin{aligned}
 P(C = 0|u; \theta) &= \frac{P(u|C = 0; \theta)P(C = 0)}{P(u)} \\
 &= \frac{P_n(u)P(C = 0)}{P(C = 1)P_m(u; \theta) + P(C = 0)P_n(u)} \\
 &= \frac{P_n(u)P(C = 0)}{P(C = 1) \left(P_m(u; \theta) + \frac{P(C=0)}{P(C=1)} P_n(u) \right)} \\
 &= \frac{P_n(u) \frac{P(C=0)}{P(C=1)}}{P_m(u; \theta) + \frac{P(C=0)}{P(C=1)} P_n(u)} \\
 &= \frac{v P_n(u)}{P_m(u; \theta) + v P_n(u)}.
 \end{aligned}$$

We are in the conditions now to see how we use the logistic regression to obtain its conditional log-likelihood form. Let $G(., \theta)$ be the log ratio between $P_m(., \theta)$ and $P_n(., \theta)$:

$$G(u; \theta) = \log \frac{P_m(u; \theta)}{P_n(u)} = \log P_m(u; \theta) - \log P_n(u). \quad (3)$$

Using this log-ratio, we obtain the following proposition. It is not a super interesting result, but we want to remark it since it connects the context with the logistic function.

Proposition 2. *Under the conditions that we have presented until this point, and naming $h(u; \theta) := P(C = 1|u; \theta)$, we have that*

$$h(u; \theta) = r_v(G(u; \theta)),$$

where

$$r_v(u) = \frac{1}{1 + v \exp(-u)}, \quad (4)$$

is the logistic function parametrized by v .

Proof. Firstly, it is easy to see that

$$\exp(-G(u; \theta)) = \exp(-\log P_m(u; \theta) + \log P_n(u)) = \frac{P_n(u)}{P_m(u; \theta)}.$$

Using this, the proof is almost pretty straightforward:

$$\begin{aligned}
 h(u; \theta) &= r_v(G(u; \theta)) \\
 &= \frac{1}{1 + v \exp(-G(u; \theta))} \\
 &= \frac{1}{1 + v \frac{P_n(u)}{P_m(u; \theta)}} \\
 &= \frac{1}{\frac{1}{P_m(u; \theta)} (P_m(u; \theta) + v P_n(u))} \\
 &= \frac{P_m(u; \theta)}{P_m(u; \theta) + v P_n(u)} \\
 &= P(C = 1 | u; \theta),
 \end{aligned}$$

as we wanted to proof. \square

Since the class labels C_t are assumed Bernoulli distributed and independent, the conditional log-likelihood has the form:

$$\ell(\theta) = \sum_{t=1}^{T_d+T_n} C_t \log P(C_t = 1 | u_t; \theta) + (1 - C_t) \log P(C_t = 0 | u_t; \theta). \quad (5)$$

Now, in the terms of the sumatory such that u_t in X , we have that $u_t = x_t$ and, hence, $P(C_t = 0 | x_t; \theta) = 0$. Because of this, we obtain that the term that adds to the sum in those specific t adopt the form:

$$1 \cdot \log P(C_t = 1 | u_t; \theta) = \log h(x_t; \theta).$$

Using the same argument for t such that $u_t \in Y$, we obtain the following form of the log-likelihood in (5):

$$\ell(\theta) = \sum_{t=1}^{T_d} \log[h(x_t; \theta)] + \sum_{t=1}^{T_n} \log[1 - h(y_t, \theta)]. \quad (6)$$

Now, optimizing $\ell(\theta)$ with respect to θ leads to an estimate $G(\cdot; \hat{\theta})$ of the log-ratio $\log(P_d/P_n)$, so we get an approximate description of X relative to Y by optimizing (6).

Remark 3. If we consider $-\ell(\theta)$, this is known as the *cross entropy function*.

Remark 4. Here, we have achieved the estimation of a p.d.f. , which is an unsupervised (not labeled data) learning problem, using logistic regression, which is supervised learning (labeled data).

Now, if we consider $P_m^0(\cdot; \alpha)$ an unnormalized (does not integrate 1) model, we can add a normalization parameter to it in order to normalize it. We can consider

$$\log P_m(\cdot; \theta) = \log P_m^0(\cdot; \alpha) + c, \quad \text{with } \theta = (\alpha, c).$$

With this model, a new estimator is defined. We consider X as before and Y an artificially generated set with $|Y| = T_n = vT_d$ independent observations

extracted from P_n , known. Then, the estimator is defined to be the argument $\hat{\theta}_T$ which maximizes

$$J_T(\theta) = \frac{1}{T_d} \left\{ \sum_{t=1}^{T_d} \log[h(x_t; \theta)] + \sum_{t=1}^{T_n} \log[1 - h(y_t; \theta)] \right\}.$$

We have to remark that in this case, we have fixed ν before T_n , so T_n will increase as T_d increases. Now, using the weak law of large numbers, $J_T(\theta) \rightarrow J$ in probability, where

$$J(\theta) = E [\log[h(x; \theta)]] + \nu E [\log[1 - h(y; \theta)]].$$

Let us rename some terms before announcing a theorem. We want to see J as a function of $\log P_m(\cdot; \theta)$ instead of only θ . In order to do this, let $f_m(\cdot) = \log P_m(\cdot; \theta)$, and consider

$$\tilde{J}(f_m) = E \{ \log[r_\nu(f_m(x) - \log P_n(x))] \} + \nu E \{ \log[1 - r_\nu(f_m(y) - \log P_n(y))] \}.$$

The following theorem states that the probability density function P_d of the data can be found by maximizing \tilde{J} , that is, learning a nonparametric classifier in *infinite data*.

Theorem 4.1. *The objective $\tilde{J}(f_m)$ achieves a maximum at $f_m = \log P_d$. Furthermore, there are not other extrema if the noise density P_n is chosen such that it is nonzero whenever P_d is nonzero.*

The proof of this theorem is beyond the scope of this work and will not be shown here. It can be found in [Gutmann & Hyvarinen \(n.d.\)](#).

Part II

INFORMATION THEORY

5 | ENTROPY

5.1 DEFINITION

Intuitively, if x is a datapoint in a dataset $\mathcal{D} \subset \mathbb{R}^d$, a *representation* of x is a vector $r \in \mathbb{R}^n$ (usually, $n \leq d$), that shares information with the datapoint x . Representations are very often used in Machine Learning (ML).

Obtaining good representations of data is one of the most important tasks in ML. Recently, it has been discovered that maximizing the *mutual information* between two elements in our data can give us good representations for our data. In this section, *information theory* notions will be presented, in order to use them in our ML models. This will provide a theoretical solid base for the notions explained later.

The *mutual information* concept is based on the *Shannon entropy*, which we will introduce first, along with some basic properties of it. The Shannon entropy is a way of measuring the uncertainty in a random variable. Given an event $\mathcal{A} \in \mathcal{A}$, P a probability measure and $P[\mathcal{A}]$ the probability of \mathcal{A} , we can affirm that

$$\log \frac{1}{P[\mathcal{A}]}$$

describes “how surprising is that \mathcal{A} occurs”. Clearly, then,

- If $P[\mathcal{A}] = 1$, then we have $\log 1 = 0$, which would indicate that it is not a surprise that \mathcal{A} occurred.
- On the other hand, if $P[\mathcal{A}] \rightarrow 0$, then we have $\log(+\infty)$, which will surely be a big number indicating great surprise that \mathcal{A} occurred.

With this motivation, we get to the following definition.

Definition 5.1.1. Let X be a discrete random variable with image \mathcal{X} . The *Shannon entropy*, or simply *entropy* $H(X)$ of X is defined as:

$$H(X) = E_X \left[\log \frac{1}{P_X(X)} \right] = \sum_{x \in \mathcal{X}} P_X(x) \log \frac{1}{P_X(x)}.$$

The *entropy* can trivially be expressed as:

$$H(X) = - \sum_{x \in \mathcal{X}} P_X(x) \log P_X(x).$$

This simple example, (Cover & Thomas, 1991), even though it is very simple, it is very illustrative for our definition:

Example 4. Let $X \sim \text{Bern}(p)$. Then, the entropy of X is:

$$H(X) = -p \log p - (1 - p) \log(1 - p) = H(p),$$

since H only depends on p . In Fig. 4 we can see a representation of this function. We appreciate that in this case, H is concave and equals 0 if $p \in \{0, 1\}$, which are the values of p that give us no uncertainty. The maximum uncertainty is obtained when $p = \frac{1}{2}$, where we do not know what to expect as an outcome from our random variable X .

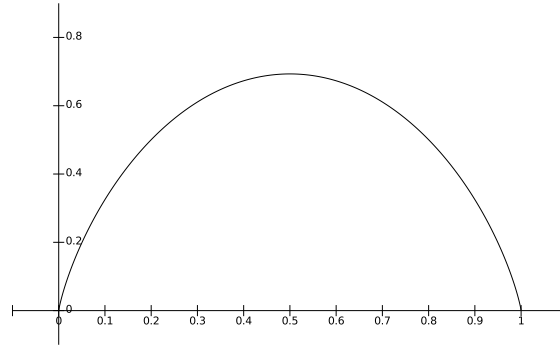


Figure 4: Representation of $H(p)$ in the example 4.

It can also be proven that, in general, the entropy is concave.

5.2 PROPERTIES OF THE ENTROPY. CONDITIONAL ENTROPY

There are some properties of the *entropy* that must be remarked, since they will extend to properties of the mutual information.

Proposition 3. *Let X be a random variable with image \mathcal{X} . If $|\mathcal{X}|$ is the cardinal of \mathcal{X} , then*

$$0 \leq H(X) \leq \log(|\mathcal{X}|).$$

Proof. Since $\log y$ is concave on \mathbb{R}^+ , by Jensen's inequality (see Appendix A, Prop. 9), we obtain:

$$H(X) = - \sum_{x \in \mathcal{X}} P_X(x) \log P_X(x) \leq \log \left(\sum_{x \in \mathcal{X}} 1 \right) = \log(|\mathcal{X}|).$$

For the lower bound we see that, since $P_X(x) \in [0, 1]$ for all $x \in \mathcal{X}$ then $\log P_X(x) \leq 0 \quad \forall x \in \mathcal{X}$. Hence , $-P_X(x) \log P_X(x) \geq 0$ for all $x \in \mathcal{X}$, so $H(X) \geq 0$. \square

We can also see that the equality on the left holds if , and only if , exists x in \mathcal{X} such that its probability is exactly one, that is $P_X(x) = 1$. The right equality holds if and only if , for all $x \in \mathcal{X}$, its probability is $P_X(x) = \frac{1}{|\mathcal{X}|}$.

Conditional entropy

We have already said that entropy measures how surprising is that an event occurs. Usually, we will be looking at two random variables and it would be interesting to see how likely is that one of them, say $X(x)$, occurred,

if we already know that $Y(y)$ occurred. This leads us to the definition of *conditional entropy*. Let us see a simpler case first:

Let A be an event, and X a random variable. The conditional probability $P_{X|A}$ defines the entropy of X conditioned to A :

$$H(X|A) = \sum_{x \in \mathcal{X}} P_{X|A}(x) \log \frac{1}{P_{X|A}(x)}.$$

If Y is another random variable and \mathcal{Y} is its image, intuitively we can sum the conditional entropy of an event with all the events in \mathcal{Y} , and this way we obtain the conditional entropy of X given Y .

Definition 5.2.1 (Conditional Entropy). Let X, Y be random variables with images \mathcal{X}, \mathcal{Y} . The *conditional entropy* $H(X|Y)$ is defined as:

$$\begin{aligned} H(X|Y) &:= \sum_{y \in \mathcal{Y}} P_Y(y) H(X|Y = y) \\ &= \sum_{y \in \mathcal{Y}} P_Y(y) \sum_{x \in \mathcal{X}} P_{X|Y}(x|y) \log \frac{1}{P_{X|Y}(x|y)} \\ &= \sum_{x \in \mathcal{X}, y \in \mathcal{Y}} P_{XY}(x, y) \log \frac{P_Y(y)}{P_{XY}(x, y)}. \end{aligned}$$

The interpretation of the conditional entropy is simple: the uncertainty in X when Y is given. Since we know about an event that has occurred (Y), intuitively the conditional entropy, or the uncertainty of X occurring given that Y has occurred, will be lesser than the entropy of X , since we already have some information about what is happening. We can prove this:

Proposition 4. Let X, Y be random variables with images \mathcal{X}, \mathcal{Y} . Then:

$$0 \leq H(X|Y) \leq H(X).$$

Proof. The inequality on the left was proved on Proposition 3. The characterization of when $H(X|Y) = 0$ was also mentioned after it. Let us look at the inequality on the right. Note that restricting to the (x, y) where $P_{XY}(x, y) > 0$ and using the definition of the conditional probability we have:

$$\begin{aligned} H(X|Y) &= \sum_{y \in \mathcal{Y}} P_Y(y) \sum_{x \in \mathcal{X}} P_{X|Y}(x|y) \log \frac{1}{P_{X|Y}(x|y)} \\ &= \sum_{x \in \mathcal{X}, y \in \mathcal{Y}} P_Y(y) P_{X|Y}(x, y) \log \frac{P_Y(y)}{P_{XY}(x, y)} \\ &= \sum_{x \in \mathcal{X}, y \in \mathcal{Y}} P_{XY}(x, y) \log \frac{P_Y(y)}{P_{XY}(x, y)}, \end{aligned}$$

and

$$H(X) = \sum_x P_X(x) \log \frac{1}{P_X(x)} = \sum_{x, y} P_{XY}(x, y) \log \frac{1}{P_X(x)}.$$

Hence,

$$\begin{aligned} H(X|Y) - H(X) &= \sum_{x,y} P_{XY}(x,y) \left(\log \frac{P_Y(y)}{P_{XY}(x,y)} - \log \frac{1}{P_X(x)} \right) \\ &= \sum_{x,y} P_{XY} \log \frac{P_Y(y)P_X(x)}{P_{XY}(x,y)}. \end{aligned} \quad (7)$$

So, using Jensen's inequality, we obtain:

$$\begin{aligned} \sum_{x,y} P_{XY} \log \frac{P_Y(y)P_X(x)}{P_{XY}(x,y)} &\leq \log \left(\sum_{x,y} \frac{\cancel{P_{XY}(x,y)} P_Y(y)P_X(x)}{\cancel{P_{XY}(x,y)}} \right) \\ &= \log \left(\left(\sum_x P_X(x) \right) \left(\sum_y P_Y(y) \right) \right) = \log 1 = 0, \end{aligned}$$

and this leads us to:

$$H(X|Y) - H(X) \leq 0 \quad \text{then} \quad H(X|Y) \leq H(X) \quad (8)$$

as we wanted. \square

It must be noted that the inequality the state of the proposition,

$$0 \leq H(X|Y) \leq H(X),$$

in the inequality of the left, equality holds if, and only if, $P_{XY}(x,y) = P_X(x)P_Y(y)$ for all (x,y) with $P_{XY}(x,y) > 0$, as it is said in Jensen's inequality. For the inequality on the right, equality holds if and only if $P_{XY}(x,y) = 0$, which implies $P_X(x)P_Y(y) = 0$ for any $x \in \mathcal{X}$, $y \in \mathcal{Y}$. It follows that $H(X|Y) = H(X)$ if and only if $P_{XY}(x,y) = P_X(x)P_Y(y)$ for all $(x,y) \in \mathcal{X} \times \mathcal{Y}$

6

MUTUAL INFORMATION

Using the entropy of a random variable we can directly state the definition of *mutual information* as follows:

Definition 6.1. Let X, Z be random variables. The *mutual information* (MI) between X and Z is expressed as the difference between the entropy of X and the conditional entropy of X and Z , that is:

$$I(X, Z) := H(X) - H(X|Z).$$

Since the entropy of the random variable $H(X)$ explains the uncertainty of X occurring, the intuitive idea of the MI is to determine the decrease of uncertainty of X occurring when we already know that Z has occurred. We also have to note that, using the definition of the *entropy* and the same argument that we used to obtain the expression in Eq. 7, we can rewrite the MI it follows:

$$\begin{aligned} I(X, Z) &= \sum_{x \in \mathcal{X}} P_X(x) \log \frac{1}{P(x)} - \sum_{x \in \mathcal{X}, z \in \mathcal{Z}} P_{XZ}(x, z) \log \frac{P_Z(x)}{P_{XZ}(x, z)} \\ &= \sum_{x, z} P_{XZ} \log \frac{P_{XZ}(x, z)}{P_Z(z)P_X(x)} \end{aligned} \quad (9)$$

and if we compare it to the formula of the KL-Divergence, we obtain:

$$I(X, Z) = \sum_{x, z} P_{XZ} \log \frac{P_{XZ}(x, z)}{P_Z(z)P_X(x)} = D_{KL}(P_{XZ} || P_X P_Z),$$

so we have obtained an expression of the mutual information using the *Kullback-Leibler* divergence. This provides with the following immediate consequences:

Corollary 1. (i) *Mutual information is non-negative. That is : $I(X, Z) \geq 0$.*
(ii) *If X, Z are random variables, then its mutual information equals zero if, and only if, they are independent.*
(iii) *Mutual information is symmetric. That is: $I(X, Z) = I(Z, X)$.*

Proof. (i) This is trivial using Prop 4 and the definition of the mutual information.

(ii) We can use the KL-Divergence formulation to see that since

$$D_{KL}(P_{XZ} || P_X P_Z) = 0 \implies P_{XZ} = P_X P_Z,$$

almost everywhere then X and Z are independent.

(iii) It is a consequence of the fact that $P_{XZ} = P_{ZX}$ and $P_X P_Z = P_Z P_X$.

□

Later in this document, we will have some sort of random variable X and would like it to maintain the mutual information with itself after being applied a function. The following proposition will be useful:

Proposition 5. *Let X, Z be random variables. Then, $I(X, Z)$ is invariant under homeomorphism.*

Proof. Let $\phi(x)$ be an homeomorphism, i.e., a continuous, monotonic function with $\phi^{-1}(x)$ also continuous and monotonic. Let X be a random variable and Y another one such $y = \phi(x)$ if $x = X(\omega)$ for some $\omega \in \Omega$. Then, if S is a particular subset we have

$$P(Y \in S) = \int_S P_Y(y) dy = \int_{\phi^{-1}(S)} P_X(x) dx \stackrel{(1)}{=} \int_S P_X(\phi^{-1}(y)) \left| \frac{d\phi^{-1}}{dy} \right| dy,$$

where in (1) we have changed from x to y . Hence,

$$P_Y(y) = P_X(\phi^{-1}(y)) \left| \frac{d\phi^{-1}}{dy} \right|.$$

As a consequence of this, $I(X, Z) = I(\phi(X), Z)$ for any homeomorphism ϕ . By symmetry, the same holds for Z .

□

Remark 5. We can set a connection between the mutual information and sufficient statistics. Let $T(X)$ be a statistic. We say that $T(X)$ is sufficient for θ if its mutual information with θ equals the mutual information between X and θ , that is:

$$I(\theta, X) = I(\theta, T(X)).$$

This means that sufficient statistics preserve mutual information and conversely.

6.1 LOWER BOUNDS ON MUTUAL INFORMATION

Although mutual information seems like a relatively intuitive concept, it is most of the times extremely hard to compute it in real life problems in which the distributions $P(x, z), P(x), P(z)$ are not known.

Example 5. Let x represent an image of size $n \times m$ pixels. Then, the dimension of the single image is $n \cdot m \cdot 3$, for RGB color channels. In these cases, there is no easy way of calculating $P(x)$.

Due to this problem related to the *Curse of Dimensionality*, we can try to compute lower bounds of it that are generally easier to calculate. We will now expose two general lower bounds, and we will focus on a third one that will be explained later in this work.

6.1.1 Variational Lower Bound

Using the expression of the mutual information in terms of entropy, $I(x, z) = H(z) - H(z|x)$, we can give a lower bound on $I(x, z)$ as a function of a probability distribution $Q_\theta(z|x)$.

Proposition 6. *Let X, Z be random variables and $Q_\theta(z|x)$ be an arbitrary probability distribution. Then,*

$$I(x, z) \geq H(z) + E_{P_X} \left[E_{P_{X|Z}} [\log Q_\theta(z|x)] \right]$$

Proof. Recalling that

$$H(z|x) = -E_{P_{XZ}} [\log P(x, z) - \log P(x)],$$

and that

$$\begin{aligned} E_{P(x,z)} \left[\log \frac{P(x,z)}{P(x)} \right] &= \sum_{x,z} P(x,z) \log \frac{P(x,z)}{P(x)} \\ &= \sum_{x,z} P(x) P(z|x) \log P(z|x) = \sum_{x,z} P(x) E_{P(z|x)} [\log P(z|x)] \\ &= E_{P(x)} \left[E_{P(z|x)} [\log P(z|x)] \right], \end{aligned}$$

we only have to use the definition of the conditional probability to see that:

$$\begin{aligned} I(x, z) &= H(z) - H(z|x) \\ &= H(z) + E_{P(x,z)} \left[\log \frac{P(x,z)}{P(x)} \right] \\ &= H(z) + E_{P(x)} \left[E_{P(z|x)} [\log P(z|x)] \right] \\ &= H(z) + E_{P(x)} \left[E_{P(z|x)} \left[\log \frac{P(z|x)}{Q_\theta(z|x)} \right] + E_{P(z|x)} [\log Q_\theta(z|x)] \right] \\ &= H(z) + E_{P(x)} \left[\underbrace{D_{KL}(P(z|x) || Q_\theta(z|x))}_{\geq 0} + E_{P(z|x)} [\log Q_\theta(z|x)] \right] \\ &\geq H(z) + E_{P(x)} \left[E_{P(z|x)} [\log Q_\theta(z|x)] \right]. \end{aligned}$$

We have taken advantage of the non-negativity of the KL-Divergence. \square

Using this bound, and combining this theoretical knowledge with machine learning methods, such as *backpropagation*, we can make Q_θ be a neural network and maximize this lower bound.

6.1.2 Donsker-Varadhan Representation lower bound

We can also give a lower bound on the mutual information using its KL-Divergence formulation. Firstly, we have to

Theorem 6.1.1 (Donsker-Varadhan). *The KL divergence admits the following dual representation:*

$$D_{KL}(P||Q) = \sup_T E_P[T] - \log E_Q[e^T],$$

where the supremum is taken over all functions $T : \Omega \rightarrow \mathbb{R}$ such that both expectations exist.

Proof. Let T be a given function. We can consider the Gibbs distribution, which is defined by

$$G(x) = \frac{1}{Z} e^{T(x)} Q(x),$$

where Z is a normalization term defined by $Z = E_{Q(x)}[e^T]$.

□

Using this representation, we reach this lower bound as a corollary.

Corollary 2. *Let \mathcal{F} be any class of functions $T : \Omega \rightarrow \mathbb{R}$ satisfying the integrability constraints of the theorem. Then,*

$$I(P, Q) = D_{KL}(P||Q) \geq \sup_{T \in \mathcal{F}} E_P[T] - \log E_Q[e^T].$$

6.1.3 Contrastive Lower Bound

In chapter 4 we presented Noise Contrastive Estimation, that tried to discriminate between elements of two different sets. One was composed of data, X , and the other one was composed of noise, Y .

Let (x, z) be a data representation drawn from a distribution $P(x, z)$ and x' be some other data drawn from the distribution $P(x)$. Using NCE, we should be able to say that (x, z) was drawn from the distribution $P(x, z)$ (which was P_d in the NCE theory) while (x', z) was drawn from the product of the marginal distributions $P(x)P(z)$ (which was P_n in the explanation of NCE). Let h_θ be a model that helps us to do this discrimination, with parameters θ .

As we did before, we want to estimate the ratio P_d/P_n of the different distributions, in this case the ratio would be $P(x, z)/P(x)P(z)$. Let (x^*, z) be a pair drawn from $P(x, z)$ and $X = \{x^*, x_1, \dots, x_{N-1}\}$, where the rest of the $N-1$ points form pairs (x_j, z) drawn from $P(x)P(z)$ the product of the marginal distribution. We can rewrite the loss 6 in a simpler expression:

$$l(\theta) = -E_X \left[\log \frac{h_\theta(x^*, z)}{\sum_{x \in X} h_\theta(x, z)} \right]. \quad (10)$$

If we maximize this objective, h_θ learns to discriminate (x^*, z) from (x_j, z) for $1 \leq j < N$ and, thus, we are learning to estimate the ratio $P(x, z)/P(x)P(z)$. Let us see how maximizing $\ell(\theta)$ we are maximizing a lower bound for $I(x, z)$.

Proposition 7. Let $X = \{x^*, x_1, \dots, x_{N-1}\}$, where $x^* \sim P(x, z)$ and the rest of them were sampled from $P(x)P(z)$. Then,

$$I(x, z) \geq -\ell(\theta) + \log N \quad (11)$$

Proof. Firstly, using Bayes' rule, $P(x^*, z) = P(x^*|z)P(z)$. Hence, since h_θ estimates $P(x^*, z)/P(x)P(z)$, it also estimates

$$\frac{P(x^*, z)}{P(x)P(z)} = \frac{P(x^*|z)P(z)}{P(x)P(z)} = \frac{P(x^*|z)}{P(x)}.$$

Using the definition of the log-likelihood that we see in (10), forgetting the sign for the moment, we see that

$$\begin{aligned} E_X \left[\log \frac{h_\theta(x^*, z)}{\sum_{x \in X} h_\theta(x, z)} \right] &= E_X \left[\log \frac{h_\theta(x^*, z)}{h_\theta(x^*, z) + \sum_{j=1}^{N-1} h_\theta(x_j, z)} \right] \\ &\approx E_X \left[\log \frac{\frac{P(x^*|z)}{P(x)}}{\frac{P(x^*|z)}{P(x)} + \sum_{j=1}^{N-1} \frac{P(x_j|z)}{P(x)}} \right]. \end{aligned}$$

Now, using that $\log(a) = -\log(a^{-1})$,

$$\begin{aligned} E_X \left[\log \frac{\frac{P(x^*|z)}{P(x)}}{\frac{P(x^*|z)}{P(x)} + \sum_{j=1}^{N-1} \frac{P(x_j|z)}{P(x)}} \right] &= E_X \left[-\log \left(\frac{\frac{P(x^*|z)}{P(x)} + \sum_{j=1}^{N-1} \frac{P(x_j|z)}{P(x)}}{\frac{P(x^*|z)}{P(x)}} \right) \right] \\ &= E_X \left[-\log \left(1 + \frac{\sum_{j=1}^{N-1} \frac{P(x_j|z)}{P(x)}}{\frac{P(x^*|z)}{P(x)}} \right) \right] \\ &= E_X \left[-\log \left(1 + \frac{(N-1)E_{X-\{x^*\}} \left[\frac{P(x|z)}{P(x)} \right]}{\frac{P(x^*|z)}{P(x)}} \right) \right] \end{aligned}$$

Now, since $E_{X-\{x^*\}} \left[\frac{P(x|z)}{P(x)} \right] = \sum_{x_j \in X-\{x^*\}} P(x_j) \frac{P(x_j|z)}{P(x_j)} = 1$, then

$$E_X \left[-\log \left(1 + \frac{(N-1)E_{X-\{x^*\}} \left[\frac{P(x|z)}{P(x)} \right]}{\frac{P(x^*|z)}{P(x)}} \right) \right] = E_X \left[-\log \left(1 + \frac{(N-1)}{\frac{P(x^*|z)}{P(x)}} \right) \right].$$

Lastly, using that if $k > 0$, then $-\log a(k+1) \geq -\log(1+ak)$, we obtain:

$$\begin{aligned}
E_X \left[-\log \left(1 + \frac{(N-1)}{\frac{P(x^*|z)}{P(x)}} \right) \right] &= E_X \left[\log \frac{1}{1 + \frac{P(x^*)}{P(x^*|z)}(N-1)} \right] \\
&\leq E_X \left[\log \left(\frac{1}{\frac{P(x^*)}{P(x^*|z)}} \frac{1}{N} \right) \right] \\
&= E_X \left[\log \left(\frac{P(x^*|z)}{P(x^*)} \frac{1}{N} \right) \right] \\
&= E_X \left[\log \left(\frac{P(x^*|z)}{P(x^*)} \right) \right] - \log N \\
&\stackrel{(1)}{=} E_X \left[\log \left(\frac{P(x^*, z)}{P(x^*)P(z)} \right) \right] - \log N \\
&\stackrel{(2)}{=} I(x, z) - \log N,
\end{aligned}$$

where, in (1), we have use Bayes' rule again and in (2) we have used the definition of the MI that we found in equation 9. Looking at the first and last equations used in this proof, and seeing that we have \leq in the middle of the chain of equalities, we have proven

$$E_X \left[\log \frac{h_\theta(x^*, z)}{\sum_{x \in X} h_\theta(x, z)} \right] = -\ell(\theta) \leq I(x, z) - \log N \implies I(x, z) \geq -\ell(\theta) + \log N,$$

as we wanted. \square

Part III

REPRESENTATION LEARNING

7 | CONTEXT

Machine learning is the field of computer science that studies algorithms that improve automatically through experience from examples. These algorithms allow computers to discover how to perform tasks without being explicitly programmed to do them. For the computers to learn, it is mandatory that a finite set of data (or dataset) \mathcal{D} is available.

Whenever a computer is provided with data, the data can be *labeled* or *unlabeled*. Labeled data is the one that, each point $x_i \in \mathcal{D}$ is related to a tag $y_i \in Y$, where Y is a set of classes. Unlabeled data is the kind of data that does not have a label or class associated to it, so it is just $x_i \in \mathbb{R}^d$.

Depending on how the data (or *signal*) is given to the computer, the machine learning approaches can be divided into three broad categories:

1. *Supervised learning*. In this category the goal is to use the labeled data in order to find a function that maps the dataset to the set of classes. That is a function $g : \mathcal{D} \rightarrow Y$. An example of supervised learning is image classification: giving a label to an image.
2. *Unsupervised learning*. In this case, the data is unlabeled, so the approach is completely different. Usually, the goal here is to discover hidden patterns in data or to learn features from it. An example of this kind of learning is K means, which consists in clustering the data in k groups. It is also known as *self-supervised* learning.
3. *Reinforcement learning*. This is the area concerned with how intelligent agents take decisions in a specific environment in order to obtain the best reward in their objective. It focuses on finding a balance between exploration of uncharted territory and exploitation of the current knowledge.

In this work, we will focus on unsupervised learning. Particularly, in representation learning.

There are many different tasks that can be performed with the data, such as linear regression, logistic regression, or classification. In any of these tasks, computers might need to do intermediate steps before giving a label to the input example. Sometimes, they must create a *representation* that contains the data's key qualities. Here is where *representation learning* is born.

Having a representation \tilde{x} of a datapoint $x \in \mathbb{R}^d$, if a machine learning model tries to make a posterior task, such as classification, the input x must be transformed to a, usually lower dimensional, vector r in order to perform the final label.

Features are parts or patterns of an datapoint $x \in \mathcal{D}$ that help to identify it. In fact, it is desirable that this attribute is shared by all independent units

that represent the same object. For instance, if we consider an image of any square, we should be able to identify 4 corners and 4 edges. These could be features of a square. When we mention feature detection, we are addressing the methods for detecting these features of a datapoint.

Representation learning is a set of techniques that allow a system to discover the representations needed for feature detection or classification. In contrast to manual feature engineering (which involves manually exploring the data and finding relationships in it), feature learning allows a machine to learn the features and to use them to perform a task.

Feature learning can be supervised or unsupervised.

- In supervised feature learning, representations are learned using labeled data. Examples of this kind of feature learning are supervised Neural Networks and Multilayer Perceptron.
- In unsupervised learning, the features are learned using unlabeled data. There are many examples of this, such as independent component analysis (ICP) and autoencoders.

In this work, we will be working with unsupervised feature learning, so we will have unlabeled data that we would like to find a representation for.

The performance of machine learning methods is heavily dependent on the choice of data features (Bengio *et al.*, 2014). This is why most of the current effort in machine learning focuses on designing preprocessing and data transformation that lead to good quality representations. A representation will be of good quality when its features produce good results when we evaluate the *accuracy* of our model.

The main goal in representation learning is to obtain features of the data that are generally good for either of the supervised tasks. These tasks are usually called *downstream tasks*. That is, we would like to obtain a representation that is either good for image classification (giving an image a label of what we can see in it) or image captioning (producing a text that describes the image).

Data's features that are invariant through time are very useful for machine learning models. In (Wiskott & Sejnowski, 2002), *slow features* are presented. Slow features are defined as features of a signal (which can be the input of a model) that vary slowly during time. That means, if \mathbf{X} is a *time series*¹, we will try to find any number of features in \mathbf{X} that vary the most slowly. These kind of features are the most interesting ones when creating representations, since they give an abstract view of the original data.

Example 6. In computer vision, the value of the pixels in an image can vary fast. For instance, if we have a zebra on a video and the zebra is moving from

¹ A time series is an ordered sequence of values of a random variable at, usually, equally spaced time intervals.

one side of the image to the other, due to the black stripes of this animal, the pixels will fast change from black to white and vice versa, so value of pixels is probably not a good feature to choose as an slow feature. However, there will always be a zebra on the image, so the feature that indicates that there is a zebra on the image will stay positive throughout all the video, so we can say that this is a slow feature.

We will be studying different models that try to learn representations from raw data without labels, as we have mentioned. We usually need a function that measures what is the penalty that the model gets for a choice of a parameter. This is called a *loss function*, that we will want to optimize.

For instance, in a regression problem, a good example of loss function is *mean squared error*, which is expressed as follows:

$$\text{MSE} = \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{n}.$$

In a classification problem, each datapoint x_i has a correct classification y_i . In this case, the score of the correct category y_i should be greater than the sum of the scores of all incorrect categories y_j with $j \neq i$, so we could use a function like *support vector machine (SVM) loss*:

$$\text{SVMLoss} = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

In the context of machine learning, a *model* is the result of running a machine learning algorithm in data. This model will represent what the computer has learned from the data using this algorithm. As an easy example, when we use a linear regression algorithm, we obtain a model that is a vector of coefficients with specific values.

In the following chapters, we will explain different kinds of neural networks and which kind of loss functions they use.

8

GENERATIVE MODELS

The vast majority of the problems in ML are usually of a discriminative nature, which is almost a synonym of supervised learning. However, there also exist problems that involve learning how to generate new examples of the data. More formally:

- Definition 8.1.**
1. *Discriminative models* estimate $p(y|x)$, the probability of a label y given an observation x .
 2. *Generative models* estimate $p(x)$, the probability of observing the datapoint x . If the dataset is labeled, a generative model can also estimate the distribution $p(x|y)$.

From now on, let \mathcal{D} be any kind of observed data. This will always be a finite subset of samples taken from a probability distribution p_{data} . There are models that, given \mathcal{D} , try to approximate the probability distribution that lies underneath it. These are called *generative models* (G.M.).

Generative models can give parametric and non parametric approximations to the distribution p_{data} . In our case, we will focus on parametric approximations where the model searches for the parameters that minimize a chosen metric (which can be a distance or other kind of metric such as K-L divergence) between the model distribution and the data distribution.

We can express our problem more formally as follows. Let θ be a generative model within a model family \mathcal{M} . The goal of generative models is to optimize:

$$\min_{\theta \in \mathcal{M}} d(p_{\text{data}}, p_{\theta}),$$

where d stands for the distance between the distributions. We can use, for instance, K-L divergence.

Generative models have many useful applications. We can however remark the tasks that we would like our generative model to be able to do. Those are:

- Estimate the density function: given a datapoint, $x \in D$, estimate the probability of that point $p_{\theta}(x)$.
- Generate new samples from the model distribution $x \sim p_{\theta}(x)$.
- Learn useful features of the datapoints.

If we have a look again at the example of the zebras, if we make our generative model learn about images of zebras, we will expect our $p_{\theta}(x)$ to be high for zebra's images. We will also expect the model to generate new images of this animal and to learn different features of the animal, such as their big size in comparison with cats.

8.1 AUTOREGRESSIVE MODELS

In time-series theory, autoregressive models use observations from previous time steps to predict values at the current time. Fixing an order of the variables x_1, \dots, x_n , the distribution for the i -th random variable depends on all the preceding values in the particular chosen order. We will make use of the name of these models to define the machine learning approach.

A very first definition of *autoregressive models (AR)* would be the following one: *autoregressive models are feed-forward models that predict future values using past values.* Let us go deeper into this concept and explain how it behaves.

Again, let \mathcal{D} be a set of n -dimensional datapoints x . We can assume that $x \in \{0, 1\}^n$ for simplicity, without losing generality. If we choose any $x \in \mathcal{D}$, using the chain rule of probability, we obtain

$$p(x) = \prod_{i=1}^n p(x_i | x_1, \dots, x_{i-1}) = \prod_{i=1}^n p(x_i | \mathbf{x}_{<i}),$$

where $\mathbf{x}_{<i} \in \mathbb{R}^{i-1}$ is a vector whose components are the previous x_j for $j = 1, \dots, i-1$, that is: $\mathbf{x}_{<i} = [x_1, \dots, x_{i-1}]$.

It is known that given a set of discrete and mutually dependent random variables, they can be displayed in a table of conditional probabilities. If K_i is the number of states that each random variable can take then $\prod K_i$ is the number of cells that the table will have. If we represent $p(x_i | \mathbf{x}_{<i})$ for every i in tabular form, we can represent any possible distribution over n random variables.

This, however, will cause an exponential growth on the complexity of the representation, due to the need of specifying 2^{n-1} possibilities for each case. In terms of neural networks, since each column must sum 1 because we are working with probabilities, we have $2^{n-1} - 1$ parameters for this conditional, and the tabular representation becomes impractical for our network to learn when n increases.

In autoregressive generative models, the conditionals are specified as we have mentioned before: parameterized functions with a fixed numbers of parameters. More precisely, we assume the conditional distributions to be Bernoulli random variables and learn a function f_i that maps these random variables to the mean of the distribution. Mathematically, we have to find

$$p_{\theta_i}(x_i | \mathbf{x}_{<i}) = \text{Bern}(f_i(x_1, \dots, x_{i-1})),$$

where θ_i is the set of parameters that specify the mean function $f_i : \{0, 1\}^{i-1} \rightarrow [0, 1]$.

Then, the number of parameters is reduced to $\sum_{i=1}^n |\theta_i|$ so we can not represent all possible distributions as we could when using the tabular form of the conditional probabilities. We are now setting the limit of its expressiveness because we are setting the conditional distributions $p_{\theta_i}(x_i | \mathbf{x}_{<i})$ to be *Bernoulli* random variables with the mean specified by a restricted class of parametrized functions.

Let us see a very simple case first in order to understand it better. Let σ be a *sigmoid*¹ non linear function and $\theta_i = \{\alpha_0^{(i)}, \alpha_1^{(i)}, \dots, \alpha_{i-1}^{(i)}\}$ the parameters of the mean function. Then, we can define our function f_i as :

$$f_i(x_1, \dots, x_{i-1}) = \sigma(\alpha_0^{(i)} + \alpha_1^{(i)}x_1 + \dots + \alpha_{i-1}^{(i)}x_{i-1}).$$

In this case, the number of parameters would be $\sum_{i=1}^n i = \frac{n(n+1)}{2}$, so using Big O notation, we would be in the case of $O(n^2)$. We will state now a more general and useful case, giving a more interesting parametrization for the mean function: *multi layer perceptrons*²(MLP).

For this example we will consider the most simple MLP: the one with one hidden layer. Let $h_i = \sigma(\mathbf{A}_i \mathbf{x}_{<i} + c_i)$ be the hidden layer activation function. Remember that $h_i \in \mathbb{R}^d$. Let $\theta_i = \{\mathbf{A}_i \in \mathbb{R}^{d \times (i-1)}, c_i \in \mathbb{R}^d, \alpha^{(i)} \in \mathbb{R}^d, b_i \in \mathbb{R}\}$ the set of parameters for the mean function f_i , that we define as:

$$f_i(\mathbf{x}_{<i}) = \sigma(\alpha^{(i)} h_i + b_i).$$

In this case, the number of parameters will be $O(n^2 d)$.

This is the simplest example. Currently, there are alternative parametrization models, such as the *Neural Autoregressive Density Estimator* (Larochelle & Murray, n.d.), that provide a more statistically and computationally efficient solution. In fact, the number of parameters is reduced from $O(n^2 d)$ to $O(nd)$. Also, *RNADE* (Uria et al., 2014) extends NADE to learn generative models over real-valued data, generalizing the case that we have just exposed. However, these models are out of the scope of this project so no further explanation will be given.

¹ A sigmoid function is a bounded, differentiable, real function which derivative is non-negative at each point and it has exactly one inflection point.

² Multi layer perceptrons are feed-forward neural networks with at least 3 layers: input, hidden and output layers; each one using an activation function.

We are now ready to connect the concepts of mutual information and generative models that we have presented. In unsupervised learning, it is a common strategy to predict future information and to try to find out if our predictions are correct. In *natural language processing*, for instance, representations are learned using neighbouring words (Mikolov *et al.*, 2013). In the field of computer vision, some studies have been able to predict color from grey-scale (Doersch *et al.*, 2016).

When we work with high-dimensional data, it is not useful to make use of an unimodal loss function to evaluate our model. If we did it like this, we would be assuming that there is only one peak in the distribution function and that it is actually similar to a Gaussian. This is not always true, so we can not assume it for our models. Generative models can be used for this purpose: they will model the relationships in the data x . However, they ignore the context c in which the data x is involved. As an easy example of this, an image contains thousands of bits of information, while the label that classifies the image contains much less information, say, 10 bits for 1024 categories. Because of this, modeling $P(x|c)$ might not be the best way to proceed if we want to obtain the real distribution that generates our data.

During the last few years, the representation learning problem has been approached using different machine learning frameworks. The most competitive ones have been self-supervised contrastive representation learning Oord *et al.* (2019); Tian *et al.* (2020); Hjelm *et al.* (2019); Gutmann & Hyvarinen (n.d.); Chen *et al.* (2020b); He *et al.* (2020a) using *contrastive losses*, and they have empirically outperformed other approaches.

In contrastive learning, different “views” of the same input are created. These are also called *positive samples*. Then, they are compared with *negative samples*, which are views created from an input that does not share information with the input of the positive sample. A very interesting idea would be to try and maximize the mutual information between positive samples and push apart the views taken from negative samples.

There are many ways of creating samples, both positive and negative, of an input. For instance:

- Randomly cropping different parts of an image. These would be examples of positive examples.
- Rotating or flipping images or crops of them would also be examples of positive samples.
- Taking different time-steps of a video would create positive samples.
- Selecting different parts of the same text would also be a positive example.

- Negative samples are created by applying one of the previous techniques to images that have nothing in common with the positive input.

In fact, if v_1, v_2 are two views of an input, we can think of the positive pairs as points coming from a joint distribution over the views $P(v_1, v_2)$, and negative samples coming from the product of the marginals $P(v_1)P(v_2)$, (Tian *et al.*, 2020).

It is important to find a way to determine how much shared information between the views is needed, in order to make the representations obtained good enough for any downstream task. Here is where the *InfoMin principle* is born.

Definition 9.1 (The InfoMin principle). A good set of views are those that share the minimal information necessary to perform well at the downstream task.

Our goal here will be to seek for a way of extracting shared information between the context c and the data x . Here is where we link the mutual information with representation learning. Remember that the mutual information of two random variables, say x and c in this case, is:

$$I(x, c) = \sum_{x, c} P(x, c) \log \frac{P(x|c)}{P(x)} \quad (12)$$

Maximizing the MI between x and c , we extract the latent variables that the inputs have in common.

9.1 CONTRASTIVE PREDICTIVE CODING

We can apply these concepts in a concrete framework, presented firstly in Oord *et al.* (2019). Let us see what information is used and how it is treated in order to train a model that tries to obtain useful representations for downstream tasks.

In this section, if x is an input signal for our network, x_t will be the value of the input at instant t . We will also make reference to x_{t+k} , meaning that x_{t+k} is k steps ahead of time to x_t .

Firstly, an *encoder* is used. An encoder is a model that, given an input x , provides a feature map or vector that holds the information that the input x had.

So, we will use an encoder g_{enc} that transforms the input sequence of observations x_t to a sequence of latent representations

$$z_t = g_{enc}(x_t).$$

After we have obtained z_t , we use it as input of an autoregressive model, explained before, to produce a context latent representation:

$$c_t = g_{ar}(z_{\leq t}).$$

In this case, c_t will summarize the information of z_i for $i \leq t$. Following the argument that we gave before, predicting the future x_{t+k} using only a

generative model (say $p_k(x_{t+k}|c)$) might not be correct, since we would be ignoring the context.

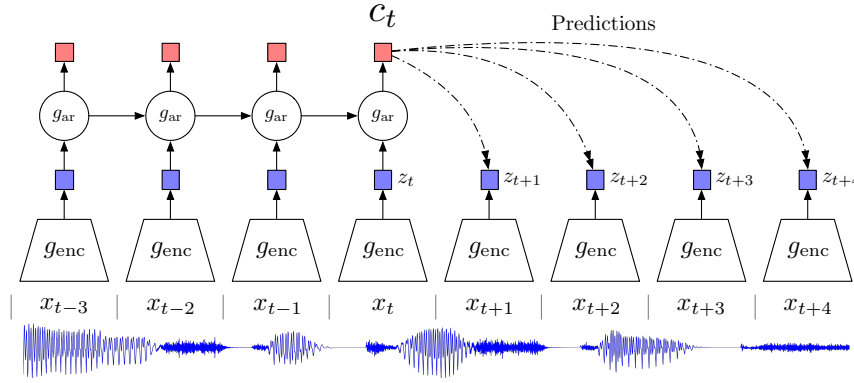


Figure 5: Image from (Oord *et al.*, 2019). Overview of Contrastive Predictive Coding framework using audio signal as input.

Let us see how we train the encoder g_{enc} and the autoregressive model g_{ar} .

In Chapter 4, we gave the notions of the general idea of Noise Contrastive Estimation. Now, we can apply those ideas to a particular case, in which one of the subsets, say X only has one element, and the other one has $N - 1$ elements. We combine both sets in X for the following argument.

Let $X = \{x_1, \dots, x_N\}$ be a set of N random samples. X will contain a positive sample taken from the distribution $P(x_{t+k}|c_t)$ and $N - 1$ negative samples from the distribution proposed $P(x_{t+k})$. With this set, we would like to optimize the following loss function, which is an alternative expression of (6):

Definition 9.1.1. The loss \mathcal{L}_N defined as

$$\mathcal{L}_N = -E_X \left[\log \frac{f_k(x_{t+k}, c_t)}{\sum_{x_j \in X} f_k(x_j, c_k)} \right], \quad (13)$$

is known as the *InfoNCE* (Information Noise Contrastive Estimation) loss.

This loss also gives the bound in Proposition ?? the name of *InfoNCE Bound*.

It is clear that this loss is based in Noise Contrastive Estimation (Gutmann & Hyvarinen, n.d.).

Let us have a look at the *categorical cross-entropy* loss function:

$$\mathcal{L}(y, s) = - \sum_i^C y_i \log(s_i)$$

where C is the number of possible classes in a classification problem, y_i are the groundtruth of each class and s_i is the score of each class.

As we remarked before, we can say that \mathcal{L}_N is no more than the categorical cross-entropy of classifying the positive sample of X correctly, with the argument of the logarithm being the prediction of the model. If we note with

$[d = i]$ as an indicator of the sample x_i being the positive sample in X , the optimal probability for this loss can be written as $P(d = i|X, c_t)$.

Now, the probability that x_i was drawn from the conditional distribution $P(x_{t+k}|c_t)$ that has the context in account, rather than the proposal distribution $P(x_{t+k})$ that does not have c_t in account, leads us to the following expression:

$$P(d = i|X, c_t) = \frac{\frac{P(x_i|c_t)}{P(x_i)}}{\sum_{j=1}^N \frac{P(x_j|c_t)}{P(x_j)}}.$$

This is the optimal case for (13).

In fact, if we denote $f(x_{t+k}, c_t)$ as the density ratio that preserves the mutual information between x_{t+k} and c_t in the mutual information definition (12), if x_{t+k} is k steps ahead on time respect to x_t , then

$$f_k(x_{t+k}, c_t) \propto \frac{P(x_{t+k}|c_t)}{P(x_{t+k})}, \quad (14)$$

where \propto means that the member on the left is proportional to the member on the right. We can see that the optimal value $f_k(x_{t+k}, c_t)$ does not depend on $N - 1$, the number of negative samples in X . Using this density ratio, we are relieved from modeling the high dimensional distribution x_{t_k} . In [Oord et al. \(2019\)](#), for instance, the following log-bilinear model expression is used:

$$f_k(x_{t+k}, c_t) = \exp(z_{t+k}^T W_k c_t).$$

In the proposed model, we can either use the representation given by the encoder (z_t) or the representation given by the autoregressive model (c_t) for downstream tasks. Clearly, the representation that aggregates information from past inputs will be more useful if more information about the context is needed. Furthermore, any type of models for the encoder and the autoregressive models can be used in this kind of framework.

It is clear how, using (13) we are using exactly the same function that we were using in Equation 10 in Chapter 6. This way, if we maximize this loss, we are also maximizing the mutual information between x_{t+k} and the context c_t .

9.2 GOOD VIEWS FOR CONTRASTIVE LEARNING

We have presented a framework in which a set $X = \{x_1, \dots, x_n\}$ contains a sample from the distribution $P(x_{t+k}, c_t)$ and the rest are samples from the distribution $P(x_{t+k})$. These samples are different views of the data.

The choice of the views affects the results in the downstream tasks ([Tian et al., 2020](#)). The views will affect on the training, hence, it will affect to the representations that are obtained. We would like to have some guarantees that the views that we are choosing provide us with good examples for our training. Let us formalize this idea.

Given two random variables v_1, v_2 , our goal was to learn a function to discriminate the samples from the joint distribution and the product of the marginal distributions, resulting on a mutual information estimator between v_1 and v_2 . In practice, v_1 and v_2 are two views of the same input x , using one of the methods that we mentioned in the introduction. We would like to have that, if y is a downstream task, the mutual information between both the inputs and the downstream task, is the same as the mutual information between the input x and the downstream task, i.e.:

$$I(v_1, y) = I(v_2, y) = I(x, y).$$

Also, we would like to remove the information that is not relevant for our downstream task. This is done by obtaining the pair of views (v_1^*, v_2^*) such that the mutual information between them is the minimum of the mutual information between all the possible views (v_i, v_j) . Formally, that is obtaining

$$(v_1^*, v_2^*) = \min_{v_1, v_2} I(v_1, v_2),$$

These two ideas form the *InfoMin Principle* that we mentioned before in Definition 9.1.

Usually, the views are encoded using an encoder f , not having it to be the same for both views. We can say that $z_i = f_i(v_i)$ for $i \in \{1, 2\}$. If an encoder is sufficient, then it has to maintain the mutual information between the random variables after one of them has been encoded. More formally,

Definition 9.2.1. We say an encoder f_i of a view v_i , with $i \in \{1, 2\}$ is *sufficient* in the contrastive learning framework if, and only if it maintains the mutual information between the pairs (v_i, v_j) and $(f_i(v_i), v_j)$ with $j \in \{1, 2\}$. That is

$$I(v_i, v_j) = I(f_i(v_i), v_j).$$

This usually means that no information was lost in the process of encoding. We want to extract only the most essential information and do not learn the “extra” information between the views.

Definition 9.2.2. We say that a sufficient encoder f_i of v_i is *minimal* if, and only if, the mutual information between $(f_i(v_i), v_j)$ is lesser than the mutual information between $(f(v_i), v_j)$ for any other sufficient encoder f . That is:

$$I(f_i(v_i), v_j) \leq I(f(v_i), v_j) \quad \text{for all sufficient } f.$$

With these notions already presented, we would like to define what representations are good for a downstream task. We get to the following definition (Tian et al., 2020):

Definition 9.2.3. For a task \mathcal{T} , whose goal is to predict a label y from the input data x , the optimal representation z^* encoded from x , that is $z^* = f(x)$ for some encoder f , is the minimal sufficient statistic with respect to y .

This means that if we use z^* to make a prediction using a machine learning model, we are using the same information that we would be using the whole

input x . What is more, since we are following the InfoMin principles and we are dismissing all the non relevant information, z^* provides with the smallest complexity.

Proposition 8. *Let f_1, f_2 be minimal sufficient encoders, and \mathcal{T} be a downstream task with label y . Then, the optimal views (v_1^*, v_2^*) from the data x are the ones that have minimal mutual information*

$$(v_1^*, v_2^*) = \arg \min_{v_1, v_2} I(v_1, v_2),$$

subject to $I(v_1, y) = I(v_2, y) = I(x, y)$.

Given the optimal views (v_1^, v_2^*) , the representation z_1^* learned by contrastive learning is optimal for \mathcal{T} .*

The proof of this proposition is out of the scope of this work so no further information will be provided. The last statement of the Proposition 8 is a consequence of the minimality and sufficiency of f_1 and f_2 .

This proposition carries the most important mathematical conclusion from this section. It will serve to prove sufficient conditions for views to be effective for contrastive learning, where effectiveness is measured as effectiveness in downstream tasks.

10 | TRIPLET LOSSES

Contrastive learning exploits the idea of comparing the input x with other different examples, either from the same class or from another class. It aims to produce closer representations for the examples of the same class and distant representations for elements of other classes.

We have been using the loss in Equation (13), which we built using noise contrastive estimation. However, this is not the only way of approaching this problem, and other types of losses have been used for similar purposes, forgetting the part of mutual information maximization and replacing it with *geometrical distance* optimization.

In this section, we will present *Triplet losses*, other kind of loss functions that also compare different views of the same input x .

10.1 FROM DEEP METRIC LEARNING TO TRIPLET LOSSES AND ITS GENERALIZATION

Distance metric learning also aims to learn an embedding representation of an input data x that preserves the distance between similar data points close and also makes the distance between different datapoints far on the embedding space (Sohn, 2016).

Let us set the notation that we will use first. We will consider sets of triplets (x, x^+, x^-) where:

- The element x is an anchor point,
- The element x^+ is a positive instance,
- The element x^- is a negative instance.

Example 7. Let us present a very simple example. If our input image is a cat, that would be the anchor x . Clearly, a positive instance would be an image of another cat or even the same cat seen from another perspective. A negative instance would be a photo of any other animal, in this case we use a dog.

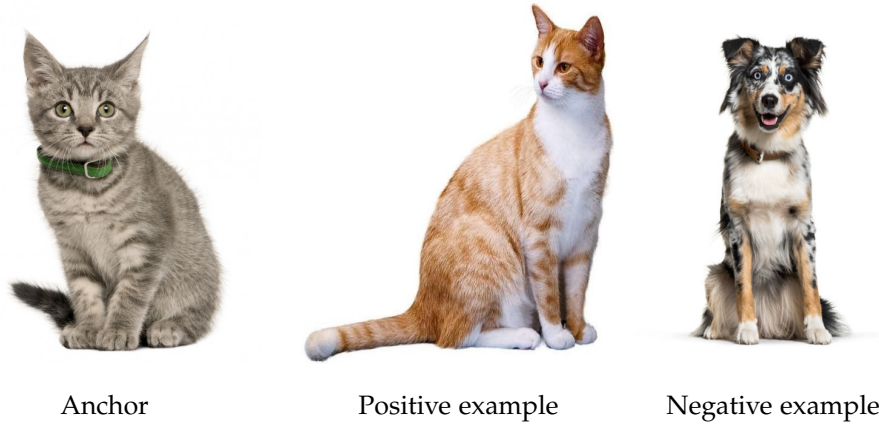


Figure 6: Example of an anchor x , a positive instance x^+ and a negative instance x^- . Images obtained from Google.

The main idea is to learn a representation of x , say $g(x)$, such that the distance of the representation of the input is closer in distance to the representation of the positive sample x^+ than the representation of the negative sample x^- . Using the norm¹, we can formally express that as follows:

$$\|g(x) - g(x^+)\|_2 \leq \|g(x) - g(x^-)\|_2,$$

for each triplet in the set.

Support-vector machines (SVMs) are supervised learning models used for classification or regression problems. They are one of the most robust prediction methods. They search for a hyperplane h in high or infinite dimensional space that separates the data as much as possible, making use of *support vectors*, the datapoints that are closest to the hyperplane. If the data is linearly separable, we can select two hyperplanes h_1, h_2 that are parallel to h and making the distance from them to h as large as possible. That region is called the *margin*.

Coming back to our triplets problem, we also want to introduce a margin between the distances of the elements of the triplets, in order to separate positive examples from negative examples as much as possible. This way, we introduce a *margin* term α , rewriting our last equation as follows:

$$\|g(x) - g(x^+)\|_2 + \alpha < \|g(x) - g(x^-)\|_2.$$

Using this inequality, we can define a hinge loss function for each triplet in the set:

$$\ell^\alpha(x, x^+, x^-) = \max\left(0, \|g(x) - g(x^+)\|_2^2 - \|g(x) - g(x^-)\|_2^2 + \alpha\right). \quad (15)$$

This loss has been defined for a single triplet. Now, we can define a global loss that accumulates the loss in Equation (15) using all the triplets in set.

¹ A definition of the norm can be found on Appendix A, Definition A.1.

Definition 10.1.1. Given a set of triplets, each containing an anchor, a positive example and a negative example, $\mathcal{T} = \{(x_i, x_i^+, x_i^-)\}_{i \in \Lambda}$, we define a triplet loss as follows:

$$\mathcal{L}(x_i, x_i^+, x_i^-) = \sum_{i \in \Lambda} \ell^\alpha(x_i, x_i^+, x_i^-). \quad (16)$$

We use this loss to train models in order to improve the representations obtained. It would be interesting to present the model non-trivial metric to the learning algorithm. When the representation g improves, this is harder to do, and this results in slow convergence and expensive data sampling methods.

10.2 GENERALIZATION OF TRIPLET LOSSES AND N —PAIRS LOSS

In a single evaluation of the loss function over a triplet during the learning process, we are comparing one positive sample to one negative sample. In practice, after looping over sufficiently many triplets, we expect the distance between positive examples and negative examples to be maximized. However, this will surely be a slow process if our dataset has many examples and also, in each step we will be separating the positive element from the specific negative element to which we are comparing it in that evaluation. Thus, the technique might be unstable (Sohn, 2016).

In order to fix this, a good idea would be to compare in each evaluation a positive sample with multiple negative samples, generalizing the case exposed before. This way, we would like the positive sample to increase its distance to *all* of the negative samples at the same time. Let us present a loss that generalizes the loss in Equation (16).

Definition 10.2.1. Let x^+ be a positive example of the anchor x , and consider the set $X^- = \{x_1^-, \dots, x_{N-1}^-\}$ of $(N-1)$ negative samples. Given an encoder g , the $(N+1)$ —tuple loss is defined as follows:

$$\mathcal{L}_{(N+1)\text{-tuple}}(x, x^+, X^-) = \log \left(1 + \sum_{i=1}^{N-1} \exp \left(g(x)^T g(x_i^-) - g(x)^T g(x^+) \right) \right) \quad (17)$$

Remark 6. If we consider the case $N = 2$, we have

$$\mathcal{L}_{(2+1)\text{-tuple}}(x, x^+, x^-) = \log \left(1 + \exp \left(g(x)^T g(x^-) - g(x)^T g(x^+) \right) \right).$$

This expression is very similar to the one in Equation (15). In fact, if the norm in Equation (15) is unit and g minimizes $\mathcal{L}_{(2+1)\text{-tuple}}$, then it minimizes ℓ^α , and hence both losses are equivalent.

Applying the $(N+1)$ —tuple loss in deep metric learning is computationally expensive. Indeed, if we apply Stochastic Gradient Descent (SGD) with batch size M , then we have to evaluate $M \times (N+1)$ times our function ℓ^α in each update. Because of this, if we increase M and N , the number of evaluations grows quadratically. We would like to avoid this.

Consider the set of N pairs of examples, with the constraint of each pair belonging to a different class, i.e. $X = \{(x_1, x_1^+), \dots, (x_N, x_N^+)\}$ with $y_i \neq y_j$ for all $i \neq j$. We now build N tuples where each tuple has all the positive samples and the i -th anchor, that is:

$$\{S_i\}_{i=1}^N, \quad \text{where } S_i = \{x_i, x_1^+, \dots, x_N^+\}.$$

We can consider that each tuple has x_i as anchor, x_i^+ as positive example and x_j^+ for $j \neq i$ as negative samples, since they were all from different classes.

Definition 10.2.2. In the last conditions, we can define the *multi class N -pair loss* as follows:

$$\mathcal{L}_{N\text{-pair-mc}} \left(\{(x_i, x_i^+)\}_{i=1}^N \right) = \frac{1}{N} \sum_{i=1}^N \log \left(1 + \sum_{j \neq i} \exp \left(g(x_i)^T g(x_j^+) - g(x_i)^T g(x_i^+) \right) \right) \quad (18)$$

This way, we are combining the $(N+1)$ -tuple loss and the N -pair construction that we presented, enabling highly scalable training. This loss has empirically proved in to have advantages if we compare it to other variations of mini-batch methods (Sohn, 2016).

10.2.1 InfoNCE Bound as a triplet loss

The InfoNCE loss on Equation (13) has proved to be useful in representation learning. Let us consider a reformulation on it. Firstly, since f_k was an exponential, we can also consider e^f and remove the exponential from f , this is just notation. Now, in Poole et al. (2019) the InfoNCE bound in (11) is rewritten as follows:

$$I(X, Y) \geq E \left[\frac{1}{N} \sum_{i=1}^N \log \frac{e^{f(x_i, y_i)}}{\frac{1}{N} \sum_{j=1}^N e^{f(x_i, y_j)}} \right] \triangleq I_{NCE}(X, Y)$$

where we have just named the right hand side of the inequality as $I_{NCE}(X, Y)$. Now, we can transform it in the following way:

$$\begin{aligned} I_{NCE} &= E \left[\frac{1}{N} \sum_{i=1}^N \log \frac{e^{f(x_i, y_i)}}{\frac{1}{N} \sum_{j=1}^N e^{f(x_i, y_j)}} \right] \\ &= E \left[\frac{1}{N} \sum_{i=1}^N \log \frac{1}{\frac{1}{N} \sum_{j=1}^N e^{f(x_i, y_j) - f(x_i, y_i)}} \right] \\ &= E \left[-\frac{1}{N} \sum_{i=1}^N \log \frac{1}{N} \sum_{j=1}^N e^{f(x_i, y_j) - f(x_i, y_i)} \right] \end{aligned} \quad (19)$$

And now, we only have to use the basic properties of the logarithm to see that

$$\begin{aligned}
 (19) &= E \left[-\frac{1}{N} \left(\sum_{i=1}^N \log \frac{1}{N} + \sum_{i=1}^N \log \sum_{j=1}^N e^{f(x_i, y_j) - f(x_i, y_i)} \right) \right] \\
 &= E \left[-\frac{1}{N} \left(N(\log 1 - \log N) + \sum_{i=1}^N \log \left(1 + \sum_{j \neq i} e^{f(x_i, y_j) - f(x_i, y_i)} \right) \right) \right] \\
 &= E \left[-\frac{1}{N} (-N \log N) \right] + E \left[-\frac{1}{N} \sum_{i=1}^N \log \left(1 + \sum_{j \neq i} e^{f(x_i, y_j) - f(x_i, y_i)} \right) \right] \\
 &= \log N - E \left[\frac{1}{N} \sum_{i=1}^N \log \left(1 + \sum_{j \neq i} e^{f(x_i, y_j) - f(x_i, y_i)} \right) \right] \tag{20}
 \end{aligned}$$

In the particular case where X, Y take values in the same space and f has the particular form

$$f(x, y) = \phi(x)^T \phi(y),$$

for some ϕ , Equation (20) is the same (up to constants and change of sign) as the expectation of the multi-class N -pair loss in Equation (18).

Hence, we have found an equivalence between representation learning by maximizing I_{NCE} using a symmetric separable critic $f(x, y)$ and an encoder g shared across views and using a multi class N -pair loss, which is a generalization of a triplet loss.

Part IV

NEW FRAMEWORKS FOR REPRESENTATION LEARNING

In order to be able to understand the chapters that will come later in this document, it is important to make a brief introduction of what *Deep Learning* (DL) refers to. Deep learning is included in the field of Machine Learning, which is also included in the field of general Artificial Intelligence.

In Chapter 7, an intuitive definition of what Machine Learning is was given. We said that ML studies the algorithms that improve from experience. Tom M. Mitchell (Mitchell, 1997) provided a more formal definition of what *learning from experience* means:

Definition 11.1. A computer program is said to *learn* from experience E with respect some class of tasks T and performance measure P , if its performance at tasks in T , as measured by P , improves with experience E .

We would also like to have a DL definition. In Deng (2014), multiple similar definitions are given. We present here the simplest of them:

Definition 11.2. *Deep Learning* is a class of ML learning techniques that exploit many layers of non-linear information processing for supervised or unsupervised feature extraction and transformation, and for pattern analysis and classification.

Usually, these techniques are based on the biologically inspired *neural networks* (NNs), which consists in several connected units: the *neurons*. Each neuron is basically a Perceptron, which is a weighted sum followed by a non-linear function, called an *activator* in the ML context. Formally, the output of each neuron is

$$y = \phi \left(w_0 + \sum_{i=1}^N w_i x_i \right).$$

There are many activation functions, but the following examples must be remarked:

- Sigmoid. The sigmoid function is defined as follows:

$$\phi(x) = \frac{1}{1 + e^{-x}}.$$

This is one of the most common used activation functions. It is differentiable, monotonic and smooth. One of its main disadvantages is that at the right part of the function, the change in the values that the function takes converges to zero, so we get to the *vanishing gradient* problem and the learning is minimal.

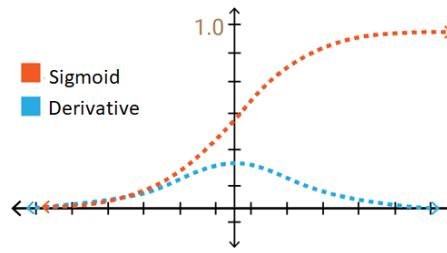


Figure 7: Sigmoid. Image from [this Medium article](#).

- Hyperbolic Tangent. This function is defined as follows:

$$\phi(x) = \tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}.$$

This activation function has a small advantage over the sigmoid: its derivative is more steep, which means it can get more value and the learning can be more efficient.

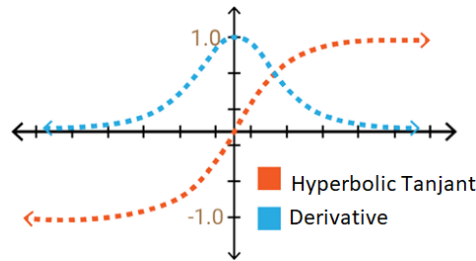


Figure 8: Hyperbolic Tangent. Image from [this Medium article](#).

- Rectified Linear Unit (ReLU). This function takes the following form:

$$\phi(x) = \max(0, x).$$

ReLU is highly computationally efficient and non-linear. Its main problem is that when the inputs approach zero or are negative, the network can not perform back propagation and can not learn.

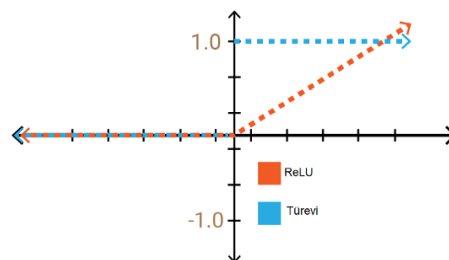


Figure 9: ReLU. Image from [this Medium article](#).

11.1 NEURAL NETWORKS

Using neurons and activation functions, we can formally define NNs. A NN with L hidden layers is a deterministic non-linear function f , parametrized by a set of matrices $W = \{W_0, \dots, W_L\}$ and non-linear activation functions

$\{\phi_0, \dots, \phi_L\}$. Given an input x , the output y of the network is calculated as follows:

$$h_0 = \phi_0(W_0^T x), \dots, h_l = \phi_l(W_l^T h_{l-1}), \dots, y = \phi_L(W_L^T h_{L-1}).$$

Having a NN, we consider it *deep* when the number of hidden layers (and, consequently, the number of matrices) is considered high.

Neural networks use loss functions, which define how well the output returned by the network matches the real output, reducing the learning problem to an optimization problem. The problem is finding W^{opt} , such that

$$W^{\text{opt}} = \arg \min_w \sum_{n=1}^N l(y_n, f_w(x_n)),$$

where $\mathcal{D} = \{(x_n, y_n)\}$ is a dataset.

This problem is solved using a variant of *stochastic gradient descent* (SGD). This algorithm involves the computation of the loss function derivatives respect to the network parameters, and updates the parameters using this derivatives. Specifically, the parameters are updated as follows:

$$W_{t+1} = W_t + \eta \nabla f(W_t),$$

where $\eta \in \mathbb{R}^+$ is a small constant called the *learning rate*. This algorithm guarantees convergence to local minimums of f and, if f is convex, the algorithm converges to a global minimum.

The last comment about neural networks is that, since the weights $W = \{W_0, \dots, W_L\}$ are constantly updated, the derivatives have to be computed repeatedly. The computational cost of this is quite high. *Backpropagation* was born to calculate the derivatives of the weights much faster. The intuitive idea is that the gradient of the layer l is computed using the gradient of the layer $l + 1$ using the chain rule.

Understanding both SGD and Backpropagation is crucial for understanding how NNs work. However, in the experimentation part of this work we will focus on researching how a few hyperparameters affect the results of the proposed frameworks, so no further explanation on these important concepts will be provided.

11.2 CONVOLUTIONAL NEURAL NETWORKS

Convolutional Neural Networks (CNNs) are a specific type of Neural Networks. The difference that we have between CNNs and NNs is that CNNs assume that the inputs have local dependencies. For instance, using an image, a CNN assumes (most of the time correctly) that given a pixel x_{ij} of an input x , the neighbours of this pixels will have similar values or intensities. This allows us to encode certain properties into the architecture of the network.

If we use colored images, the input of the CNNs are 3-dimensional volumes, which will be transformed in some layers to other 3-dimensional volumes.

In order to do this we have different types of layers, remarking *convolutional* layers, *pooling* layers and *fully connected* layers. The most important ones are convolutional layers.

Convolutional layers receive a *tensor* with shape $k \times n \times m \times c$, which means that the layer receives $k \in \mathbb{N}$ inputs of sizes $n \times m \times c$. After the convolutional layer, the input has shape $k \times n' \times m' \times c'$, where n' can be different from n (respectively m', c'). One convolutional layer can apply one or several filters to the same input, producing different 2-dimensional activation maps that are stacked along the depth dimension to produce the output.

Formally, a convolution is the process of adding each element of the image to its local neighbours, weighted by a kernel. That is in our case performing a dot product between the input and the filter. We obtain

$$g(x, y) = \omega \star f(x, y) = \sum_{dx=-a}^a \sum_{dy=-b}^b \omega(dx, dy) f(x + dx, y + dy),$$

where $g(x, y)$ is the pixel (x, y) of the filtered image, $f(x, y)$ is the same pixel in the original image, and ω is the filter kernel. Depending on the filter kernel, the result will be different.

Example 8. If we want to produce noise reduction in an image, we use Gaussian Blur, which kernel is calculated by a Gaussian function like the one in Equation 1. For instance, a 3×3 gaussian kernel would approximately be:

$$\frac{1}{16} \begin{bmatrix} 1 & 2 & 1 \\ 2 & 4 & 2 \\ 1 & 2 & 1 \end{bmatrix}.$$

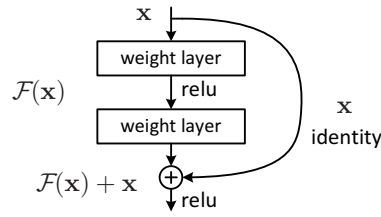
11.2.1 ResNet

To finish with this introduction to CNNs, we will present a widely used CNN architecture. Since CNNs appeared, people tried to build such deep networks that were very hard to train. When deeper networks start converging, a degradation problem occurred: with the network depth increasing, accuracy gets saturated and then it degraded rapidly. *ResNet* (He *et al.*, 2015) brought an end to this problem, allowing the training of very deep models with up to hundreds of layers.

In the paper, they mentioned that if $\mathcal{H}(x)$ is the underlying map of a sequence of layers and knowing that it can be asymptotically approximated, we can also approximate its residuals

$$\mathcal{F}(x) = \mathcal{H}(x) - x.$$

Then, the original function becomes $\mathcal{F}(x) + x$. In order to achieve this, they introduced the *Residual Block*, a block of layers that introduces skip or short-cut connections that makes it easy for networks to represent the identity mapping. We can see one of this blocks depicted in Figure 10.

Figure 10: Residual Block. Image from He *et al.* (2015)

Usually, ResNet comes accompanied by a number (e.g. *Resnet-50*), that indicates the number of layers it has. Usually, a *batch normalization* is used after each convolution and ReLU is used at the end of each group of layers.

Layers	Output Size	Resnet 18
conv1	112×112	7×7 , 64, stride 2
conv2	56×56	3×3 max pool, stride 2
		$\begin{bmatrix} 3 \times 3, & 64 \\ 3 \times 3, & 64 \end{bmatrix} \times 2$
conv3	28×28	$\begin{bmatrix} 3 \times 3, & 128 \\ 3 \times 3, & 128 \end{bmatrix} \times 2$
conv4	14×14	$\begin{bmatrix} 3 \times 3, & 256 \\ 3 \times 3, & 256 \end{bmatrix} \times 2$
conv5	7×7	$\begin{bmatrix} 3 \times 3, & 512 \\ 3 \times 3, & 512 \end{bmatrix} \times 2$
	1×1	average pool, 1000-d fc, softmax

Table 1: Resnet 18 architecture.

For instance, having a look at the Table ??, in the convolution block *conv2*, we have two convolutions with kernel size 3×3 and depth size 64. This would be one of the blocks, represented in Figure 10. However, since at the right hand side of the matrix we have a $\times 2$, it means that we will have two of this blocks in this convolution group.

11.3 DATA AUGMENTATION

In the general problem of performing a DL task, it may happen that the amount of data, either labeled or unlabeled, is not enough to give the model the number of examples that it needs to learn.

Definition 11.3.1. *Data augmentation* are techniques used to increase the amount of data by adding modified copies of the already existing data or newly created synthetic data also from the existing one.

In our self-supervised problem, data augmentation gains even more importance. It has been empirically shown that choosing the appropriate techniques

and the order of the application of them to the images can cause a huge improvement on the results (Chen *et al.*, 2020b).

Depending on the kind of data that we are facing, the techniques change. For instance, if we are working in the field of Natural Language Processing (NLP), which involves applying ML to texts, we should use techniques such as *Back translation* or *Synonym replacement*. If we try to process audio we can use *Noise Injection* or changing the *Speed* of the audio.

In our case, we will be working with images. Let us present some examples that can be used for data augmentation in the computer vision field.

Rotations and flips

A very common data augmentation technique is the *rotation* of the image a certain amount of degrees, usually one of the following: $\{\frac{\pi}{2}, \pi, \frac{3}{2}\pi\}$, using the center of the image as a rotation center. To do this, we only have to consider the following transformation:

$$\begin{pmatrix} x' \\ y' \end{pmatrix} = \begin{pmatrix} \cos(\theta) & -\sin(\theta) \\ \sin(\theta) & \cos(\theta) \end{pmatrix} \begin{pmatrix} x \\ y \end{pmatrix},$$

where (x, y) is the pixel of the image that we will rotate and θ are the degrees that the image will be rotated.

This rotation can be composed (or not) with a *flip* operation. That operation consists in applying a symmetry respect the central column of an image. Clearly, a flip can be applied on its own and it is still generating a new example.



(a) Tulips



(b) Rotated, flipped tulips

Figure 11: The image on the left is the original image, and the image on the right is a new example generated by first a rotation of a random angle and then a flip.

Random crop with resize

Another type of example generation is taking a crop randomly out of the image, making sure that the whole crop stays inside the image, and then resizing it back to the original image size.

If an image is sized $n \times m$ with $n, m \in \mathbb{N}$, this method consists in selecting a rectangle of size $n' \times m'$ with $n' \leq n, m' \leq m$ of the original image and then resizing it back via upsampling¹.



Figure 12: The image on the left is the original image, and the image on the right is a new example generated by performing a random crop of the image and then resizing the crop to the original image size.

This can be useful to obtain images that are similar to each other, for instance, taking two crops of the same image. However, this could also lead to the model to learn an specific kind of data that is not relevant in general: the image data histogram. Let us see the image data histograms for both the original image and the random cropped image:

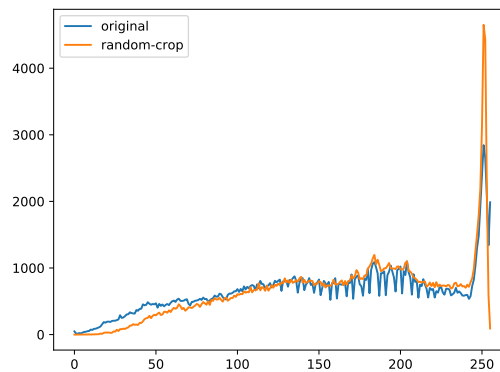


Figure 13: Histogram for both the original image on blue and the random crop on orange.

As we can see, the histograms of both the images are quite similar, so a DL model could learn to identify similar objects making use only on this data histogram, which would not be very effective since the same object can have many different histograms. Because of this, this random crops are useful but the random cropped images have to be applied more functions in order to be really useful for the training of our models.

¹ Upsampling is increasing the size of an image by inserting new rows and columns in the image matrix and interpolating the values of the introduced pixels using the pixels that we already had in the image.

Random color jitter

This kind of data augmentation is not always used, but it is very interesting in our case. The random crops produced very similar data histograms since the colors in the original image and the cropped image are very similar. We can try to change this by randomly modifying the value of the pixels by changing the brightness, saturation or the contrast of the image.

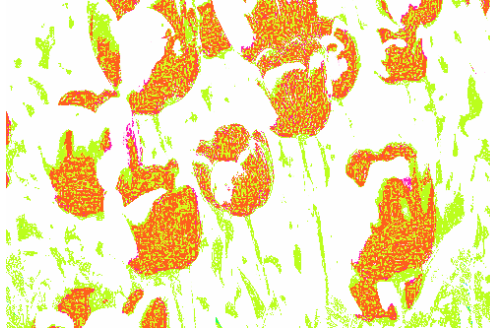


Figure 14: Color jitter applied to tulips image.

We would like to see if there is a difference between the histograms of the cropped image compared with the color jittered image. The result is the following:

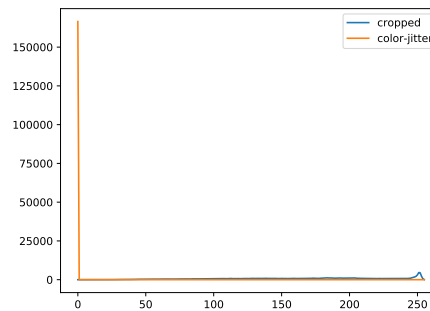


Figure 15: Histograms of the color jitter image in orange and cropped image in blue.

As we can see, since in the color jittered image most of the pixels turn white, the histograms get a big difference this time specially in the left side. This way, if we use both augmentation composed, we will avoid our model to learn only from the data histogram of our images.

12 | SIMCLR

12.1 INTRODUCTION

Until this point of the work, we have been presenting the theoretical basis of representation learning using contrastive learning. Previous approaches, such as the framework presented in Oord *et al.* (2019), use a generative approach as a part of the representation learning process. Although this can be beneficial at some points and, in fact, achieved the *state-of-art*¹ empirical results, we have to consider that generative models have some drawbacks.

Let us set in the case of learning representations of images to present a very simple example. In this case, generative models must *generate* each pixel on the image. This can be extremely computationally expensive.

Until now, we had been trying to minimize the loss in Equation (13), which we proved that maximizes a lower bound in the mutual information. However, some papers such as Chen *et al.* (2020b), Tschannen *et al.* (2020), suggest that it is unclear if the success of their methods is caused by the maximization of mutual information between the latent representations, or by the specific form that the contrastive loss has.

In fact, in Tschannen *et al.* (2020) they provide empirical proof for the loose connection between the success of the methods that use MI maximization and the utility of the MI maximization in practice. They also empirically proof that the encoder architecture can be more important than the estimator used to determine the MI.

Even with the empirically proved disconnection between MI maximization and representation quality, recent works that have used the loss function defined in Equation (13) have obtained state-of-art results in practice.

12.2 THE FRAMEWORK

SimCLR (Chen *et al.*, 2020b) presents a framework that achieved state-of-art results when it was presented in July 2020. It also uses contrastive learning in an specific way that we will present later.

This framework learns representations by maximizing agreement between examples of the same input obtained by using data augmentation on the input example and a contrastive loss in the latent space.

¹ *State-of-art* refers to the best results that have been achieved at some point of time.

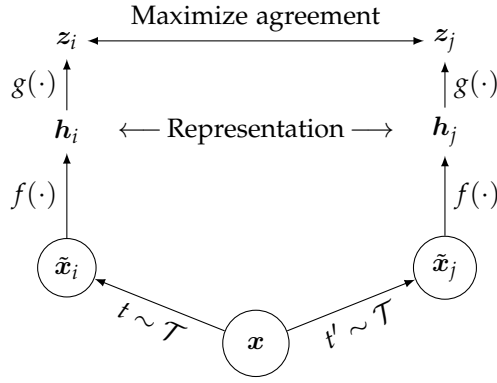


Figure 16: Figure obtained from [Chen et al. \(2020b\)](#). A simple framework for contrastive learning of visual representations.

The framework that is presented follows a linear structure. Figure 16 depicts it. Let us provide with deeper explanation. We will present the steps in a general way and later we will remark the specific considerations that were used in the implementation of the framework for the experiments.

The steps followed are:

1. Firstly, using the input x and two augmentation functions $t, t' \in \mathcal{T}$, two different views \tilde{x}_i, \tilde{x}_j are obtained using data augmentation. They are both sampled from the same family of data augmentations \mathcal{T} . They are *both* considered as positive views.
2. Secondly, a NN base encoder $f(\cdot)$ is used to extract representations for the two different views, obtaining $h_i = f(\tilde{x}_i)$, where $h_i \in \mathbb{R}^d$.
3. Then, a *small* neural network projection $g(\cdot)$ is used. This neural network maps the representations h_i to the space where contrastive loss is applied. Hence, we obtain

$$z_i = g(h_i) = W^{(2)}\sigma(W^{(1)}h_i),$$

where σ is a nonlinear function and $W^{(i)}$ are the weights matrix.

4. Lastly, the contrastive loss is used for a contrastive prediction task. Using a set $\{\tilde{x}_k\}$ that includes a pair of positive examples \tilde{x}_i, \tilde{x}_j , the contrastive loss will (as we have already been doing theoretically) try to identify \tilde{x}_j in the set for a given \tilde{x}_i . It is important to remark how the contrastive loss is used in this framework:
 - a) A minibatch of N samples is randomly taken from the training set.
 - b) Using the N samples, we augment each pair to obtain $2N$ data points. The idea is, given a positive pair, use the other $2(N - 1)$ as negative examples.
 - c) We define

$$\text{sim}(u, v) = \frac{u^T v}{\|u\| \|v\|},$$

the normalized dot product between u and v . This function is also known as the *cosine similarity*. Then, the loss function for a positive pair of examples takes the form:

$$\ell_{i,j} = -\log \frac{\exp(\text{sim}(z_i, z_j)/\tau)}{\sum_{k=1}^{2N} \mathbb{1}_{k \neq i} \exp(\text{sim}(z_i, z_k)/\tau)}, \quad (21)$$

where $\mathbb{1}_{k \neq i} \in \{0, 1\}$ produces 0 if $k = i$ and 1 elsewhere. τ is a temperature parameter that has to be adjusted.

- d) The final loss is computed across all the positive pairs, both (i, j) and (j, i) in the minibatch.

$$\mathcal{L} = \frac{1}{2N} \sum_{k=1}^{2N} (\ell(2k-1, 2k) + \ell(2k, 2k-1)). \quad (22)$$

Remark 7. The loss in Equation (22) is just another formulation of the one in Equation (13), applied to this particular way of obtaining positive and negative views.

We can summarize all this steps in the following algorithm:

Algorithm 1: SimCLR's learning algorithm.

```

input: batch size  $N$ , temperature  $\tau$ , structure of  $f, g, \mathcal{T}$ .
for sampled minibatch  $\{\mathbf{x}_k\}_{k=1}^N$  do
  for all  $k \in \{1, \dots, N\}$  do
    draw two augmentation functions  $t \sim \mathcal{T}, t' \sim \mathcal{T}$ 
    # use augmentations
     $\tilde{\mathbf{x}}_{2k-1} = t(\mathbf{x}_k)$ 
     $\tilde{\mathbf{x}}_{2k} = t'(\mathbf{x}_k)$ 
    # use  $f$ 
     $\mathbf{h}_{2k-1} = f(\tilde{\mathbf{x}}_{2k-1})$  # representation
     $\mathbf{h}_{2k} = f(\tilde{\mathbf{x}}_{2k})$  # representation
    # use  $g$ 
     $\mathbf{z}_{2k-1} = g(\mathbf{h}_{2k-1})$  # projection
     $\mathbf{z}_{2k} = g(\mathbf{h}_{2k})$  # projection
  end for
  for all  $i \in \{1, \dots, 2N\}$  and  $j \in \{1, \dots, 2N\}$  do
     $s_{i,j} = \mathbf{z}_i^\top \mathbf{z}_j / (\|\mathbf{z}_i\| \|\mathbf{z}_j\|)$  # compute similarity
  end for
  update networks  $f$  and  $g$  to minimize  $\mathcal{L}$  defined in Eq. (22)
end for
return encoder network  $f(\cdot)$ , and throw away  $g(\cdot)$ 

```

Figure 17: Algorithm that summarizes the learning process that SimCLR follows.

Considerations:

1. In SimCLR, the data augmentation is applied sequentially, and using three techniques of image data augmentation that are very simple and common:

- a) Random cropping followed by resize to original size,
 - b) Random color distortions, consisting in changing the value of certain amount of pixels randomly, like color jittering or changing from colored images to black and white images.
 - c) Random Gaussian Blur: Gaussian Blur consists in convolving an image with a *Gaussian function*.
2. For the base encoder $f(\cdot)$, a *ResNet* network is used because of its simplicity.
 3. For the projection head $g(\cdot)$, a MLP with one hidden layer is used, and *ReLU* is used as the activation function σ .

12.3 FINDINGS OF SIMCLR

Using the algorithm structure that we have presented, the experiments focused on finding what made the biggest impact on the representation learning task.

Intuitively, since we are comparing $2N$ images, either positive or negative, the batch size that we will use for the experiment should be as big as possible. Comparing a positive sample and trying to push it apart from as much negative samples as possible in the same iteration sounds like a good idea for the quality of our representations.

A few of the most relevant findings using SimCLR are:

- In the SimCLR framework, multiple data augmentations are used and it is empirically shown that this improves the contrastive prediction tasks that yield effective representations.
- The introduction of a nonlinear transformation that it can be learnt during the learning process. This linear transformation is between the representation and the contrastive loss, so before evaluating the loss function, the representation is applied the nonlinear function.
- The contrastive loss benefits from normalized embeddings and also from a temperature parameter that has to be adjusted.
- This loss also benefits from larger batch sizes and longer training, as well as from deeper and wider networks.

Later, *SimCLRv2* was presented (Chen *et al.*, 2020a), achieving a new state-of-art in semi-supervised learning on the ImageNet dataset. In this new framework, semi-supervised learning is use . This approach involves unsupervised learning followed by supervised fine-tuning². After this stages that are similar to SimCLR, a final stage is added. It consists in *distillation* with unlabeled examples for refining the task-specific knowledge.

In the self-supervised (or unsupervised) part, the images are used without its labels, so that the representations of them that are obtained by the framework are not related to an specific task. Using this technique, Chen *et al.* (2020a) shows how using deeper and wider neural networks for both self-supervised pretraining and fine-tuning improves accuracy greatly. Also, the

importance of the projection head $g(\cdot)$ is remarked in this new framework.

² *Fine-Tuning* refers to the process of re-training certain parts of an already trained NN so that it focuses its learn on specific examples, such as a different dataset.

13 | BOOSTRAP YOUR OWN LATENT

13.1 MOTIVATION

We have shown how contrastive methods rely on comparing different views of the same image with views of other images, considering the ones from the same image as positive and the rest as negative samples. This is why they are called *self-supervised* methods.

Self-supervised methods build upon the cross-view prediction framework, i.e., learning representations by predicting different views (or data augmentations) of the same image. This could lead the frameworks to collapsed representations, such as a constant representation for any view of the image can easily help us to identify objects, but it is useless for downstream tasks.

The methods presented earlier, which use contrastive learning, avoid this problem by trying to discriminate between positive and negative views, as we have already explained.

However, some papers (e.g. [Caron et al. \(2019\)](#)) have already raised the following question:

Is the use of negative samples necessary to avoid collapsing? .

This question is studied in [Grill et al. \(2020\)](#), paper that presents an algorithm that we will be deeply explaining.

The first solution that comes to mind to prevent the collapsing problem would be to use a randomly initialized network to produce the targets of the predictions. As it is probably expected, due to its randomness, it does not produce good representations for downstream tasks. Nonetheless, the representations obtained were empirically much better than initial fixed representation, so it could be interesting to refine this representation in order to make it better for the later tasks. This is the intuitive idea behind *Bootstrap your own latent* (BYOL) ([Grill et al., 2020](#)).

13.2 BYOL ALGORITHM

BYOL's algorithm has certain similarities with the SimCLR framework that we presented in Chapter [12](#). The goal of this framework is to learn a representation for an input. In this case, the representation will be noted as y_θ .

For this purpose, two neural networks are used:

- An *online* network defined by a set of weights θ .

- A *target* network, defined by a different set of weights ξ .

They both have the same structure, composed of three stages:

1. An encoder f_γ ,
2. A projector g_γ ,
3. A predictor q_γ ,

where $\gamma \in \{\theta, \xi\}$. In the online network, despite their different name, the projector g_θ and the predictor q_θ have the same architecture.

Remark 8. The projector g_γ is used because in SimCLR (Chen *et al.*, 2020b) is proven empirically that this projection improves the general performance of the framework.

The *difference* between them is that the target network provides the regression targets to train the online network, and its parameters ξ are an exponential moving average¹ of the online parameters θ . Mathematically, given a rate decay $\tau \in [0, 1]$, after each training step ξ is updated as follows:

$$\xi \leftarrow \tau \xi + (1 - \tau) \theta$$

Having presented the networks and its structure, the following steps are followed in BYOL's framework:

1. The input that both networks receive is different, even if it comes from the same image. As in SimCLR, given an input image x and two distributions of data augmentation for images, $\mathcal{T}, \mathcal{T}'$, two views are produced from x to get $v = t(x), v' = t'(x)$ where $t \sim \mathcal{T}$ and $t' \sim \mathcal{T}'$.
2. Each produced view is passed to one of the networks. In particular, the first view v is passed to the online network, and follows the next steps:

$$x \mapsto v = t(x) \mapsto y_\theta = f_\theta(v) \mapsto z_\theta = g_\theta(y_\theta)$$

where f_θ, g_θ are the ones that we mentioned before in the structure of each networks. This online network outputs y_θ and z_θ .

In a similar process, the the target network is passed the second view v' , which follows the next steps:

$$x \mapsto v' = t'(x) \mapsto y'_\xi = f_\xi(v') \mapsto z'_\xi = g_\xi(y'_\xi)$$

where, again, f_ξ, g_ξ are the ones mentioned before.

3. Then, using the online network, a prediction $q_\theta(z_\theta)$ is produced. Remark that the prediction is *only* applied to the online network.
4. Having $q_\theta(z_\theta)$ in the online network and z'_ξ in the target network, they are both ℓ_2 -normalized to

$$\overline{q_\theta}(z_\theta) = \frac{q_\theta(z_\theta)}{\|q_\theta(z_\theta)\|} \quad \text{and} \quad \overline{z'_\xi} = \frac{z'_\xi}{\|z'_\xi\|}.$$

¹ A moving average is a calculation to analyze data points by creating a series of averages in different subsets of the data.

5. Now, we can define the mean squared error between the normalized prediction $\overline{q_\theta}(z_\theta)$ and the normalized projection $\overline{z'_\xi}$:

$$\mathcal{L}_{\theta,\xi} = \left\| \overline{q_\theta}(z_\theta) - \overline{z'_\xi} \right\|_2^2.$$

6. If we stopped in the step 5, the framework would be asymmetric between the two networks, since the projection is performed in one of the views, v but not in the other one, v' . To fix this, the process described is repeated except that now v' is the input of the online network and v is the input of the target network, producing a new loss $\tilde{\mathcal{L}}_{\theta,\xi}$. The final loss is computed as:

$$\mathcal{L}_{\theta,\xi}^{\text{BYOL}} = \mathcal{L}_{\theta,\xi} + \tilde{\mathcal{L}}_{\theta,\xi} \quad (23)$$

The loss in Equation (23) is the one that must be optimized stochastically. Furthermore, since ξ depends on θ , if η is the learning rate that we want to apply, the optimization problem can be expressed as follows:

$$\begin{cases} \theta \leftarrow \text{optimizer}(\theta, \nabla_{\theta} \mathcal{L}_{\theta,\xi}^{\text{BYOL}}, \eta), \\ \xi \leftarrow \tau \xi + (1 - \tau) \theta \end{cases}$$

The framework can be summarized in the following figure:

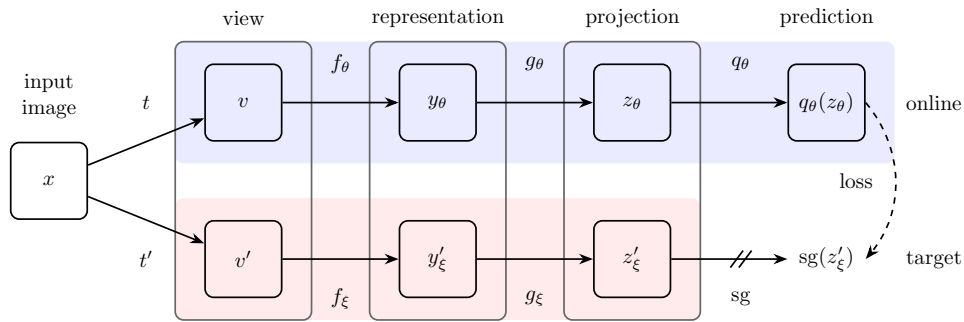


Figure 18: Image from (Grill *et al.*, 2020). Overview of Bootstrap Your Own Latent Framework.

The *sg* (stop gradient) indicates that, at that point, the gradient is not propagated back.

After the whole model has been fully trained, both the projection g_θ and the prediction q_θ are discarded, since what it is interesting for us is the representation of the input, and it is what it will be used in downstream tasks.

Some considerations have to be made about this framework:

1. About data augmentation: In BYOL, it is not important if the model learns about the histogram of the image, since we want to keep any information captured by the target representation into its online network. In the original paper, it is empirically shown how BYOL is more robust to the choice of the image augmentations than SimCLR.

2. Since the weights of the target network ξ depend on the weights of the online network θ by a parameter τ , the first mentioned weights ξ represent a delayed and stable version of θ . Actually, if the decay rate τ is 1, we never update ξ , and if $\tau = 0$, then the weights of the target network are always being updated. This way, a trade-off is tablished between updating ξ too often or too slowly. Empirically, it has been shown that the most adequate values that yield to more stable results are $\tau \in [0.9, 0.999]$.

13.3 RESULTS OBTAINED BY THIS FRAMEWORK

In this framework, the impact of the batch size varies respect to the importance that it had in SimCLR. In contrast with BYOL's framework, in SimCLR we were using negative samples in our network to learn how to produce the representations. Because of this, BYOL's is expected to be more robust to smaller batch sizes.

It is empirically shown that the regularization of the weights is crutial in the self-supervised setting. The experiments in the paper showed that removing the weight decay regularization led to model divergence.

BYOL achieved the state-of-art performance results in the linear evaluation. This is, again, one of the most interesting parts since it helps us to evaluate how good the representations that our model create are.

Method	Architecture	Params	Top-1	Top 5
SimCLR	ResNet-50(2×)	94M	74.2	92.0
BYOL	ResNet-50(2×)	94M	77.4	93.6
SimCLR	ResNet-50(4×)	375M	76.5	93.2
BYOL	ResNet-50(4×)	375M	79.6	94.8

Table 2: Comparison between SimCLR and BYOL Top1 and Top5 accuracies using the same architectures on the ImageNet dataset.

BYOL also outperformed other models such as *MoCo* (He *et al.*, 2020b) and a second version of the Contrastive Predictive Coding framework (?), but they have not been included in the table since they are not deeply studied in this work.

Part V

EXPERIMENTS

14 | INTRODUCTION

In this chapter, we will explain the fundamentals and technologies that have been used for the experimentation. We will focus on three main aspects:

1. The used dataset.
2. The used libraries for the development of the code.
3. Analysis of the code itself.
4. The used metrics to evaluate the obtained results.

The idea of the experimentation part is to test and compare the frameworks that we have presented in Chapters 12 and 13. Their architectures have already been explained, and the original code for both backbones has not been done by me.

This work will focus on testing how changing the training hyperparameters of the model affects the final results, since the original papers [Chen *et al.* \(2020b\)](#); [Grill *et al.* \(2020\)](#) already mention that using their structure, the results are affected by those hyperparameters, such as batch size, network depth or network width.

The implementations that have been used can be found in:

- SimCLR implementation: Official implementation from Google in <https://github.com/google-research/simclr/tree/master/tf2>
- BYOL implementation: not official, found in <https://github.com/gardner14/byol-tensorflow2>.

Although there exists an official implementation of BYOL, non-official one has been chosen because it uses *Tensorflow*, which makes the implementations easier to understand and modify, which we needed to do. Also, the idea is to make use of *Tensorboard*, a Tensorflow utility that helps with visualization and graph generation of the training and final results.

14.1 HARDWARE AND BASIC LIBRARIES

When running training experiments in DL, the used hardware is one of the most determining factors for the results obtained by the models. There are many reasons for this, such as the time spent on the training a model or the amount of data that we can fit in the GPUs (which will be our case) or TPUs.

For the experiments of this project, the DECSAI¹ department of the University of Granada generously provided us access to a server that has a few NVIDIA GeForce RTX 3090. This model of GPU is one of the best in the market, having a *compute capability* of 8.6/10, rated by the NVIDIA company. IT also has a memory of 24GB, which allows us to experiment with rela-

tively large batch sizes using a single GPU and not having to parallelize the experiment.

To be able to use this GPUs in our experiments, two basic libraries have to be installed in the server:

1. `tensorflow-gpu`. This is a variant of the `tensorflow` library that was developed for using GPUs while using `tensorflow`.
2. `CUDA 11.4`. `CUDA` is a parallel computing platform and API created by `NVIDIA` which allows to use a `CUDA-enabled` GPU for general purpose processing. In other words, `CUDA` is needed to be able to use the GPU in our python scripts.

These two libraries, along with other packages that are needed for creating graphics (`Matplotlib`, `Seaborn`) or computing metrics (`sklearn`) are installed using a Conda enviroment in our server's user.

14.2 THE DATASET: CIFAR10

The computational resources that we have for the experiment are limited. Due to this, we must fix a dataset that, having enough and representative examples, allows us to achieve feasible training time and successful results.

One of the ever most used dataset, which was also used in both `SimCLR` and `BYOL` papers, is `CIFAR10` ([Krizhevsky, n.d.](#)). This dataset will be used to test the overall performance of our representation learning methods.

`CIFAR10` contains 60.000 images divided in 10 classes, where each class contains 10.000 images. The size of the images is $32 \times 32 \times 3$, so the size of the images is not very large. This helps us to have faster trainings.

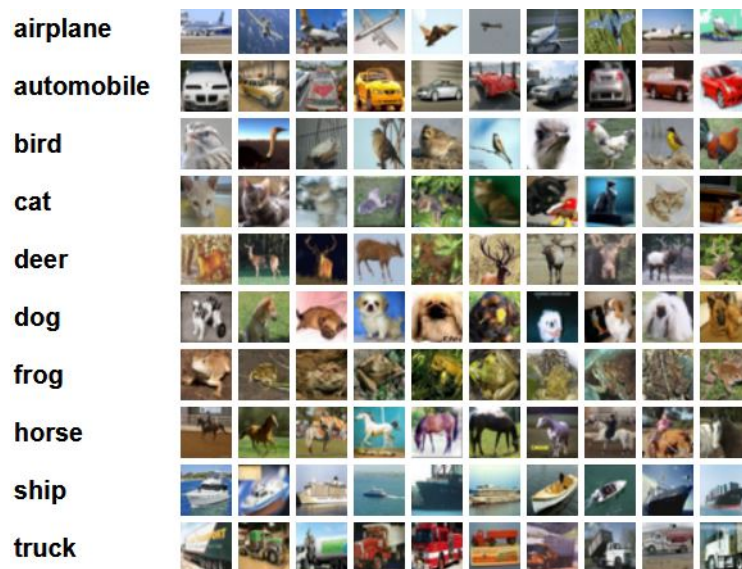


Figure 19: Ten examples of each class in the `CIFAR10` dataset.

¹ The DECSAI's website is <https://decsai.ugr.es/>.

This dataset has 50.000 samples for training and 10.000 for test. The test batch contains the same number of examples of each of the 10.000 classes in the dataset, that is, it contains 1.000 examples of each class.

It is important to remark that the classes are *completely mutually exclusive*. That means that there is no overlap between the classes even if they have similar images, such as *Cars* and *Truck*, which are two of the classes of the dataset.

14.3 TENSORFLOW

Tensorflow² is an open source library for developing machine learning frameworks.

It can be used for many tasks, but it focuses on training and inference of deep neural networks. It is used for both research and production at *Google*, since it was also developed by the *Google Brain* team for internal use. However, it was later released as open source.

Figure 20: Tensorflow logo.



The creation of new models is very simple, offering multiple abstraction levels. This is why it is suitable for our experiments. Also, the code is most of the times easily understandable.

There are other libraries that are widely used in machine learning algorithm development, such as *Pytorch* or *Jax*. However, Tensorflow has been chosen because I was a little bit more familiarized than with the other and because of its simplicity and how common it is.

There are a few ways to define a NN or a framework using tensorflow. The most classic one is defining a sequential model using Keras, a tensorflow API that defines layers of a neural network and helps with the implementation of simple NN structures. Let us see how to implement a very simple example of a NN with three *Dense* layers:

```
1 model = keras.Sequential(
2     [
3         layers.Dense(2, activation="relu", name="layer1"),
4         layers.Dense(3, activation="relu", name="layer2"),
5         layers.Dense(4, name="layer3"),
6     ]
7 )
```

Another way of creating models using tensorflow is by defining a single step of training using `tf.GradientTape()` and then executing the single step multiple times in a loop. Using `GradientTape`, tensorflow performs automatic differentiation, which is needed for the minimization process. Let us see

² Tensorflow documentation can be found at <https://www.tensorflow.org/>.

the simplest example, consider the function $f(x) = x^2$, and imagine that we want to obtain $f'(3)$. We can obtain it using `tf.GradientTape()` as follows:

```
1 x = tf.constant(3.0)
2 with tf.GradientTape() as g:
3     g.watch(x)
4     y = x * x
5 dy_dx = g.gradient(y, x)
```

In our case, the gradient is obtained and then applied to the optimizer by using:

```
1 with tf.GradientTape() as tape:
2     grads = tape.gradient(loss, model.trainable_variables)
3     optimizer.apply_gradients(zip(grads, model.trainable_variables))
```

14.3.1 Tensorboard

Tensorboard is a Tensorflow's visualization kit. It provides the visualization and tooling needed for machine learning experimentation. Among its more important utilities, we can find:

- Tracking and visualizing metrics (such as loss, accuracy, entropy) not only during the training but also when the training time has ended.
- Visualizing the model graph: ops and layers.
- Visualizing histograms of weights, biases and how tensors change during the training.
- Projecting high-dimensional data to a lower dimensional space.
- Displaying images, text and audio data.

Also, it is very easy to integrate with tensorflow. Actually, in most of the cases it is as simple as adding the following *callback* when we fit the model:

```
1 tensorboard_callback = tf.keras.callbacks.TensorBoard
2                       (log_dir=log_dir, histogram_freq=1)
3 model.fit(x=x_train, y=y_train,
4           epochs=5,
5           validation_data=(x_test, y_test),
6           # The added callback produces the magic!
7           callbacks=[tensorboard_callback])
```

In our case, there are some differences, since we are not using the standard `fit` function to train our models. Because of this, we have to log the information that we have obtained in each step. To do this, we can use the `metrics` python package to group them (as it is done in the SimCLR code), or we can just directly save the information creating a *file writer* and writing the desired variables on this file. We do this in the modification that we have done to BYOL's original code as follows:

```
1 train_summary_writer = tf.summary.create_file_writer(args.log_dir)
2 with train_summary_writer.as_default():
3     tf.summary.scalar('top_1_acc', float(acc), step=epoch)
4     tf.summary.scalar('top_5_acc', float(top_5_acc), step=epoch)
5     tf.summary.scalar('loss', float(losses[-1]), step=epoch)
```

14.4 METRICS

As we have seen, firstly, our models create a representation of the input image and then this representation is evaluated using a supervised linear head. This is the most interesting part, since we can see if the representation obtained was really useful for the classification task. We need to present the measures that we will use to measure how good the representations that we are producing are.

Notation 1. We will address the true positives (the positive samples of a class classified correctly) as TP , the true negatives (the negative samples classified correctly) as TN , the false positives (the negative samples classified as positive ones, which is a mistake of our model) as FP and the false negatives (the positive samples classified incorrectly as negative samples) as FN .

Using this notation, the measures that we will be using are the following:

- *Accuracy.* The classic measure for models. It measures the number of successes our network has obtained producing the correct label for the representation created during the unsupervised part of the network. Formally,

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN}.$$

In SimCLR, we will use two kinds of accuracies:

1. *Top1* accuracy, which is the ordinary accuracy.
2. *Top5* accuracy, which measures if any of the 5 highest probability answers matches the true label.

15 | EXPERIMENTATION

We are now ready to perform the experiments. We will begin exploring SimCLR implementation and results, later we will explore BYOL's implementation and lastly we will compare them in order to see how BYOL tries to improve SimCLR and we will check if it successess or not.

The code used for the experimentations, as well as some files with the results, can be found on the Github repository for this work.

15.1 SIMCLR EXPLORATION

We will perform an iterative exploration with this framework. We will explore a few range for a subset of the hyperparameters and then we will go deeper into some hyperparameters to try and obtain better results.

15.1.1 First approach

The first thing we did to experiment with this framework is to explore a wide range of hyperparameters to see which set of them performed better for us. The script that can be found in `code/SimCLR/run.py`.

What we did for this first exploration was to define a *parameter grid* and execute the whole framework using each combination of the parameters. The parameters that were firstly considerered are:

1. `batch_size`. This is one of the most important parameters of SimCLR. In the original paper, it was proven that the higher this parameter, the better the results obtained for the linear classification. This parameter is also important since it has to adapt to our GPU's memory. The options that we have considered for the first experiments are: `batch_size = {512,1024}`.
2. `temperature`. The temperature parameter τ plays an important role in the individual loss seen in Equation (21). It is suggested to try values in the range $[0,1]$ and the one that had better performance for the original experiments was around 0.5, so we chose the following values: `temperature = {0.25,0.5,0.75,1}`.
3. `color_jitter_strength`. This parameter measures how hard is the color variation in the data augmentation. Previous results show that this parameter is important for the success of the network, so we provide with a big range of values. We include the following values: `color_jitter_strength = {0.25,0.5,0.75,1}`.

Using python and the python function `itertools.product` we straightforwardly generate all possible different combinations of unions of the parameters so we only have to append them to a general string and execute the `run.py` script mentioned before to obtain the results. In total, we obtain 32 possible combinations, so we will obtain 32 models.

The script was executed and took approximately 24 hours to train and evaluate all the different models, obtaining the results in Table ?? in Appendix B. We have to remark, for each of the two possible batch sizes, the best results obtained. These are:

batch_size	temperature	color_jitter	regularization_loss	top_1_accuracy	top_5_accuracy	steps
512	0.25	0.25	0.0093	0.833	0.994	9800
1024	0.25	0.75	0.0093	0.841	0.995	4900

Table 3: Best results for the grid search experiment with SimCLR.

Remark 9. On the original paper, the Top 1 accuracy score reported for CIFAR10 in 100 epochs was $\sim 84.6\%$ accuracy for batch size 512 and $\sim 85.1\%$ accuracy score for 1024 batch size. If we have a look at the results obtained in our experiment, we can say that our first attempt was *successfull*.

Let us make some observations. The clearest is that, since the second batch size doubles the first one, the training ends in half the steps. We can see that the temperature obtained in both cases is the same, and it is equal to 0.25. However, there is a big change in the color jitter parameter: while using 512 as batch size obtains the best result with the value of 0.25 for color jitter, using 1024 as batch size this value changes to 0.75, which implies stronger color changes on the data augmentation.

In general, having a look again at the Table ?? in Appendix B, we can see that:

- As it also happens in the experiments made on the original SimCLR paper, lower temperature τ parameters cause higher accuracies on the linear heads of our models.
- The models perform better when the `color_jitter` parameter is in the range $[0.5, 0.75]$. This means that, in general, a generous amount of color jittering to our picture is beneficial for the models.
- The *Top 5 accuracy* is above 99% in each model. This says that, since the models almost always give the correct label one of the 5 highest values, but they obtain the correct tag around 84% of the times, they might be creating some similar representations of the data for inputs that do not belong to the same class.

Observations about the batch size

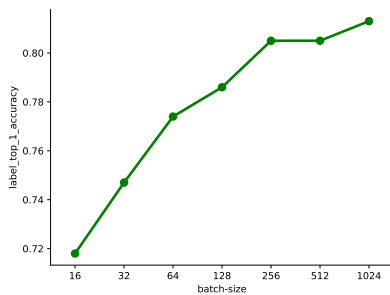
Although both models get to the same value of regularization loss, the model with a batch size of 1024 obtains a higher *top 1 accuracy*. This seems to follow the intuitive idea that we presented before: SimCLR benefits from bigger batch sizes since it allows a positive sample to be compared and pushed appart from negative samples. Ideally, we would keep pushing the batch size

parameter forward, doubling it again to try to find out if the performance of the model keeps benefiting from bigger batch sizes. However, this requires either using multiple GPUs(or TPUs) working in parallel on the training or using GPUs with bigger memory, which we do not have access to.

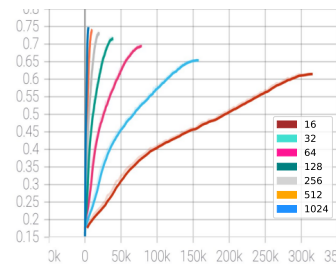
Using the resources we have, the following experiment was performed: using the hyperparameters that we found out in our *GridSearch* to be the ones that achieve the best results except for the batch-size, we train models moving the batch size in the following range:

$$\text{batch-size} = \{16, 32, 64, 128, 256, 512, 1024\},$$

where the last two values had already been computed in the GridSearch so we do not have to repeat the training.



(a) Evolution of the Top 1 accuracy with the batch size.



(b) Evolution curves of the accuracy with the steps taken in the training.

Figure 21: Results of the batch-size experiment.

Figure 21a shows the clear improvement of the Top 1 accuracy score when increasing the size of the batch that is used to compute the loss of our model. Increasing from 16 to 1024 gives an increase of more than 8% in Top 1 accuracy, which is a lot. Further increase can be obtained if we keep increasing the batch size, but this is beyond our computation possibilities.

Figure 21b supports what we have just presented: not only with a smaller batch size we have to do thousands of extra steps in the training, but also we obtain much better accuracy performance of the linear head of the model.

Observations other hyperparameters

We have seen that in our first experiment, batch size has been really relevant in the final results. We would like to see if the rest of the parameters play such an important role as well.

As we have seen before, the parameters that we want to see how they affect the models are the `color_jitter_strength` and the temperature τ parameter. Let us study their impact on the model one by one.

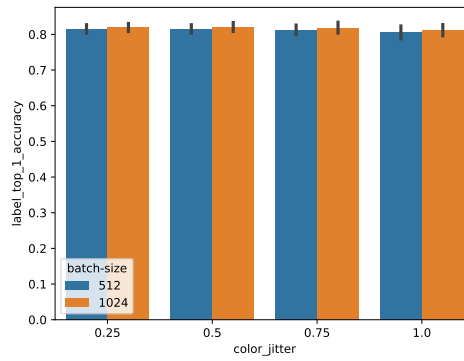


Figure 22: Accuracy score following the color jitter parameter and both batch sizes considered.

As we can see in Figure ??, the `color_jitter_strength` parameter does not have a huge impact on the performance of the linear head of the model. The black bars on the center of each orange/blue bar indicate the variance of the accuracy respect the rest of the parameters. Surely, there are centesimal differences between the different values of this parameter. However, more finetuning might be needed in order to make this parameter relevant for our model. The low influence may as well be caused by the small encoder network (ResNet 18), so we will explore if there are any changes when we change the encoder network later in the document.

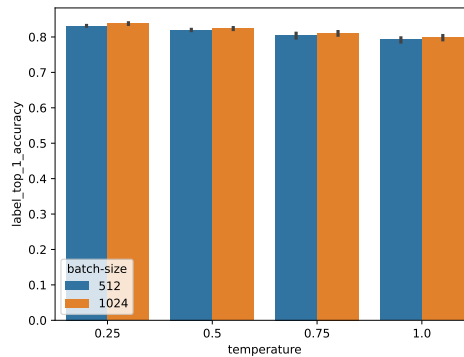


Figure 23: Accuracy score following the temperature and both batch sizes considered.

Part VI

APPENDIX

This appendix will be used to set forth some theoretical results that might not always be relevant but are needed to understand some details during this thesis. Not all of them will be proven.

Proposition 9 (Jensen's Inequality). *Let $f : \mathcal{D} \rightarrow \mathbb{R}$ be a concave function and $n \in \mathbb{N}$. For any $p_1, \dots, p_n \in \mathbb{R}_0^+$ with $\sum p_i = 1$ and any $x_1, \dots, x_n \in \mathcal{D}$, it holds that:*

$$\sum_{i=1}^n p_i f(x_i) \leq f\left(\sum_{i=1}^n p_i x_i\right).$$

Furthermore, if f is strictly concave and $p_i \geq 0$ for all $i = 1, \dots, n$, then the equality holds if, and only if, $x_1 = \dots = x_n$.

In Chapter 10, the norm $\|\cdot\|_2$ is mentioned. Norm theory is a very extensive field, so we will only mention the definition and the norm that we will use in the text.

Definition A.1. Given a vector space X over a subfield F of the complex numbers \mathbb{C} , a *norm* is a real valued function $\|\cdot\| : X \rightarrow \mathbb{R}$ with the following properties:

1. Triangle inequality, that is: $\|x + y\| \leq \|x\| + \|y\|$ for all $x, y \in X$.
2. Absolute homogeneity, that is: $\|sx\| = |s| \|x\|$ for all $x \in X$ and any scalar s .
3. Positive definiteness, that is $\|x\| \geq 0$ for all $x \in X$ and $\|x\| = 0$ if, and only if, $x = 0$.

In particular, the $\|\cdot\|_2$ that we used in the euclidean space \mathbb{R}^n , is defined as follows:

$$\|x\|_2 := \left(\sum_{i=1}^n x_i^2\right)^{\frac{1}{2}}, \quad \forall x \in \mathbb{R}^n.$$

B | APPENDIX B

In this appendix we will insert some experimental results that have been obtained during this work, but have not been presented before since not all of them are completely relevant. The most relevant ones have already been presented, so this section is not needed for complete understanding of this work.

B.1 SIMCLR EXPERIMENTS

The results obtained for the first experiment, which was presented in Section 15.1.1, are presented in table ??

batch_size	temperature	color_jitter	regularization_loss	top_1_accuracy	top_5_accuracy	steps
512	0.25	0.25	0.0093	0.833	0.994	9800
		0.5	0.0089	0.832	0.993	9800
		0.75	0.0086	0.831	0.994	9800
		1.0	0.008	0.83	0.992	9800
	0.5	0.25	0.0136	0.819	0.993	9800
		0.5	0.0124	0.822	0.993	9800
		0.75	0.0121	0.821	0.993	9800
		1.0	0.0118	0.817	0.992	9800
	0.75	0.25	0.0161	0.809	0.993	9800
		0.5	0.015	0.812	0.993	9800
		0.75	0.0141	0.805	0.99	9800
		1.0	0.0137	0.793	0.99	9800
	1.0	0.25	0.017	0.798	0.99	9800
		0.5	0.0163	0.797	0.99	9800
		0.75	0.016	0.793	0.99	9800
		1.0	0.0155	0.782	0.989	9800
1024	0.25	0.25	0.0103	0.836	0.993	4900
		0.5	0.0097	0.839	0.995	4900
		0.75	0.0093	0.841	0.995	4900
		1.0	0.009	0.835	0.993	4900
	0.5	0.25	0.0149	0.822	0.993	4900
		0.5	0.0133	0.827	0.994	4900
		0.75	0.0132	0.826	0.994	4900
		1.0	0.0128	0.82	0.993	4900
	0.75	0.25	0.0168	0.814	0.991	4900
		0.5	0.0166	0.816	0.992	4900
		0.75	0.0157	0.813	0.993	4900
		1.0	0.0153	0.802	0.989	4900
	1.0	0.25	0.0175	0.806	0.991	4900
		0.5	0.0174	0.802	0.992	4900
		0.75	0.017	0.794	0.99	4900
		1.0	0.0164	0.79	0.989	4900

Table 4: All results for first experiment using SimCLR.

BIBLIOGRAPHY

- Bengio, Yoshua, Courville, Aaron, & Vincent, Pascal. 2014. Representation Learning: A Review and New Perspectives. *arXiv:1206.5538 [cs]*, Apr. arXiv: 1206.5538.
- Caron, Mathilde, Bojanowski, Piotr, Joulin, Armand, & Douze, Matthijs. 2019. *Deep Clustering for Unsupervised Learning of Visual Features*.
- Chen, Ting, Kornblith, Simon, Swersky, Kevin, Norouzi, Mohammad, & Hinton, Geoffrey. 2020a. *Big Self-Supervised Models are Strong Semi-Supervised Learners*.
- Chen, Ting, Kornblith, Simon, Norouzi, Mohammad, & Hinton, Geoffrey. 2020b. A Simple Framework for Contrastive Learning of Visual Representations. *arXiv:2002.05709 [cs, stat]*, June. arXiv: 2002.05709.
- Cosma, Rohilla Shalizi. 2021. *Advanced Data Analysis from an Elementary Point of View*.
- Cover, T. M., & Thomas, Joy A. 1991. *Elements of information theory*. Wiley series in telecommunications. New York: Wiley.
- Deng, Li. 2014. Deep Learning: Methods and Applications. *Foundations and Trends® in Signal Processing*, 7(3-4), 197–387.
- Doersch, Carl, Gupta, Abhinav, & Efros, Alexei A. 2016. Unsupervised Visual Representation Learning by Context Prediction. *arXiv:1505.05192 [cs]*, Jan. arXiv: 1505.05192.
- Grill, Jean-Bastien, Strub, Florian, Altché, Florent, Tallec, Corentin, Richemond, Pierre H., Buchatskaya, Elena, Doersch, Carl, Pires, Bernardo Avila, Guo, Zhaohan Daniel, Azar, Mohammad Gheshlaghi, Piot, Bilal, Kavukcuoglu, Koray, Munos, Rémi, & Valko, Michal. 2020. *Bootstrap your own latent: A new approach to self-supervised Learning*.
- Gutmann, Michael U, & Hyvarinen, Aapo. Noise-Contrastive Estimation of Unnormalized Statistical Models, with Applications to Natural Image Statistics. 55.
- He, Kaiming, Zhang, Xiangyu, Ren, Shaoqing, & Sun, Jian. 2015. *Deep Residual Learning for Image Recognition*.
- He, Kaiming, Fan, Haoqi, Wu, Yuxin, Xie, Saining, & Girshick, Ross. 2020a. Momentum Contrast for Unsupervised Visual Representation Learning. *arXiv:1911.05722 [cs]*, Mar. arXiv: 1911.05722.
- He, Kaiming, Fan, Haoqi, Wu, Yuxin, Xie, Saining, & Girshick, Ross. 2020b. *Momentum Contrast for Unsupervised Visual Representation Learning*.

- Hjelm, R. Devon, Fedorov, Alex, Lavoie-Marchildon, Samuel, Grewal, Karan, Bachman, Phil, Trischler, Adam, & Bengio, Yoshua. 2019. Learning deep representations by mutual information estimation and maximization. *arXiv:1808.06670 [cs, stat]*, Feb. arXiv: 1808.06670.
- Hénaff, Olivier J., Srinivas, Aravind, Fauw, Jeffrey De, Razavi, Ali, Doersch, Carl, Eslami, S. M. Ali, & van den Oord, Aaron. 2020. *Data-Efficient Image Recognition with Contrastive Predictive Coding*.
- Krizhevsky, Alex. Learning Multiple Layers of Features from Tiny Images. 60.
- Larochelle, Hugo, & Murray, Iain. The Neural Autoregressive Distribution Estimator. 9.
- Löwe, Sindy, O'Connor, Peter, & Veeling, Bastiaan. 2019. Putting an End to End-to-End: Gradient-Isolated Learning of Representations. *Pages 3039–3051 of: Advances in Neural Information Processing Systems*.
- Mikolov, Tomas, Chen, Kai, Corrado, Greg, & Dean, Jeffrey. 2013. Efficient Estimation of Word Representations in Vector Space. *arXiv:1301.3781 [cs]*, Sept. arXiv: 1301.3781.
- Mitchell, Tom M. 1997. *Machine Learning*. McGraw-Hill series in computer science. New York: McGraw-Hill.
- Oord, Aaron van den, Li, Yazhe, & Vinyals, Oriol. 2019. Representation Learning with Contrastive Predictive Coding. *arXiv:1807.03748 [cs, stat]*, Jan. arXiv: 1807.03748.
- Poole, Ben, Ozair, Sherjil, Oord, Aaron van den, Alemi, Alexander A., & Tucker, George. 2019. On Variational Bounds of Mutual Information. *arXiv:1905.06922 [cs, stat]*, May. arXiv: 1905.06922.
- Sohn, Kihyuk. 2016. Improved Deep Metric Learning with Multi-class N-pair Loss Objective. *In: NIPS*.
- Tian, Yonglong, Sun, Chen, Poole, Ben, Krishnan, Dilip, Schmid, Cordelia, & Isola, Phillip. 2020. What Makes for Good Views for Contrastive Learning? *arXiv:2005.10243 [cs]*, Dec. arXiv: 2005.10243.
- Tschannen, Michael, Djolonga, Josip, Rubenstein, Paul K., Gelly, Sylvain, & Lucic, Mario. 2020. On Mutual Information Maximization for Representation Learning. *arXiv:1907.13625 [cs, stat]*, Jan. arXiv: 1907.13625.
- Uria, Benigno, Murray, Iain, & Larochelle, Hugo. 2014. RNADE: The real-valued neural autoregressive density-estimator. *arXiv:1306.0186 [cs, stat]*, Jan. arXiv: 1306.0186.
- Wiskott, Laurenz, & Sejnowski, Terrence J. 2002. Slow Feature Analysis: Unsupervised Learning of Invariances. *Neural Computation*, **14**(4), 715–770.