# Intrinsic dimension and density profile of neural representations

Diego Doimo and Aldo Glielmo
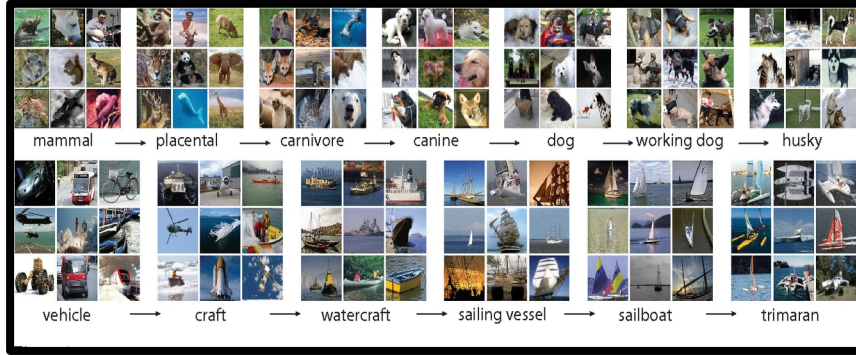
International School for Advanced Studies (SISSA), Trieste
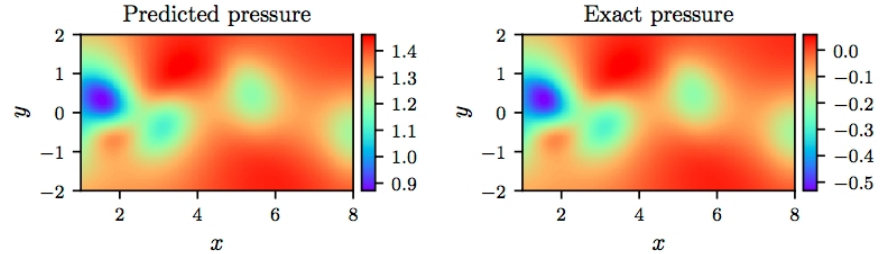
SISSA

UniTS, 27 April 2021

# What can we do with deep learning models?
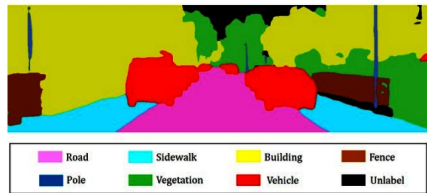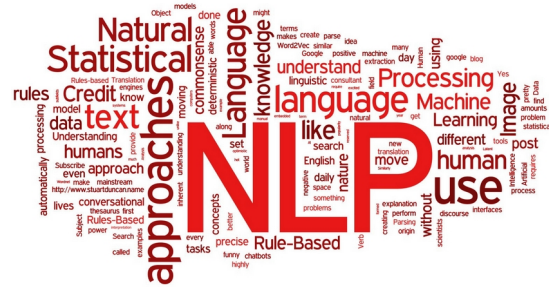
Image classification

Modeling physical systems

Translation, test generation, sentiment analysis, …

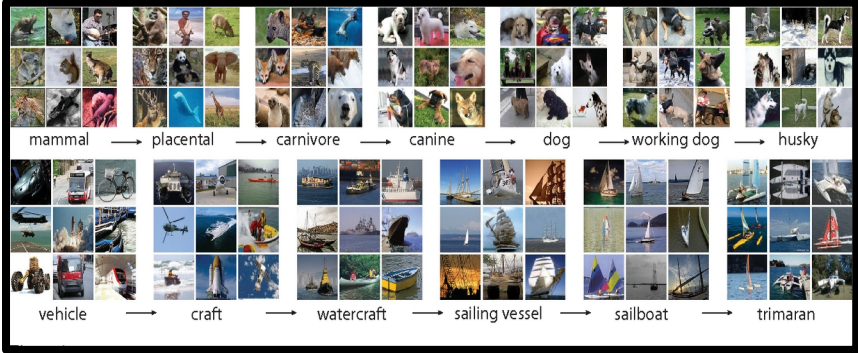Solving constrained optimization
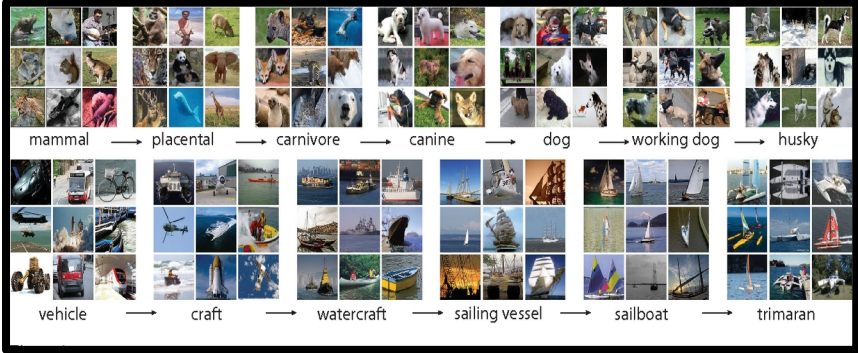
Generative models

Semantic segmentation

# What do we learn with deep learning models?

Image classification

# What do we learn with deep learning models?

Image classification



Convolutional neural networks

# What do we learn with deep learning models?

Image classification



mammal → placental → carnivore → canine → dog → working dog → husky

vehicle → craft → watercraft → sailing vessel → sailboat → trimaran

**Representations**



Input $i$

Pool2

Output

Convolutional neural networks



input    conv1    pool1    conv2    pool2    hidden4    output

Convolution    Subsample    Convolution    Subsample    Full Connection    Convolution

# The importance of representations in neural networks

Representations arise automatically  ➡️  Need to understand their meaning

- To  make NN more interpretable

  Q1:  When / How do interpretable representations arise ?

- To transfer efficiently the learned concepts

  Q2:  Which information is encoded in a given representation?
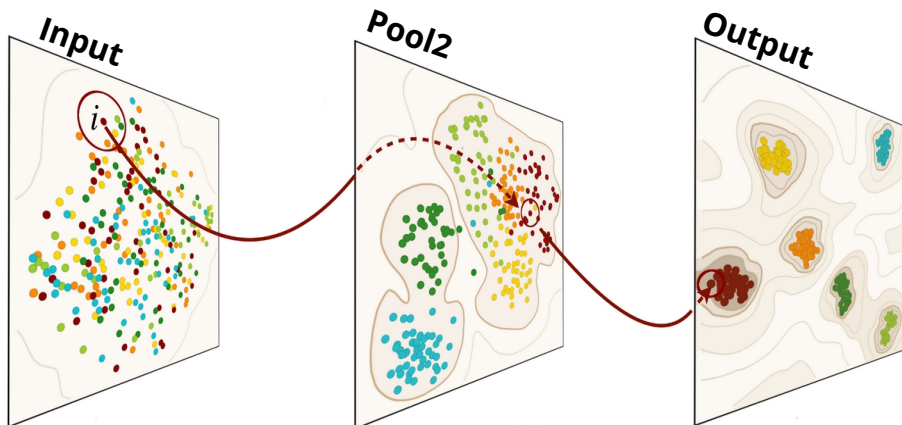
- To improve the architecture design

  What is the depth required to achieve a given performance?

1) **Intrinsic dimension**  [Ansuini et al., NeurIPS 2019]

2) **Probability density**  [Doimo et al., NeurIPS 2020]

**Colab:**

https://colab.research.google.com/drive/1fTxE0GWb5BobZhL3j6G6Ra5hBj__c9X-#scrollTo=VrIL_J3FLQab

# What is a representation?

Representation = function of the input data

$$X_L = f(X_0)$$

f = neural network

-1.61
-1.54
-1.53
-1.51
-1.47
-1.54
-1.56

=

0.00
0.00
0.00
0.95
0.00
0.00
0.00

# What is a representation?



Representation = function of the input data

$$X_L = f(X_0)$$

f = neural network

Input representation

$X_0$

# What is a representation?

Representation = function of the input data

$$X_L = f(X_0)$$

f = neural network

| | |
|---|---|
| -1.61 | |
| -1.54 | |
| -1.53 | |
| -1.51 | |
| -1.47 | |
| -1.54 | |
| -1.56 | |

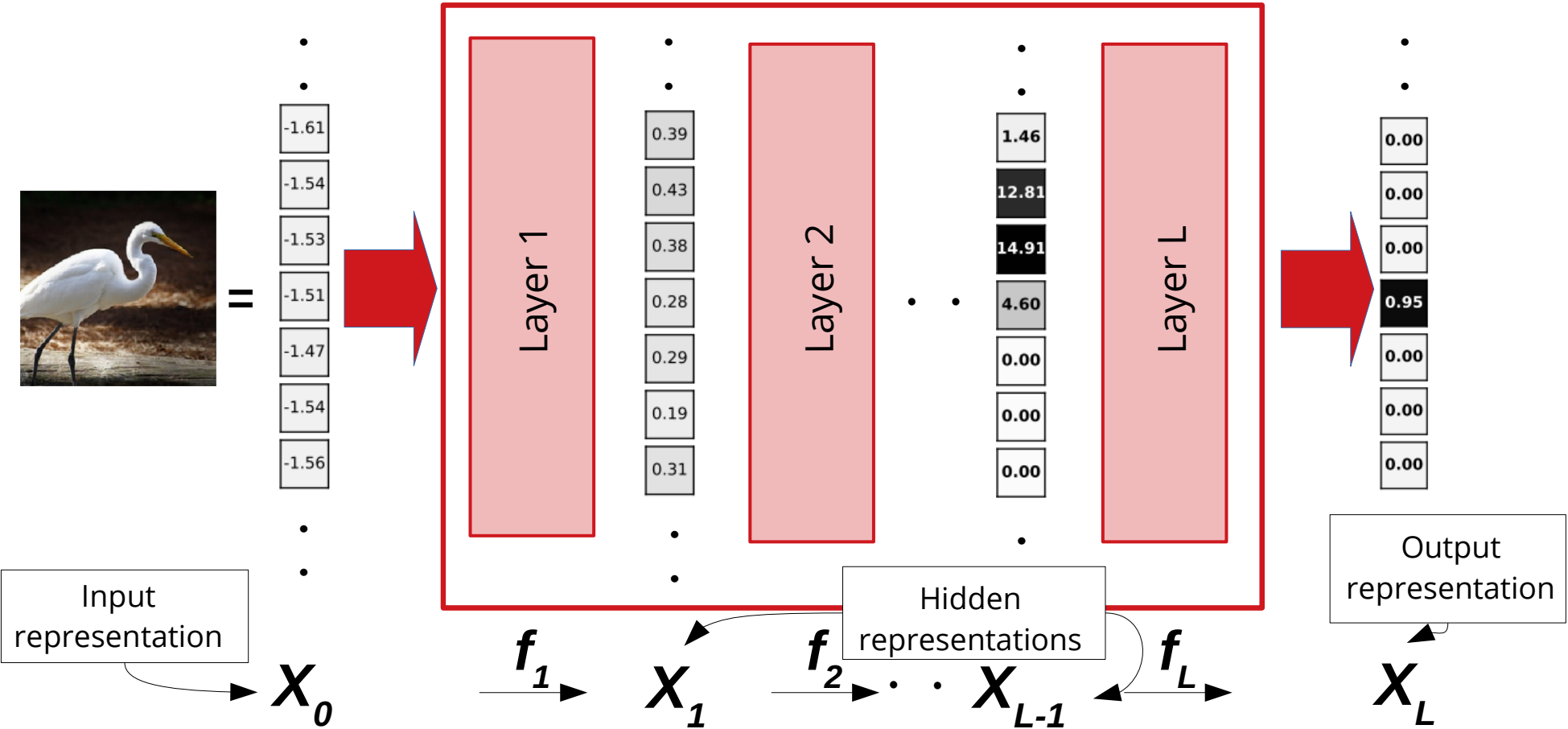| |
|---|
| 0.00 |
| 0.00 |
| 0.00 |
| 0.95 |
| 0.00 |
| 0.00 |
| 0.00 |

Input representation

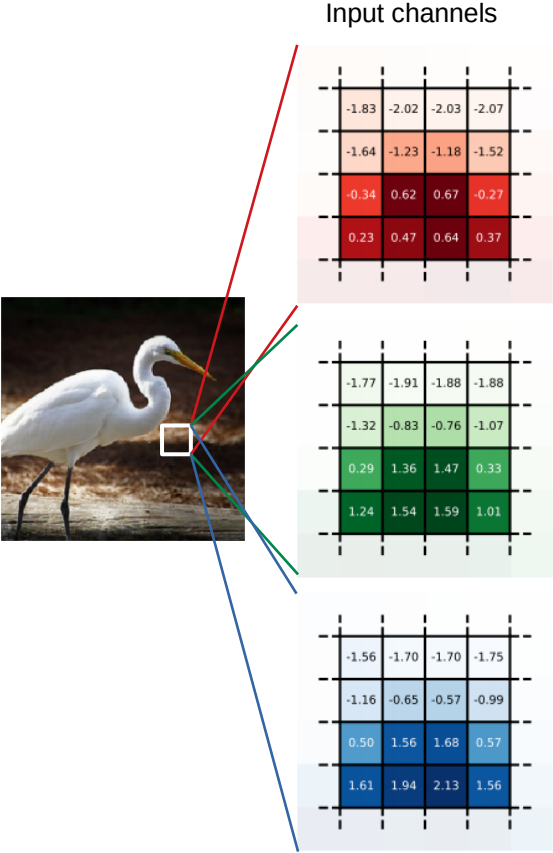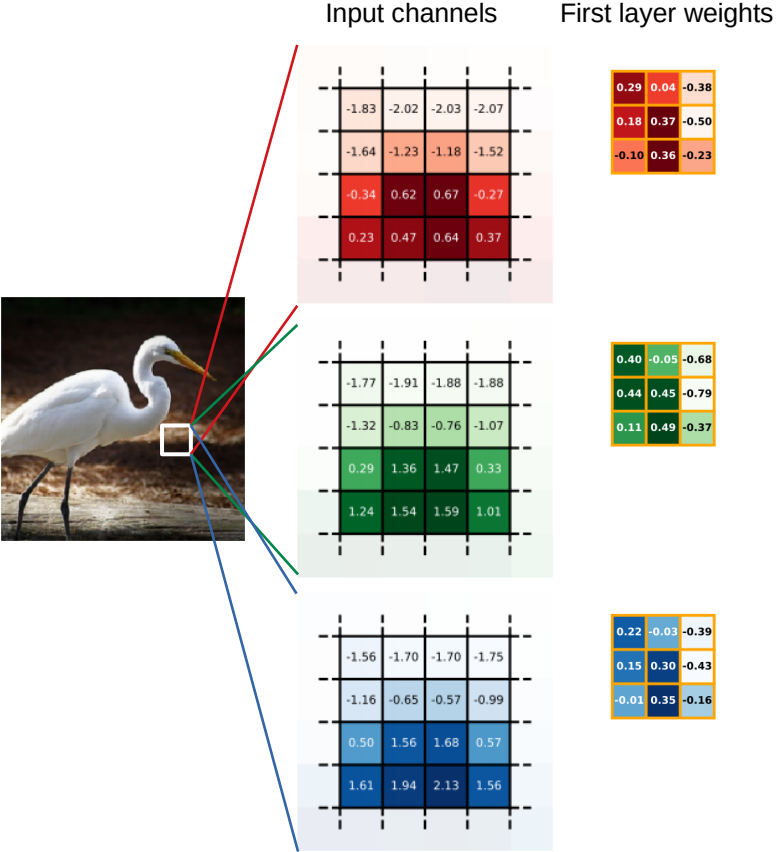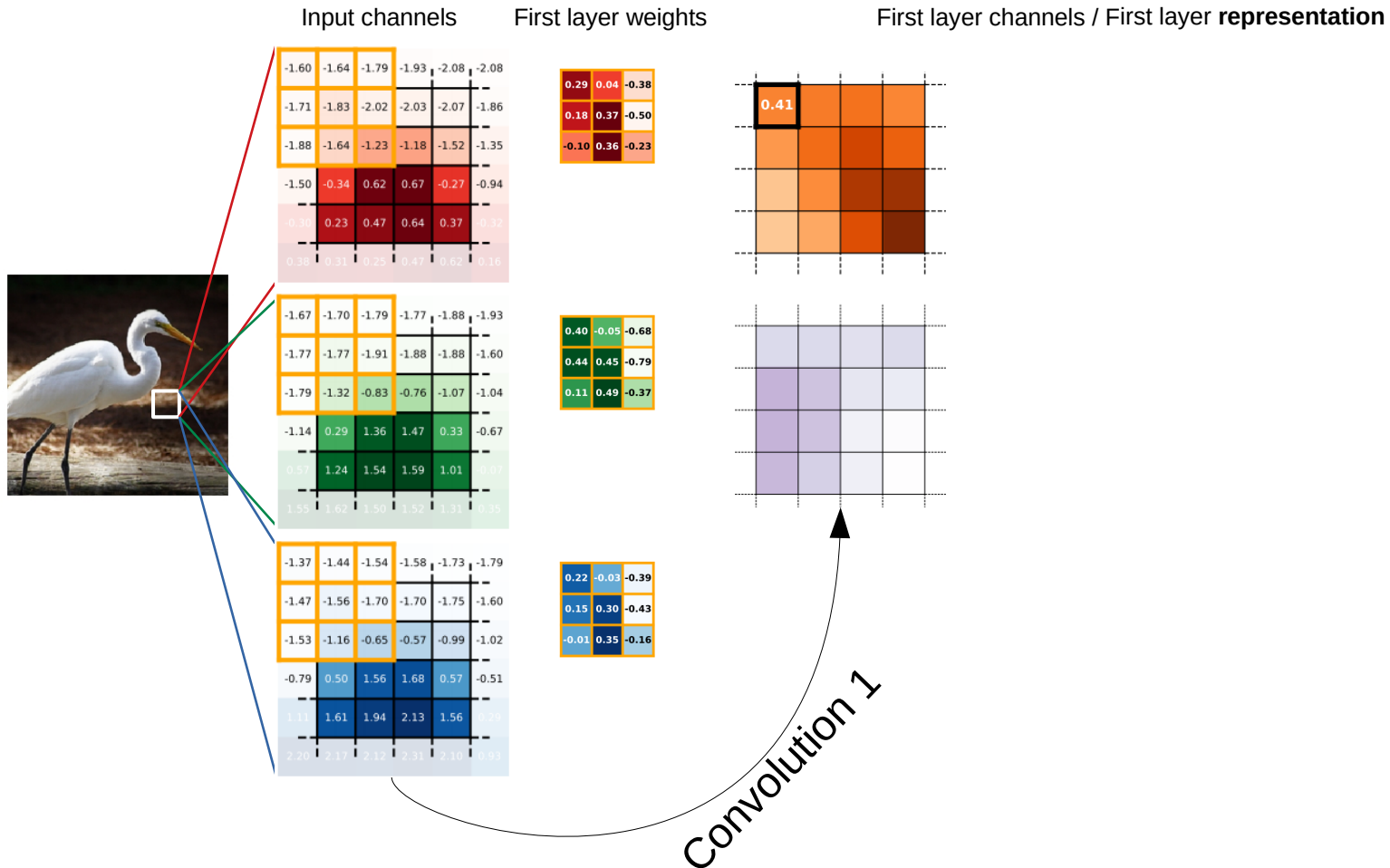$X_0$

Output representation

$X_L$

# What is a representation?

# What is a representation in a convolutional neural network?

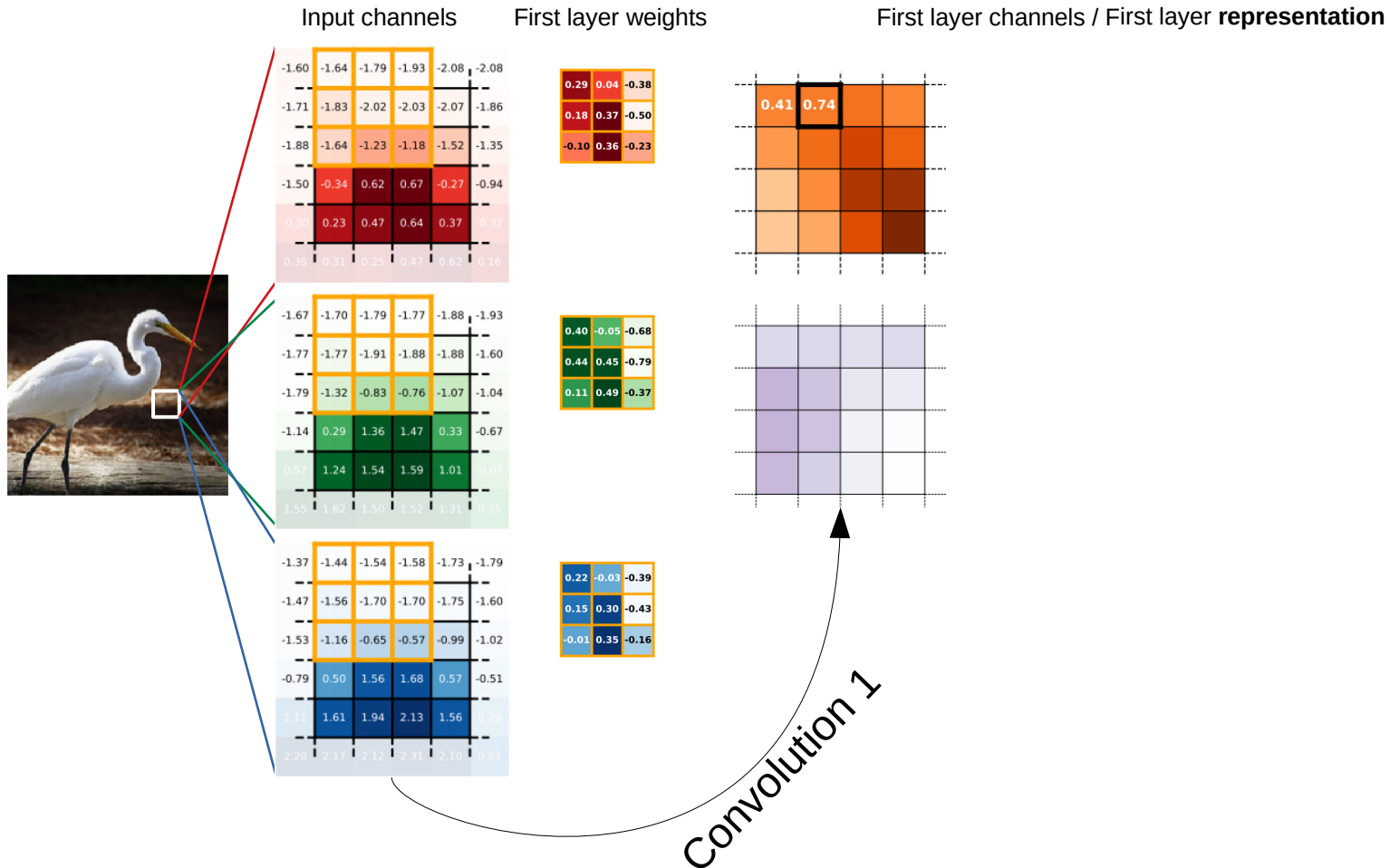# What is a representation in a convolutional neural network?

Input channels

| -1.83 | -2.02 | -2.03 | -2.07 |
|-------|-------|-------|-------|
| -1.64 | -1.23 | -1.18 | -1.52 |
| -0.34 | 0.62  | 0.67  | -0.27 |
| 0.23  | 0.47  | 0.64  | 0.37  |

| -1.77 | -1.91 | -1.88 | -1.88 |
|-------|-------|-------|-------|
| -1.32 | -0.83 | -0.76 | -1.07 |
| 0.29  | 1.36  | 1.47  | 0.33  |
| 1.24  | 1.54  | 1.59  | 1.01  |

| -1.56 | -1.70 | -1.70 | -1.75 |
|-------|-------|-------|-------|
| -1.16 | -0.65 | -0.57 | -0.99 |
| 0.50  | 1.56  | 1.68  | 0.57  |
| 1.61  | 1.94  | 2.13  | 1.56  |

# What is a representation in a convolutional neural network?

Input channels

First layer weights

# What is a representation in a convolutional neural network?

Input channels          First layer weights                    First layer channels / First layer **representation**



Convolution 1

# What is a representation in a convolutional neural network?

Input channels      First layer weights      First layer channels / First layer **representation**



Convolution 1

# What is a representation in a convolutional neural network?

Input channels          First layer weights          First layer channels / First layer **representation**



Convolution 1

# What is a representation in a convolutional neural network?

Input channels    First layer weights    First layer channels / First layer **representation**

# What is a representation in a convolutional neural network?

Input channels

First layer weights

First layer channels / First layer **representation**



Relu 1

Convolution 1

# What is a representation in a convolutional neural network?

Input channels

First layer weights

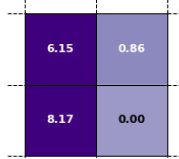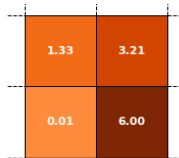First layer channels / First layer **representation**

Convolution 1

Relu 1
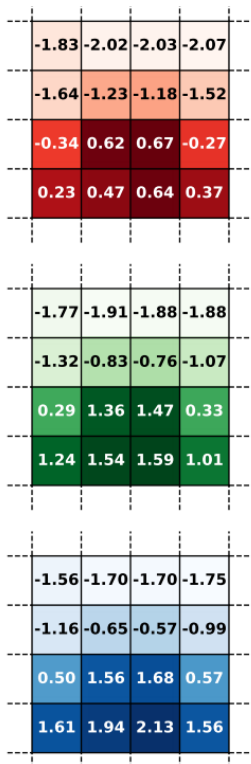
Maxpool

# What is a representation in a convolutional neural network?



Input channels

First layer weights

First layer channels / First layer **representation**

Relu 1

Maxpool

Convolution 1

**First convolutional layer**
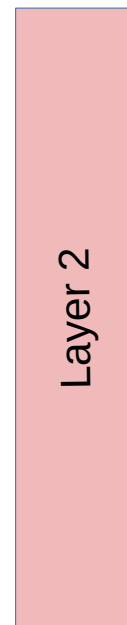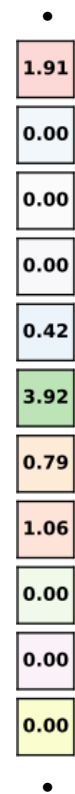
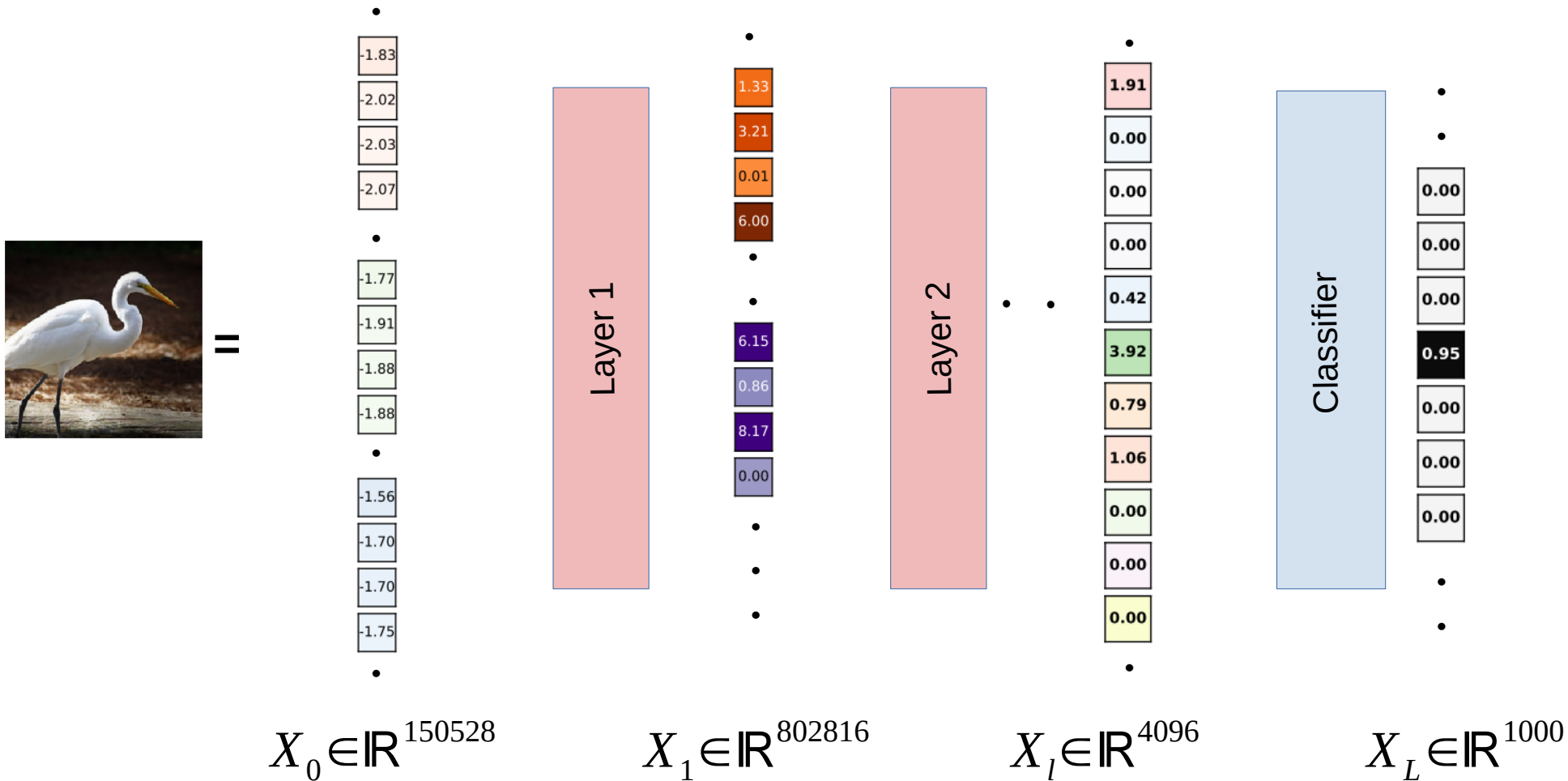# What is a representation in a convolutional neural network?



$$X_0 \in \mathbb{R}^{3 \times 224 \times 224} \qquad X_1 \in \mathbb{R}^{64 \times 112 \times 112} \qquad X_l \in \mathbb{R}^{4096 \times 1 \times 1} \qquad X_L \in \mathbb{R}^{1000}$$
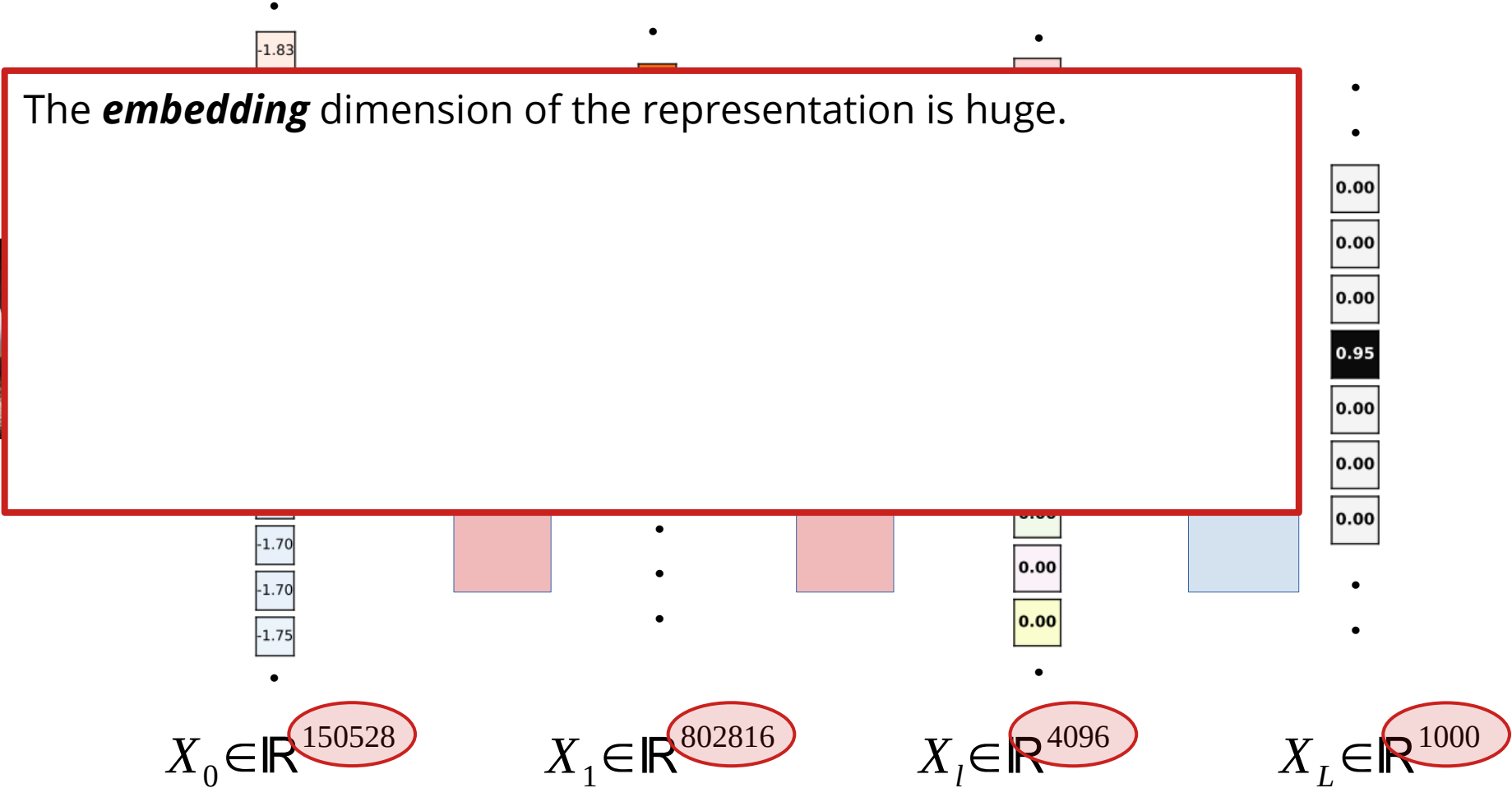
# What is a representation in a convolutional neural network?



$$X_0 \in \mathbb{R}^{150528} \qquad X_1 \in \mathbb{R}^{802816} \qquad X_l \in \mathbb{R}^{4096} \qquad X_L \in \mathbb{R}^{1000}$$

# What is a representation in a convolutional neural network?



The **embedding** dimension of the representation is huge.

$X_0 \in \mathbb{R}^{150528}$

$X_1 \in \mathbb{R}^{802816}$

$X_l \in \mathbb{R}^{4096}$

$X_L \in \mathbb{R}^{1000}$

# What is a representation in a convolutional neural network?


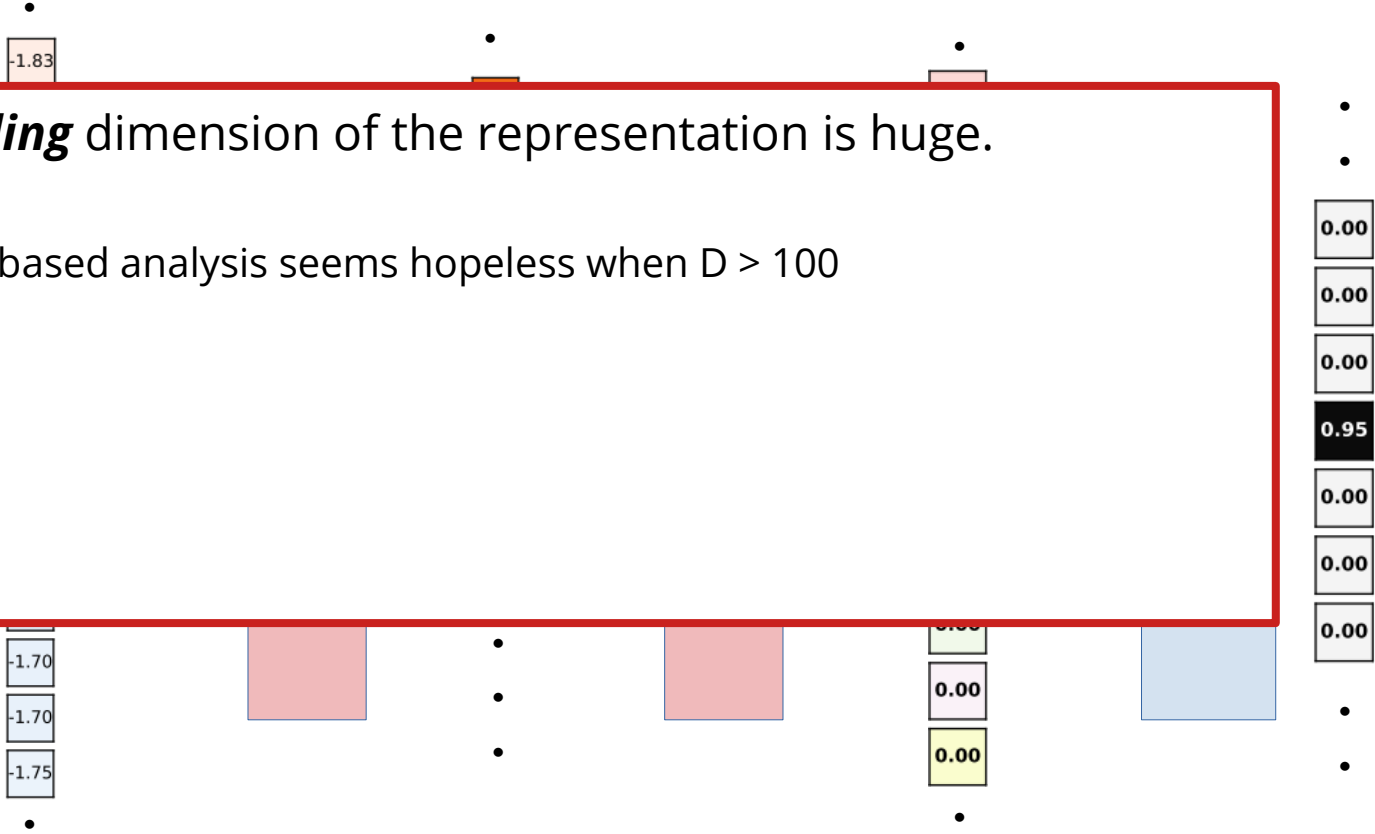
The ***embedding*** dimension of the representation is huge.

a) Any density-based analysis seems hopeless when D > 100
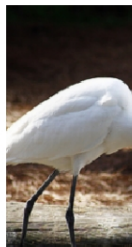
$X_0 \in \mathbb{R}^{150528}$

$X_1 \in \mathbb{R}^{802816}$

$X_l \in \mathbb{R}^{4096}$

$X_L \in \mathbb{R}^{1000}$

# What is a representation in a convolutional neural network?

The **embedding** dimension of the representation is huge.

a) Any density-based analysis seems hopeless when D > 100

b) Neural networks take advantage of the low dimensional structure of the data.
   This is **not** true for other classification approaches (kernels,... )

   Chizat & Bach, *Implicit bias of gradient descent…* Conference on Learning Theory (2020)
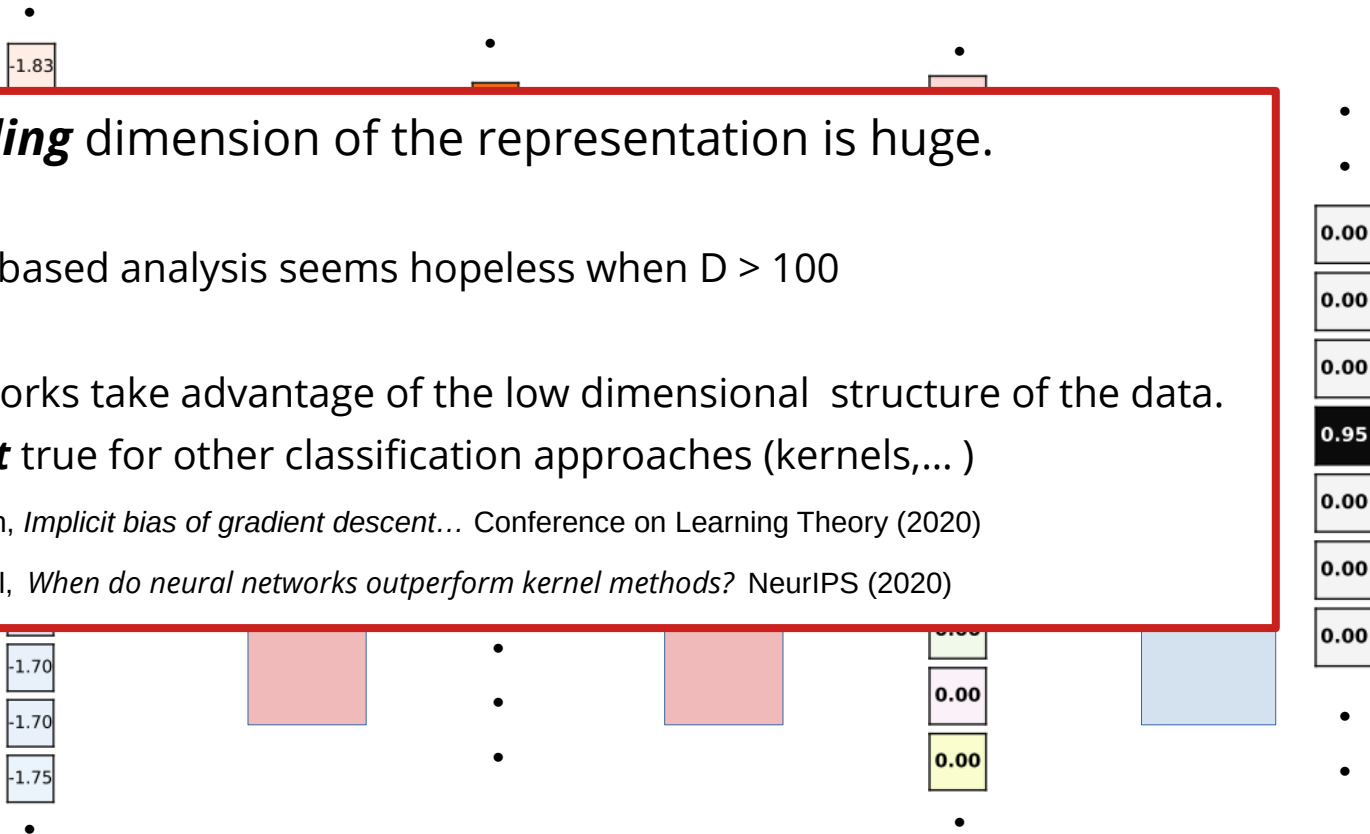
   Ghorbani et al, *When do neural networks outperform kernel methods?* NeurIPS (2020)

$X_0 \in \mathbb{R}^{150528}$        $X_1 \in \mathbb{R}^{802816}$        $X_l \in \mathbb{R}^{4096}$        $X_L \in \mathbb{R}^{1000}$
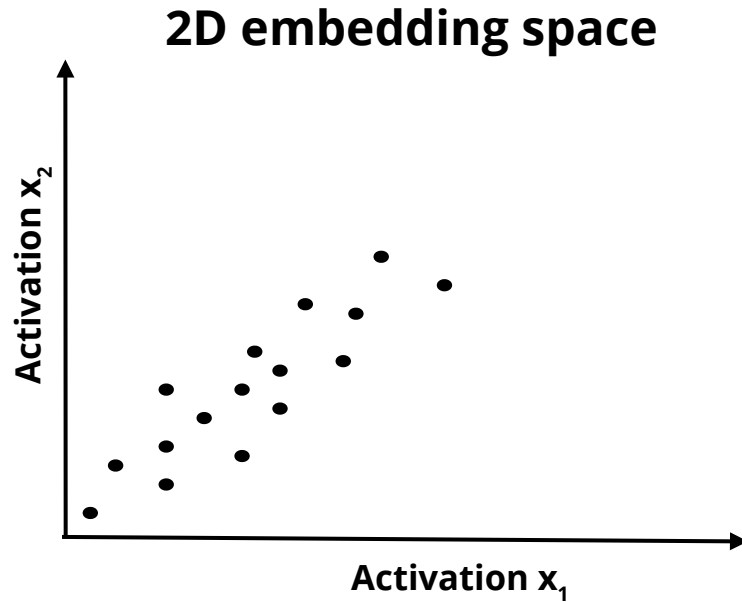
# Intrinsic dimension estimation

**Intrinsic dimension of a data representation**:
minimum number of coordinates to describe the data without significant information loss
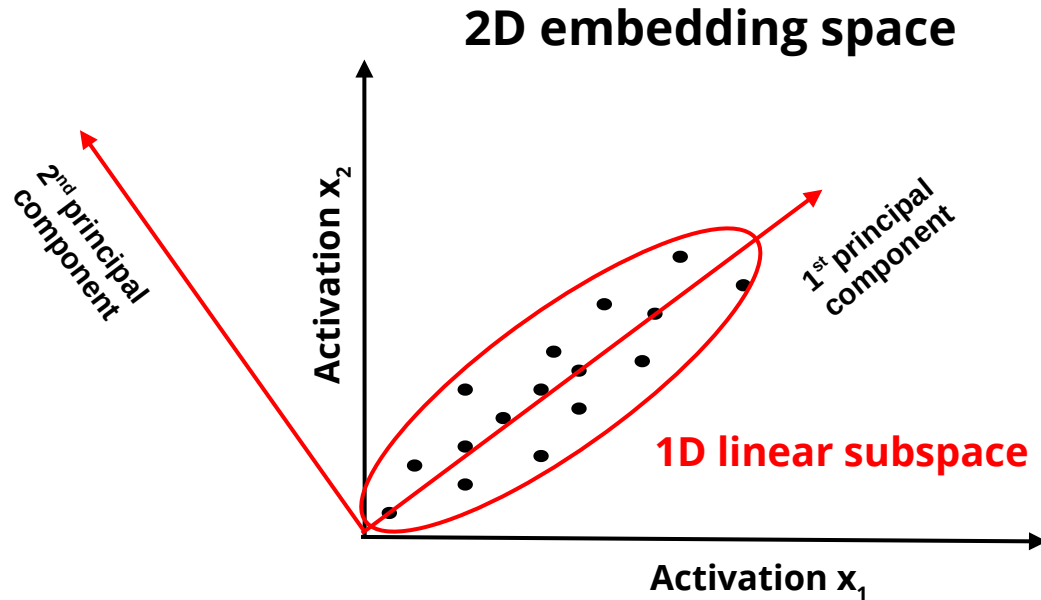
Linear case: Principal Component Analysis (PCA)



**2D embedding space**

# Intrinsic dimension estimation

**Intrinsic dimension of a data representation**:
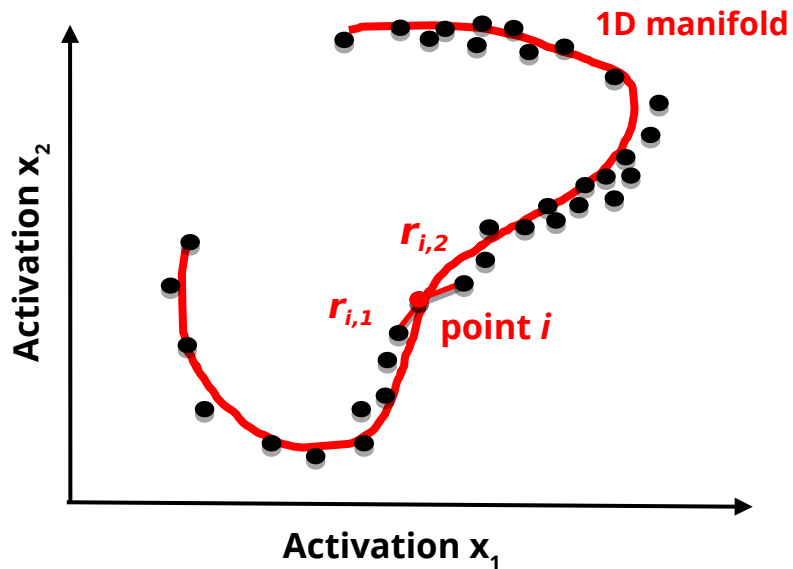minimum number of coordinates to describe the data without significant information loss

Linear case: Principal Component Analysis (PCA)

# Intrinsic dimension estimation

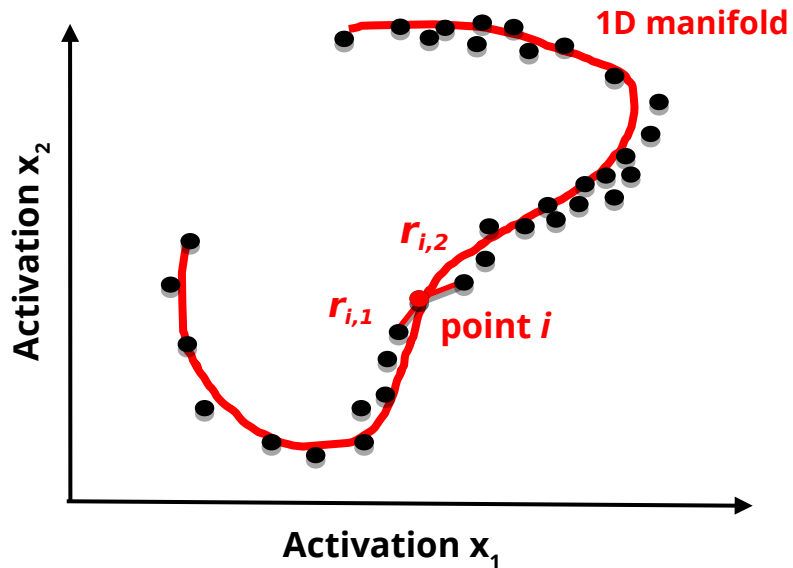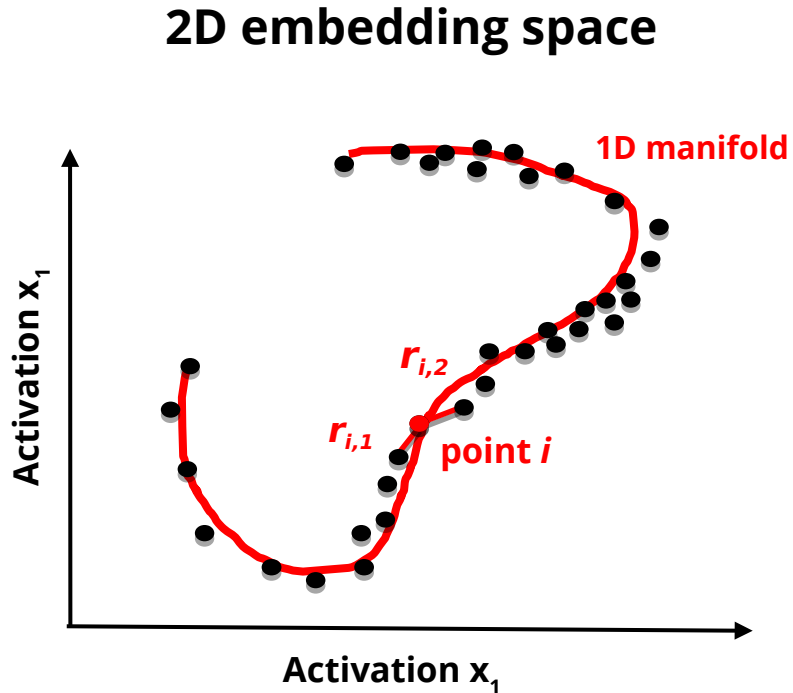The general (non linear) case: TwoNN (Facco et al, 2017)

# Intrinsic dimension estimation

The general (non linear) case: TwoNN (Facco et al, 2017)

1) For each data point $i$ compute the distance to its first and second neighbors ($r_{i,1}$ and $r_{i,2}$)



**2D embedding space**

# Intrinsic dimension estimation

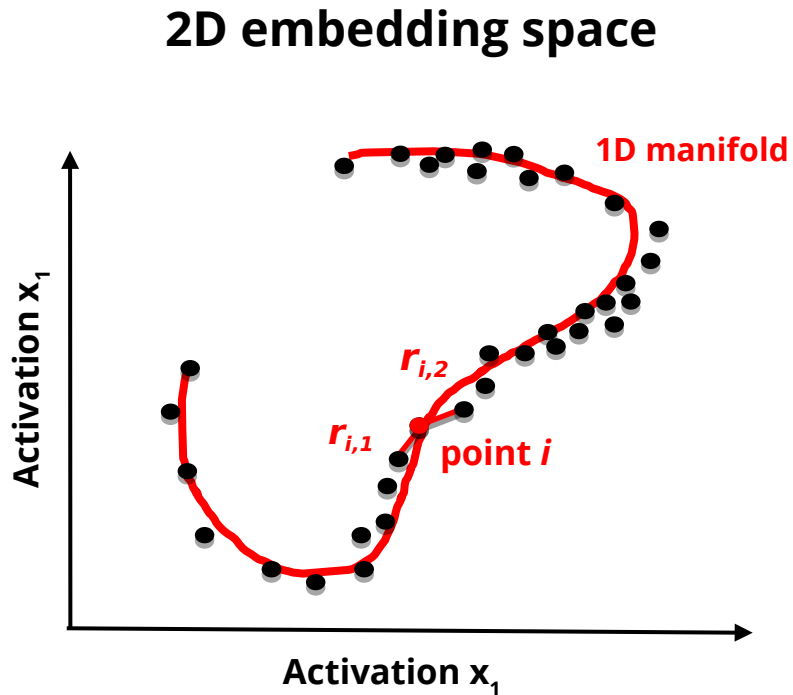The general (non linear) case: TwoNN (Facco et al, 2017)

**2D embedding space**



1D manifold

$r_{i,2}$

$r_{i,1}$

point $i$

Activation $x_1$

Activation $x_1$

1) For each data point $i$ compute the distance to its first and second neighbors ($r_{i,1}$ and $r_{i,2}$)

2) For each $i$ compute $\quad \mu_i = \dfrac{r_{i,2}}{r_{i,1}}$

# Intrinsic dimension estimation

The general (non linear) case: TwoNN (Facco et al, 2017)



**2D embedding space**

1) For each data point $i$ compute the distance to its first and second neighbors ($r_{i,1}$ and $r_{i,2}$)

2) For each $i$ compute $\quad \mu_i = \dfrac{r_{i,2}}{r_{i,1}}$
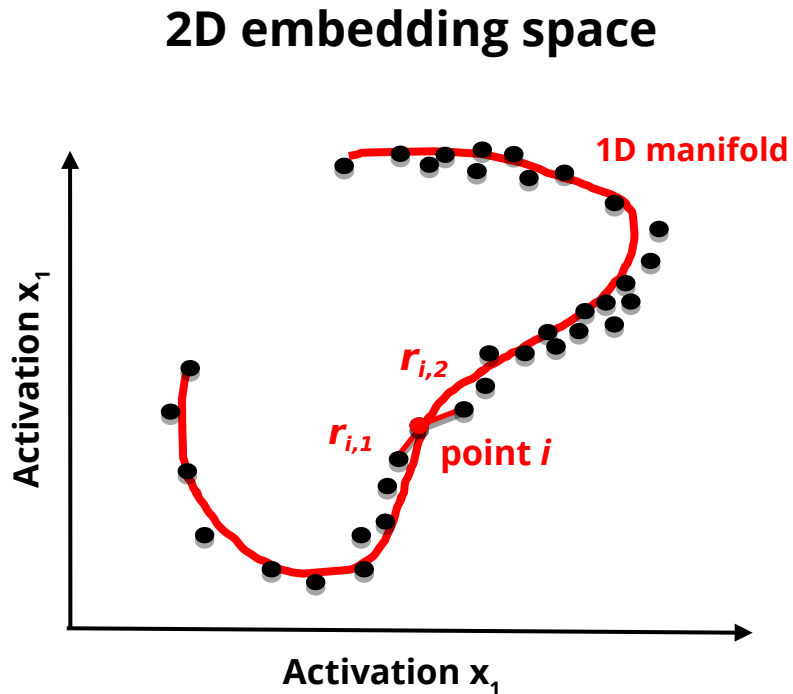
The probability distribution of $\mu$ is

$$p(\mu|d) = \frac{d}{\mu^{d+1}}$$

where $d$ is the ID.

# Intrinsic dimension estimation

The general (non linear) case: TwoNN (Facco et al, 2017)



**2D embedding space**

1D manifold

$r_{i,2}$

$r_{i,1}$  point $i$

Activation $x_1$

Activation $x_1$

1) For each data point $i$ compute the distance to its first and second neighbors ($r_{i,1}$ and $r_{i,2}$)

2) For each $i$ compute $\quad \mu_i = \dfrac{r_{i,2}}{r_{i,1}}$

The probability distribution of $\mu$ is

$$p(\mu|d) = \frac{d}{\mu^{d+1}}$$

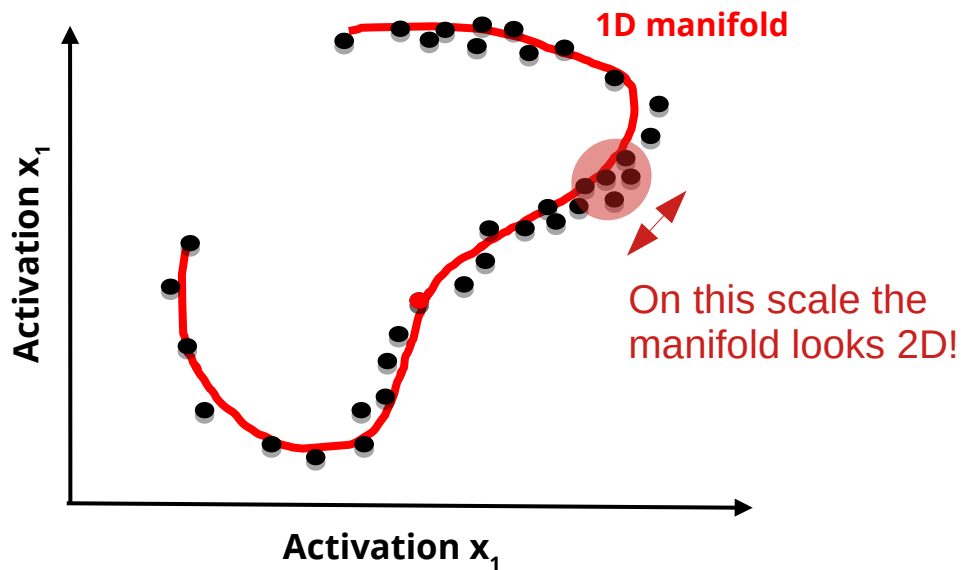where $d$ is the ID.

3) Infer $d$ e.g via maximum likelihood

$$L(\mu_i|d) = \log \prod^{N} p(\mu_i|d)$$

$$\partial_d L(\mu_i|d) = 0 \rightarrow \hat{d} = \frac{N}{\sum \log \mu_i}$$

# Intrinsic dimension estimation of a noisy manifold

When the data are noisy TwoNN can overestimate the ID due to its **local** nature

**2D embedding space**


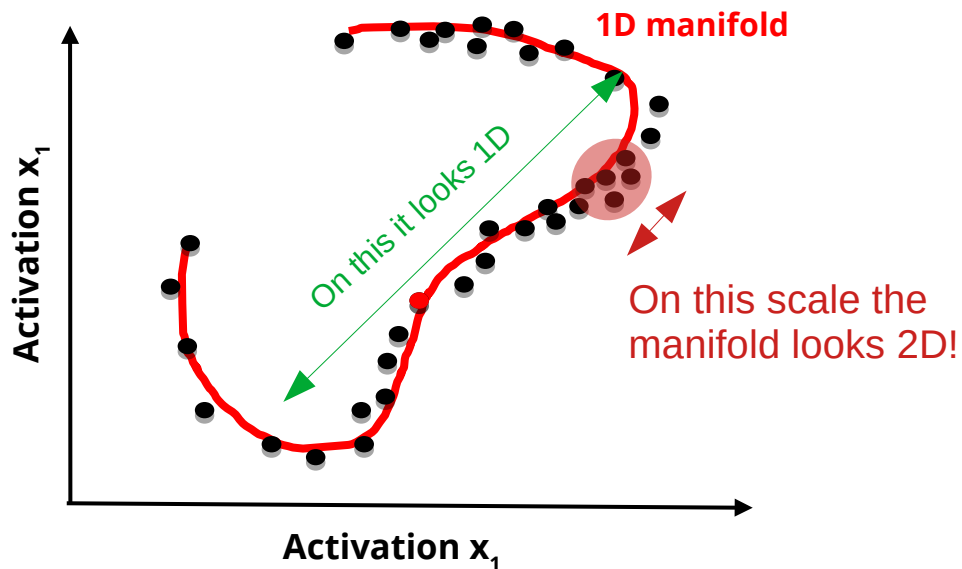
1D manifold

On this scale the manifold looks 2D!

# Intrinsic dimension estimation of a noisy manifold

When the data are noisy TwoNN can overestimate the ID due to its **local** nature

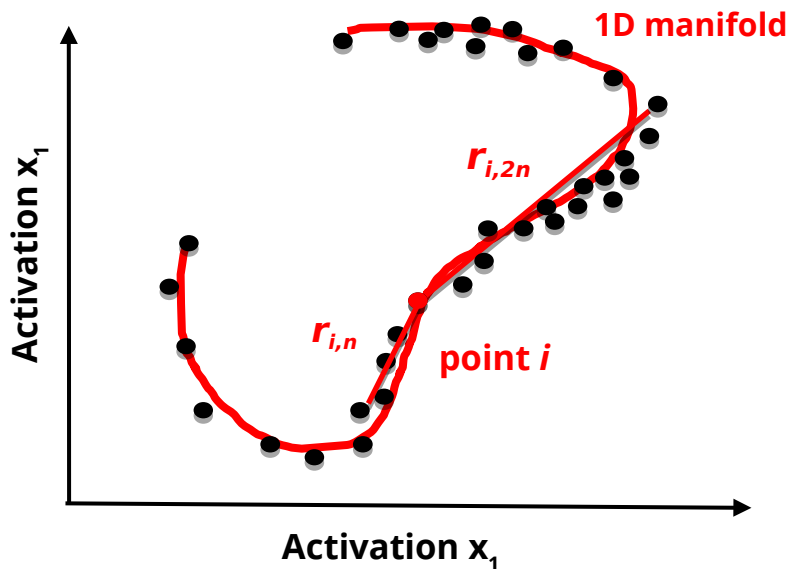Enlarge the neighborhood range to find the actual 'soft directions' of the data

**2D embedding space**

# Intrinsic dimension estimation of a noisy manifold

When the data are noisy TwoNN can overestimate the ID due to its **local** nature
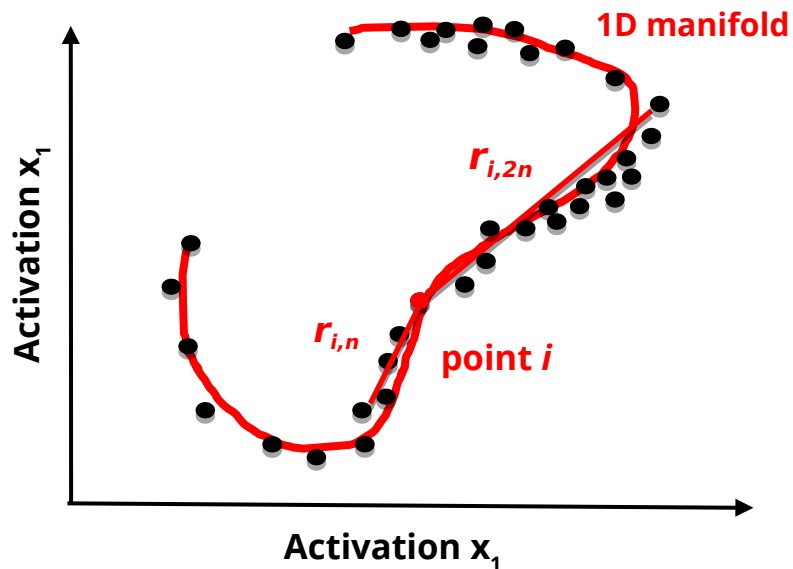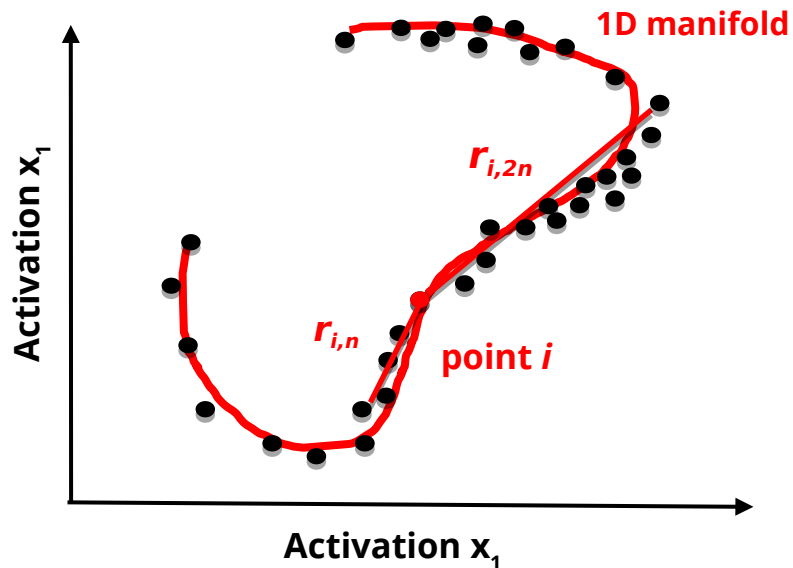
**2D embedding space**



**1D manifold**

$r_{i,2n}$

$r_{i,n}$

**point $i$**

Activation $x_1$

Activation $x_1$

1) For each data point $i$ compute the distance to its **nth** and **2*nth** neighbors ($r_{i,n}$ and $r_{i,2n}$)

# Intrinsic dimension estimation of a noisy manifold

When the data are noisy TwoNN can overestimate the ID due to its **local** nature

**2D embedding space**



1) For each data point $i$ compute the distance to its **nth** and **2*nth** neighbors ($r_{i,n}$ and $r_{i,2n}$)

2) For each $i$ compute $\quad \mu_{n,i} = \dfrac{r_{i,2n}}{r_{i,n}}$

# Intrinsic dimension estimation of a noisy manifold

When the data are noisy TwoNN can overestimate the ID due to its **local** nature

## 2D embedding space



1) For each data point $i$ compute the distance to its **nth** and **2*nth** neighbors ($r_{i,n}$ and $r_{i,2n}$)

2) For each $i$ compute $\quad \mu_{n,i} = \dfrac{r_{i,2n}}{r_{i,n}}$
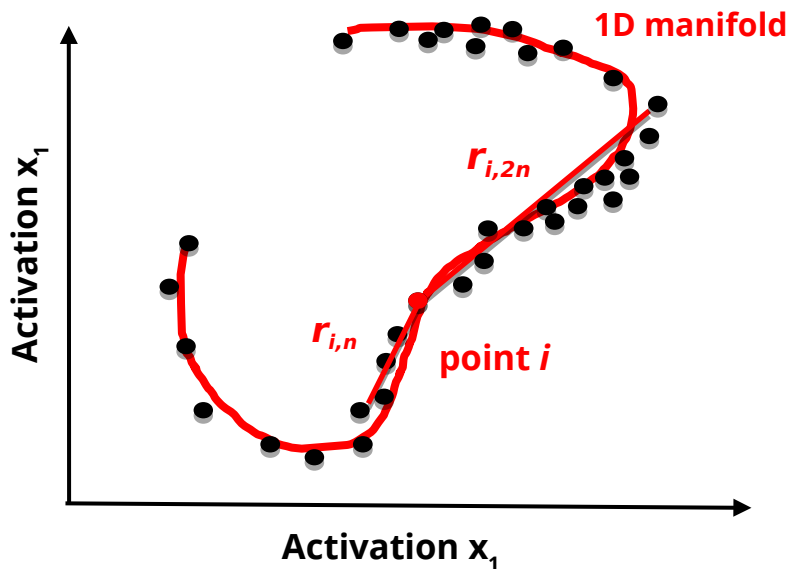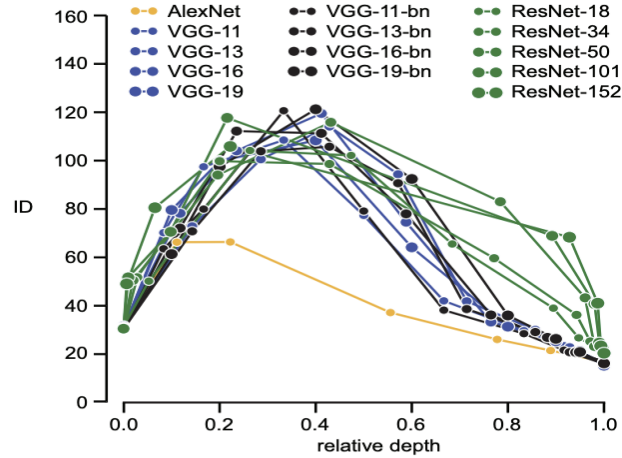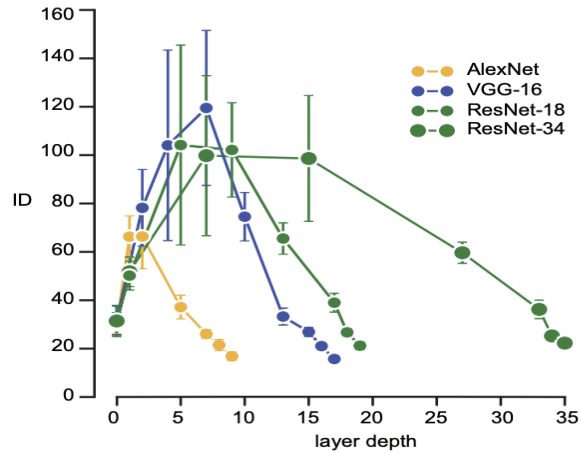
The probability distribution of $\mu_n$ is

$$p(\mu_n | d) \propto \frac{d(\mu^d - 1)^{n-1}}{\mu^{(2n-1)d+1}}$$
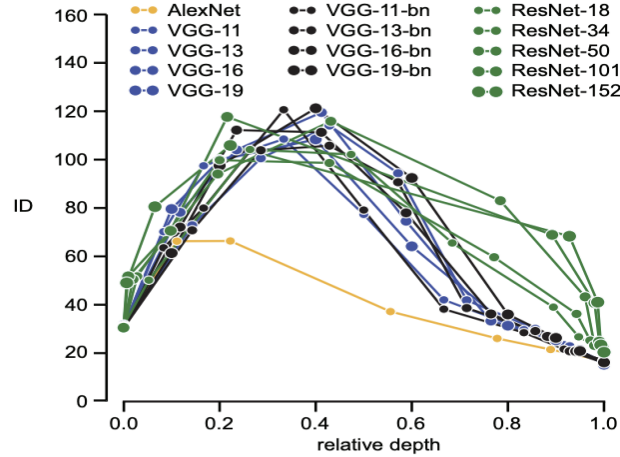
where $d$ is the ID.

# Intrinsic dimension estimation of a noisy manifold

When the data are noisy TwoNN can overestimate the ID due to its **local** nature

## 2D embedding space



Activation $x_1$

Activation $x_1$

1D manifold

$r_{i,2n}$

$r_{i,n}$

point $i$

1) For each data point $i$ compute the distance to its **nth** and **2*nth** neighbors ($r_{i,n}$ and $r_{i,2n}$)

2) For each $i$ compute $\quad \mu_{n,i} = \dfrac{r_{i,2n}}{r_{i,n}}$

The probability distribution of $\mu_n$ is

$$p(\mu_n|d) \propto \frac{d(\mu^d - 1)^{n-1}}{\mu^{(2n-1)d+1}}$$
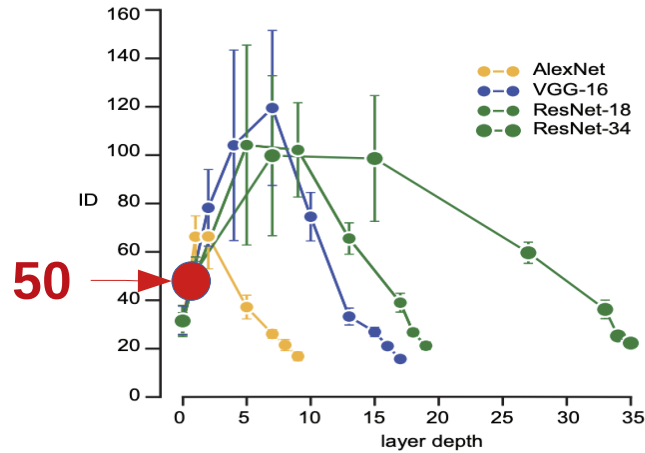
where $d$ is the ID.

3) Infer $d$ via maximum likelihood
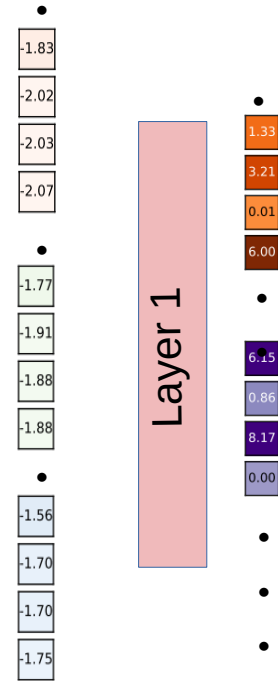
# Expansion and compression of the ID



The ID is always much smaller than the embedding dimension
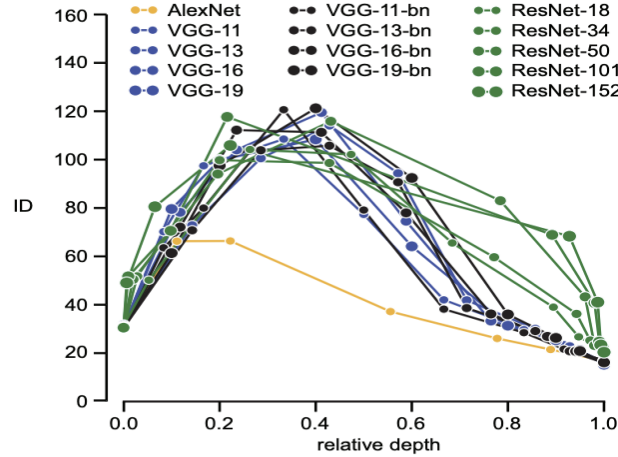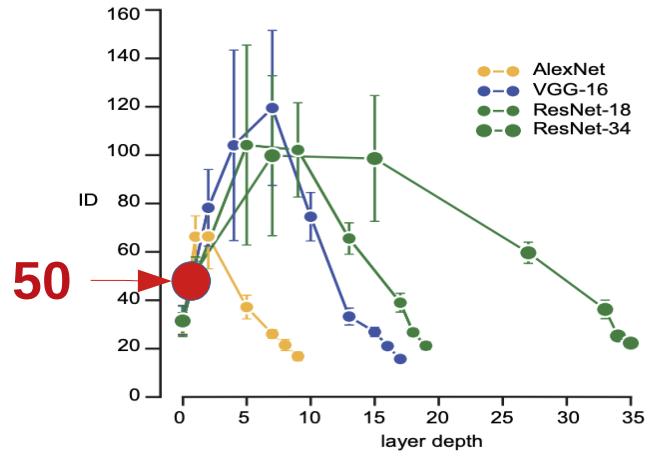
# Expansion and compression of the ID



The ID is always much smaller than the embedding dimension



$$X_1 \in \mathbb{R}^{800\,000}$$

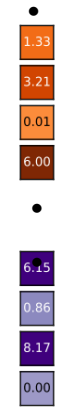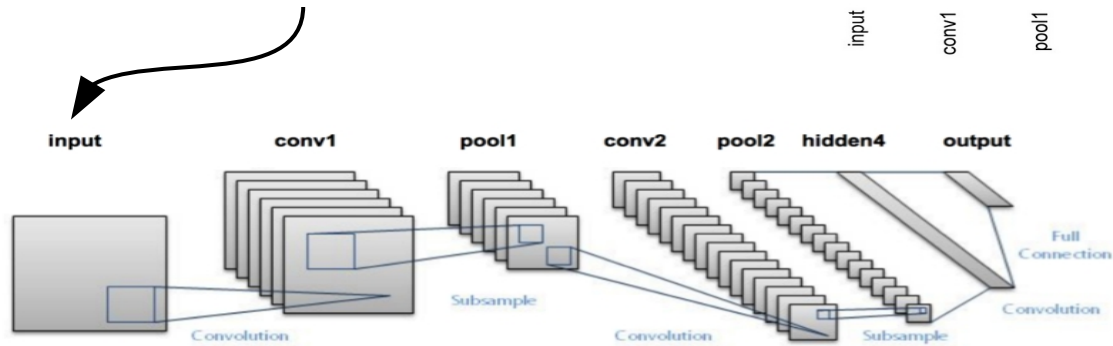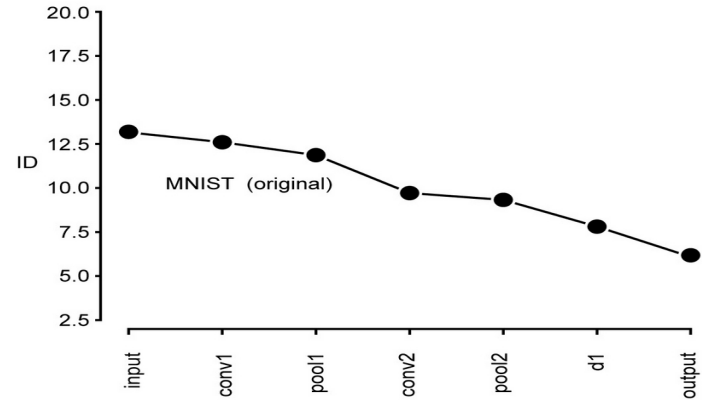# Expansion and compression of the ID



The ID is always much smaller than the embedding dimension

ID evolution across layer has a hunchback shape

# Discarding useless features

**MNIST (original)**
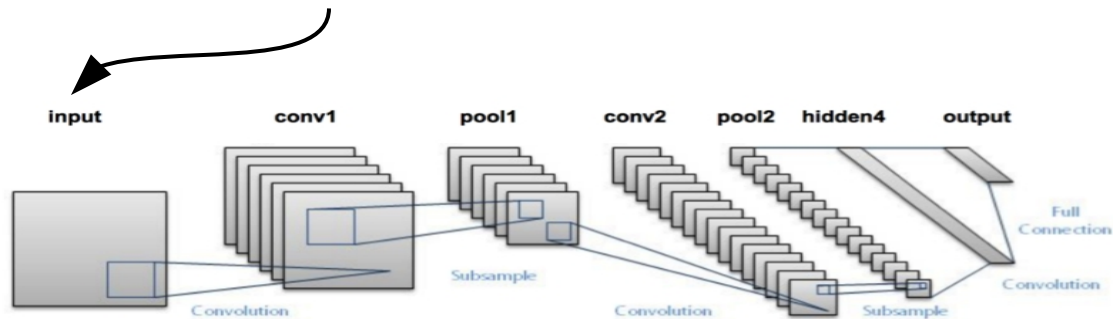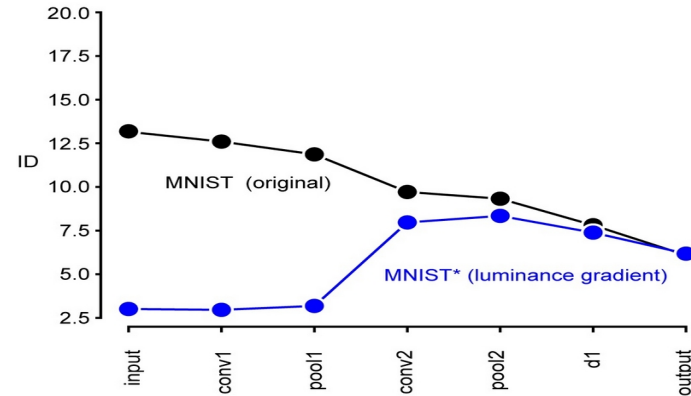
# Discarding useless features



**MNIST\* (luminance gradient)**

ORIGINAL MNIST DATA

MNIST DATA PERTURBED WITH A LUMINANCE GRADIENT (MNIST\*)

average image pixel value (for MNIST\*)

MNIST (original)

MNIST\* (luminance gradient)

input    conv1    pool1    conv2    pool2    d1    output

input    conv1    pool1    conv2    pool2    hidden4    output

Convolution    Subsample    Convolution    Subsample    Full Connection    Convolution

In a trained network, the initial ID expansion reflects the pruning of low-level visual features that carry no information about the correct labeling