#### Assignment 2

# Learning latent space representations

and application to image generation

Alexandre Vérine, Constant Bourdrez (Benjamin Negrevergne)

PSL University - ENS PSL

## **Outline**

- ① What is a good representation?
- 2 Learning representations with Deep Learning
- 3 Intriguing properties of learned representations
- 4 Synthetic data generation

Task: classify pictures of crocodiles and alligators



Crocodile



Alligator

Task: classify pictures of crocodiles and alligators

**■** Representation 1: features

beast\_color = Light beast\_size = Large ightarrow f

 $\rightarrow$ 

crocodile

Task: classify pictures of crocodiles and alligators

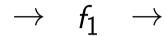
■ Representation 1: features

beast\_color = Dark beast\_size = Small ightarrow  $f_1$  ightarrow

alligator

Task: classify pictures of crocodiles and alligators

■ Representation 1: features



alligator

■ Representation 2: Use raw pixel data



$$\rightarrow$$
  $f_2$   $\rightarrow$ 

crocodile

Task: classify pictures of crocodiles and alligators

■ Representation 1: features

$$ightarrow$$
  $f_1$   $ightarrow$ 

alligator

■ Representation 2: Use raw pixel data

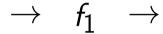


$$\rightarrow$$
  $f_2$   $\rightarrow$ 

alligator

Task: classify pictures of crocodiles and alligators

■ Representation 1: features



alligator

■ Representation 2: Use raw pixel data



$$\rightarrow$$
  $f_2$   $\rightarrow$ 

alligator

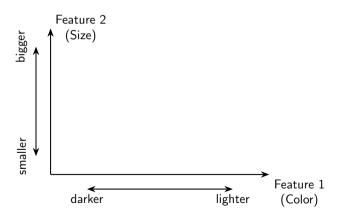
Which **representation** of the input is easier to work with? Why?

 $f_1: \mathbb{R}^2 
ightarrow \mathbb{R}$  Input space 1:  $\mathbb{R}^2$ 

 $f_2: \mathbb{R}^{w \times h \times 3} \to \mathbb{R}$ Input space 2:  $\mathbb{R}^{w \times h \times 3}$ 

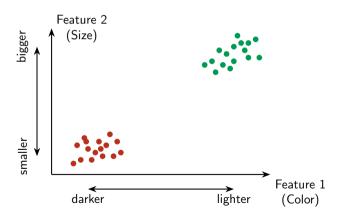
 $f_1: \mathbb{R}^2 o \mathbb{R}$ Input space 1:  $\mathbb{R}^2$ 

 $f_2$ :  $\mathbb{R}^{w \times h \times 3} \rightarrow \mathbb{R}$ Input space 2:  $\mathbb{R}^{w \times h \times 3}$ 



 $egin{array}{ll} f_1 &:& \mathbb{R}^2 & 
ightarrow \mathbb{R} \ & & \\ \hline ext{Input space 1: } \mathbb{R}^2 \end{array}$ 

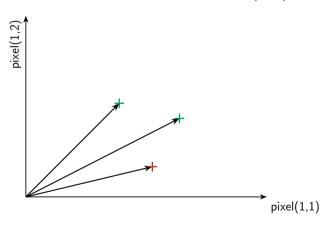
 $f_2: \mathbb{R}^{w \times h \times 3} \rightarrow \mathbb{R}$ Input space 2:  $\mathbb{R}^{w \times h \times 3}$ 



Data points are linearly separable in  $\mathbb{R}^2$ 

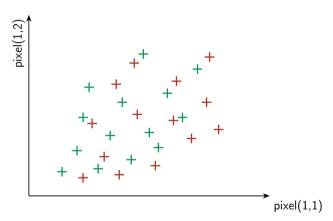
 $f_1: \mathbb{R}^2 o \mathbb{R}$  Input space 1:  $\mathbb{R}^2$ 

 $f_2: \mathbb{R}^{w \times h \times 3} \rightarrow \mathbb{R}$ Input space 2:  $\mathbb{R}^{w \times h \times 3}$ 



 $f_1: \mathbb{R}^2 
ightarrow \mathbb{R}$ Input space 1:  $\mathbb{R}^2$ 

 $f_2: \mathbb{R}^{w \times h \times 3} \rightarrow \mathbb{R}$ Input space 2:  $\mathbb{R}^{w \times h \times 3}$ 



Data points are not linearly separable in  $\mathbb{R}^{w \times h \times 3}$ , Why?

# Pro/cons

#### ■ Representation 1: hand-crafted features

- + Can be processed with simple (linear) models
  - Individual features are highly discriminant
  - Input data points are (almost) linearly separable
- Requires expertise and manual labor to be built
- No extra information in case of ties

# Pro/cons

#### ■ Representation 1: hand-crafted features

- + Can be processed with simple (linear) models
  - Individual features are highly discriminant
  - Input data points are (almost) linearly separable
- Requires expertise and manual labor to be built
- No extra information in case of ties

#### ■ Representation 2: raw pixel data

- + Contains all the information available
- Input data points are not (nearly) linearly separable
   Features are individually non-discriminant
- Difficult to process with simple models

# Pro/cons

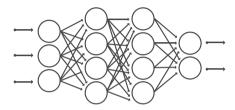
#### ■ Representation 1: hand-crafted features

- + Can be processed with simple (linear) models
  - Individual features are highly discriminant
  - Input data points are (almost) linearly separable
- Requires expertise and manual labor to be built
- No extra information in case of ties

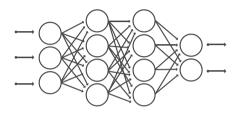
#### ■ Representation 2: raw pixel data

- + Contains all the information available
- Input data points are not (nearly) linearly separable
   Features are individually non-discriminant
- Difficult to process with simple models

Can we build high level features automatically?



$$f = f_1 \circ f_2 \circ \ldots \circ f_{n-1} \circ f_n$$



$$f = f_1 \circ f_2 \circ \ldots \circ f_{n-1} \circ f_n$$

- ullet  $f_1:\mathcal{X} o\mathcal{Z}_1$
- $\bullet \ f_i: \mathcal{Z}_{i-1} \to \mathcal{Z}_i$
- $\bullet \ f_n: \mathcal{Z}_{n-1} \to \mathcal{Y}$

$$f = \underbrace{f_1 \circ f_2 \circ \ldots \circ f_{n-1}}_{g} \circ \underbrace{f_n}_{h}$$

#### Remarks

- h is a linear classifier:  $\mathcal{Z}_{n-1} \to \mathcal{Y}$ 
  - ▶ Data points **must be** linearly separable in  $\mathcal{Z}_{n-1}$
- ullet Data points x are not linearly separable in the input space  ${\mathcal X}$

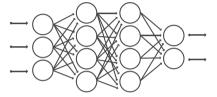
$$f = \underbrace{f_1 \circ f_2 \circ \ldots \circ f_{n-1}}_{g} \circ \underbrace{f_n}_{h}$$

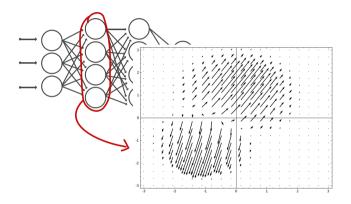
#### Remarks

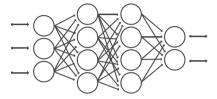
- h is a linear classifier:  $\mathcal{Z}_{n-1} \to \mathcal{Y}$ 
  - ▶ Data points **must be** linearly separable in  $\mathcal{Z}_{n-1}$
- ullet Data points x are not linearly separable in the input space  ${\mathcal X}$

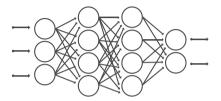
#### What is g?

- ullet a function  $g:\mathcal{X} o\mathcal{Z}_{n-1}$
- such that data points are linearly separable in  $\mathcal{Z}_{n-1}$
- lacktriangle a representation of the point in  ${\mathcal X}$  that is adequate for the task at hand

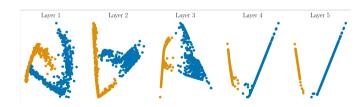












## Deep representations: first lessons

- DNN learn how to projet inputs into a latent space
- The structure of the latent space is useful for the task at hand.
- Given an input  $x \in \mathcal{X}$  we say that g(x) is an **embedding** of x, i.e. continuous vector representations of input data (image, text, graph . . . )

## **Outline**

- What is a good representation?
- 2 Learning representations with Deep Learning
- 3 Intriguing properties of learned representations
- 4 Synthetic data generation

## **Outline**

- ① What is a good representation?
- 2 Learning representations with Deep Learning
- 3 Intriguing properties of learned representations
- 4 Synthetic data generation

# Image embeddings

- *D*: a database with 80 000 pictures.
- $f = g \circ h$ : a classifier trained on object recognition g non-linear function, h linear classifier
- x a random picture from the internet

$$x_1 = \arg\min_{x' \in D} ||g(x) - g(x')||$$
 and  $x_2 = \arg\min_{x' \in D \setminus \{x_1\}} ||g(x) - g(x')||$ 

# Image embeddings

- *D*: a database with 80 000 pictures.
- $f = g \circ h$ : a classifier trained on object recognition g non-linear function, h linear classifier
- x a random picture from the internet

$$x_1 = \underset{x' \in D}{\operatorname{arg\,min}} ||g(x) - g(x')||$$
 and  $x_2 = \underset{x' \in D \setminus \{x_1\}}{\operatorname{arg\,min}} ||g(x) - g(x')||$ 



X

# Image embeddings

- *D*: a database with 80 000 pictures.
- $f = g \circ h$ : a classifier trained on object recognition g non-linear function, h linear classifier
- x a random picture from the internet

$$x_1 = \arg\min_{x' \in D} ||g(x) - g(x')||$$
 and  $x_2 = \arg\min_{x' \in D \setminus \{x_1\}} ||g(x) - g(x')||$ 







 $X_1$ 

 $X_2$ 



X









X



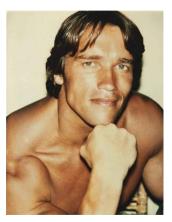




 $x_1$ 

*X*<sub>2</sub>





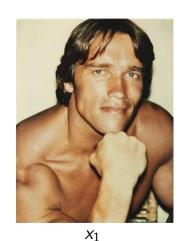


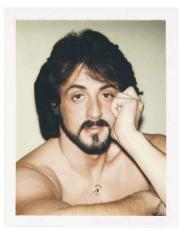
X

 $x_1$ 

*X*<sub>2</sub>







*X*<sub>2</sub>

X

Result: The distance in the latent space seems to be meaningful!

# Word embeddings

We can train (monolingual) *word embeddings*, i.e. representations trained to predict well words that appear in its context (ref here)

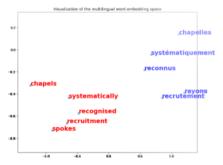


Image generated using pretrained embeddings available here. En avant guiGuan team, 2021.

# Word embeddings

We can train (monolingual) *word embeddings*, i.e. representations trained to predict well words that appear in its context (ref here)

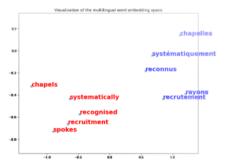
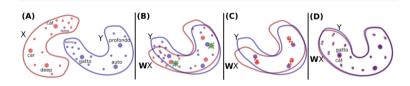


Image generated using pretrained embeddings available here. En avant guiGuan team, 2021.

▶ The structure between word embeddings is preserved across languages

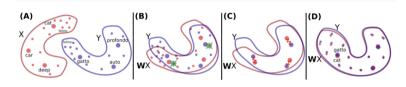
# Word to word translation using word embeddings



Exploit linear transformations and rotations to translate a word.

$$Y = \mathbf{W}X$$

# Word to word translation using word embeddings



Exploit linear transformations and rotations to translate a word.

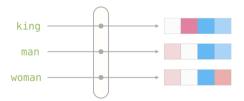
$$Y = \mathbf{W}X$$

### Learn W

- with a parallel corpus (e.g. supervised dataset FR-EN)
- without a parallel corpus using a minmax and a discriminator (GAN-like)

# Text manipulation with word embeddings

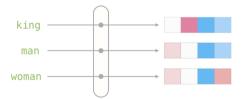
The embedding space is geometric!



The embeddings space is geometric: 
$$v(king) - v(man) + v(woman) = ?$$

# Text manipulation with word embeddings

The embedding space is geometric!



The embeddings space is geometric: 
$$v(king) - v(man) + v(woman) = queen$$

### Visualising the data manifold

Dimensionality reduction by learning an invariant mapping.

Hadsell, R., Chopra, S., LeCun, Y. CVPR (2006)

Learns a mapping g that maps inputs with few controlled variations to a low dimensional space

- all pictures are pictures of planes with different poses
- 9 different elevations and 18 different azimuth (orientation).
- input pictures are projected into a low 3-dimensional space

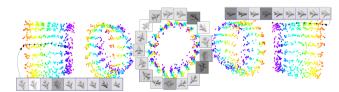
# Visualising the data manifold

Dimensionality reduction by learning an invariant mapping.

Hadsell, R., Chopra, S., LeCun, Y. CVPR (2006)

Learns a mapping g that maps inputs with **few controlled variations** to a **low dimensional space** 

- all pictures are pictures of planes with different poses
- 9 different elevations and 18 different azimuth (orientation).
- input pictures are projected into a low 3-dimensional space

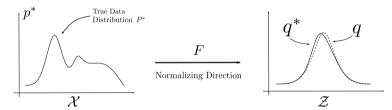


### Result:

- Most input data lies on a well defined subspace of the output space
- There is a clear relation between spacial coordinates and features (elevation, azimuth)

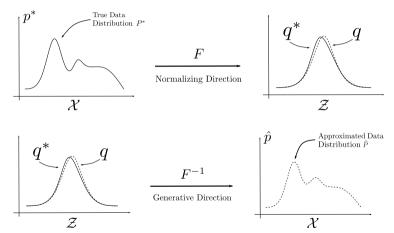
# **Continuous Representation**

What if we enforce continuity and a specific structure in the latent space?



# **Continuous Representation**

What if we enforce continuity and a specific structure in the latent space?



We can generate new data points by sampling in the latent space!

# Visualising the latent space

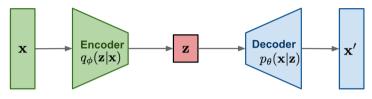
Interpolating between two points in the latent space:

and more examples in this blog post: StyleGAN2 latent space exploration

### **Outline**

- ① What is a good representation?
- 2 Learning representations with Deep Learning
- 3 Intriguing properties of learned representations
- 4 Synthetic data generation

### Data generation with VAE

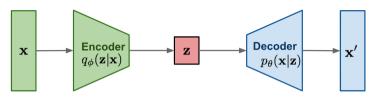


VAE, image from https://lilianweng.github.io

### ■ Main differences with AE

- Use a probabilistic encoder
- Regularize the latent space during training

### Data generation with VAE



VAE, image from https://lilianweng.github.io

### Main differences with AE

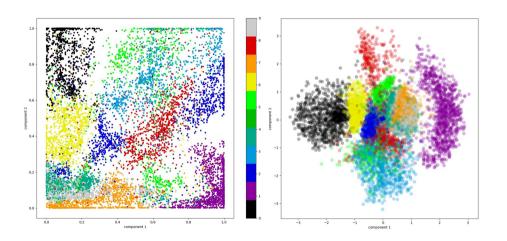
- Use a probabilistic encoder
- Regularize the latent space during training

### ■ Training procedure

- **1** take a training data point x, obtain  $\mu_x$  and  $\sigma_x$  from the encoder
- ② sample  $z \sim \mathcal{N}(\mu_x, \sigma_x)$ , decode z into x'
- 3 compute the loss and update parameters

$$loss(x, x') = ||x - x'|| + KL(\mathcal{N}(\mu_x, \sigma_x), \mathcal{N}(0, I))$$

### AE vs. VAE



# Data generation with VAE

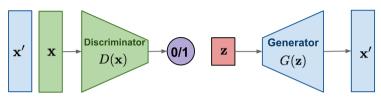
### Pros:

- The latent space is continuous and well structured
- Easy to sample from the latent space

### Cons:

- Generated images are often blurry
- Difficult to train
- Need to tune the weight of the KL term

### **Generative Adversarial Networks**



GAN, image from https://lilianweng.github.io

#### ■ Main difference with VAE

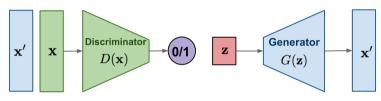
- VAE reconstructs inputs; GANs instead match the data distribution: train a generator G(z) so its samples look real.
- Representation view: G maps a low-dimensional code z to data; a discriminator D(x) provides the learning signal so G(z) and real x share the same statistics.
- Adversarial training (min-max): G tries to fool D, and D learns to tell real from fake.

### **■** Training procedure

G and D are trained simultaneously to solve the following min-max problem:

$$\min_{G} \max_{D} \mathbb{E}_{x_r \sim P}[\log D(x)] + \mathbb{E}_{x_g \sim \hat{P}_G}[\log 1 - D(x)]$$

### **Generative Adversarial Networks**



GAN, image from https://lilianweng.github.io

#### References:

- Generative Adversarial Nets https://arxiv.org/pdf/1406.2661
- Unsupervised Representation Learning with Deep Convolutional Generative Adversarial Networks https://arxiv.org/abs/1511.06434

### Goal of Assignment 2

- Train a GAN on MNIST.
- The structure of the Generator is fixed.
- Use different one or two possible improvements to improve the data generation.

```
0123456789
0123456789
0123456789
0123456789
0123456789
0123456789
```

### Requirements Assignment 2

- Train a Generator (decoder) model to generate MNIST digits.
- Write a script generate.py that generate 10000 samples in the folder samples (use mine).
- Based on these 10k samples, you will be evaluated on FID, Precision and Recall. Precision/Recall

Base repo of the Assignment 2: Link

## Requirements Assignment 2

### Program of next week session:

- Review of different GAN implement for trading off quality of the samples with diversity.
- How to compute FID, Precision and Recall to measure the quality and the diversity of generated samples.
- How to run training on GPU cluster.
- Group session to help you start on Assignment 2.

#### To do list for next week:

Read the GAN paperGenerative Adversarial Nets