

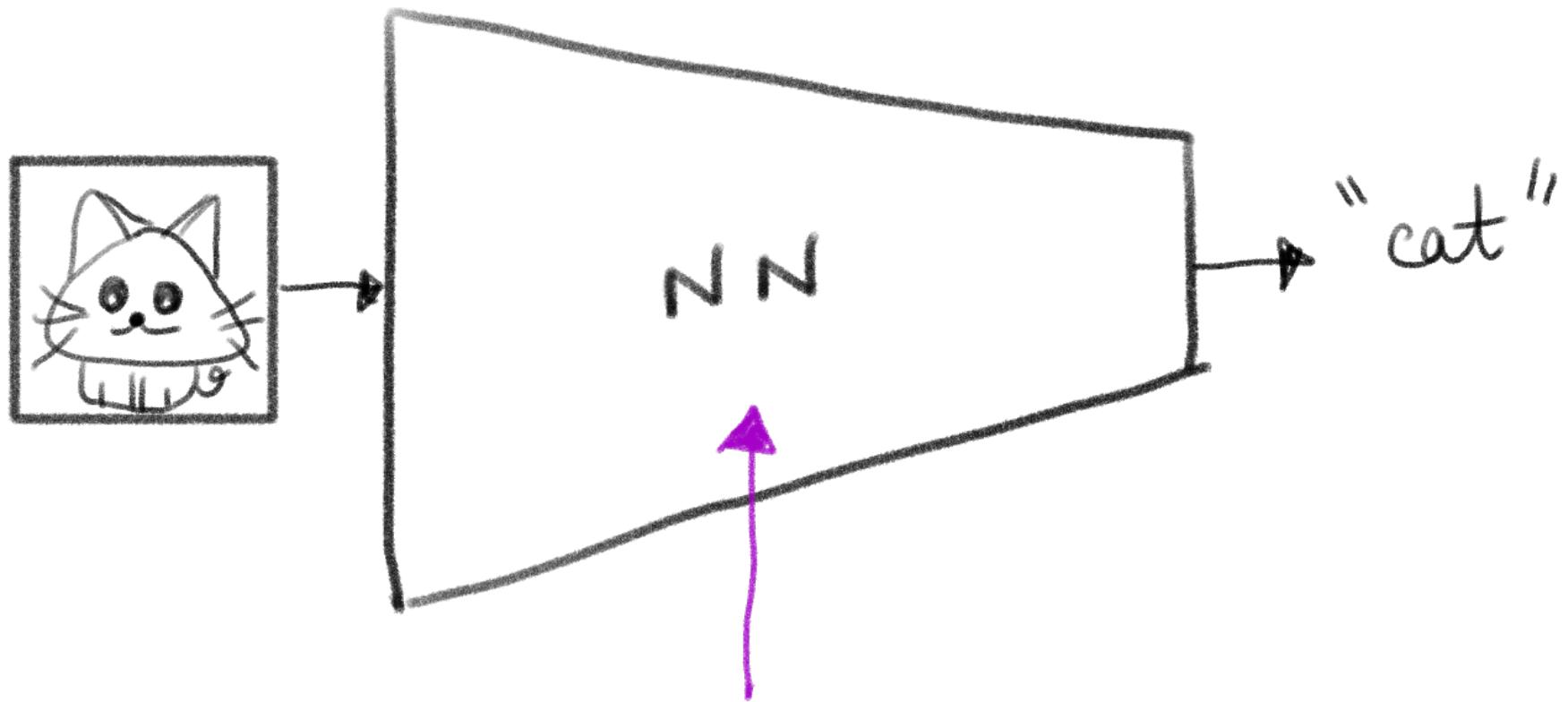
HDDL - LESSON 3

AUTOENCODERS

SELF - SUPERVISED LEARNING

JOSEBA DALMAU

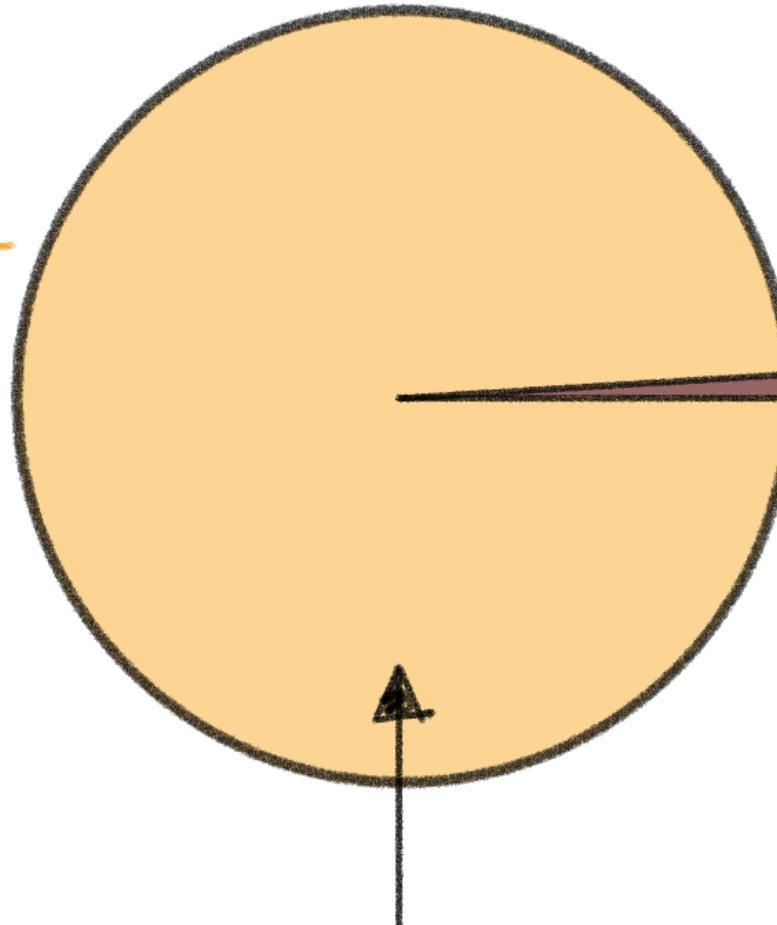
SUPERVISED LEARNING



What do we need in order to
train this neural network?

DATA

unlabeled



labeled

how do we leverage all this data ?

WHAT IS AN AUTOENCODER ?

Autoencoder

An unsupervised learning technique
for the task of representation learning.

Unsupervised learning

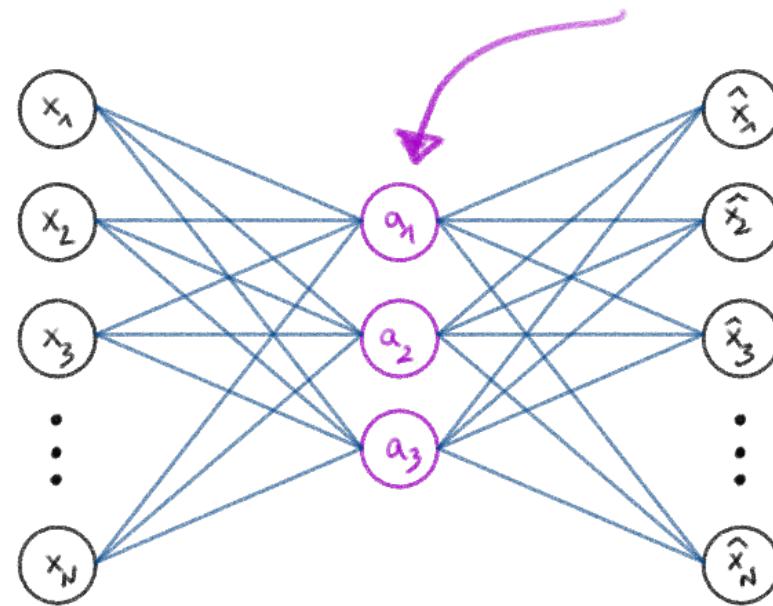
- no labels,
only input!

Representation learning

- Latent representation
(lower dimensions)
- Transferable
representations

MAIN IDEA

Architecture: NN with a bottleneck



Loss: reconstruction error

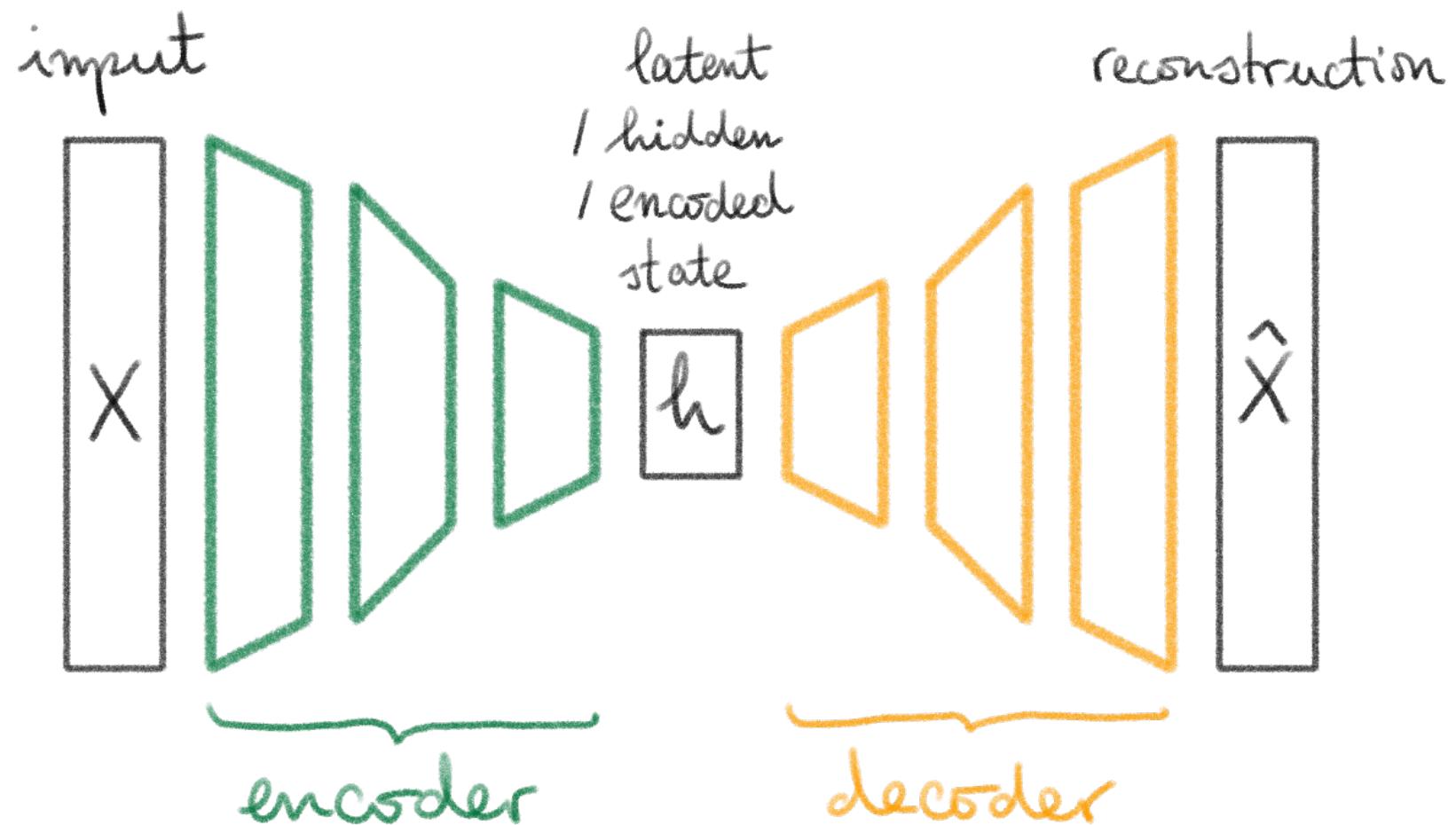
$$L(x, \hat{x})$$

QUESTIONS AND REMARKS

- Why is the bottleneck important?
- This architecture is:
 - unsupervised
 - provides latent representations
 - are the latent representations transferable?
- Remark: linear autoencoder = PCA
- What is the risk with a reconstruction error
 $L(x, \hat{x}) \approx 0$?

NAMES AND ABBREVIATIONS

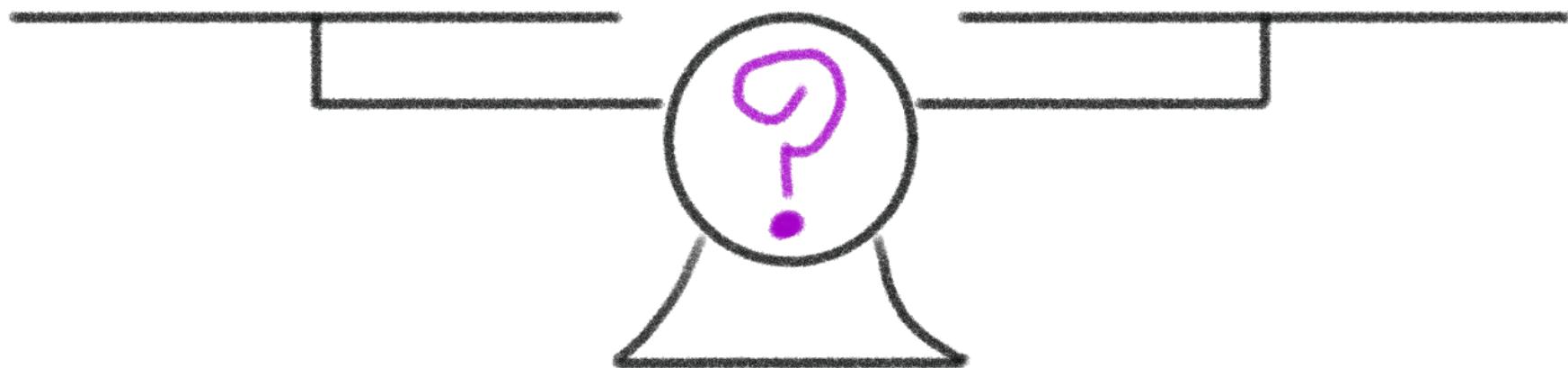
- Autoencoder \rightarrow AE



THE IDEAL AUTOENCODER

Sensitive enough
to the inputs
⇒ reconstruct accurately

Insensitive enough
to the inputs
⇒ avoid memorizing



$$\text{Loss} = L(x, \hat{x}) + \lambda \text{Regularizer}$$

UNDERCOMPLETE AUTOENCODER

IDEA : Rely on the bottleneck

If $\dim(\text{latent representation}) \ll \dim(\text{input})$

we
expect $\Rightarrow \begin{cases} - \text{Learn the most important attributes} \\ - \text{Be able to reconstruct input} \end{cases}$

Loss = $L(x, \hat{x})$ (no Regularizer)



WAIT A MINUTE...

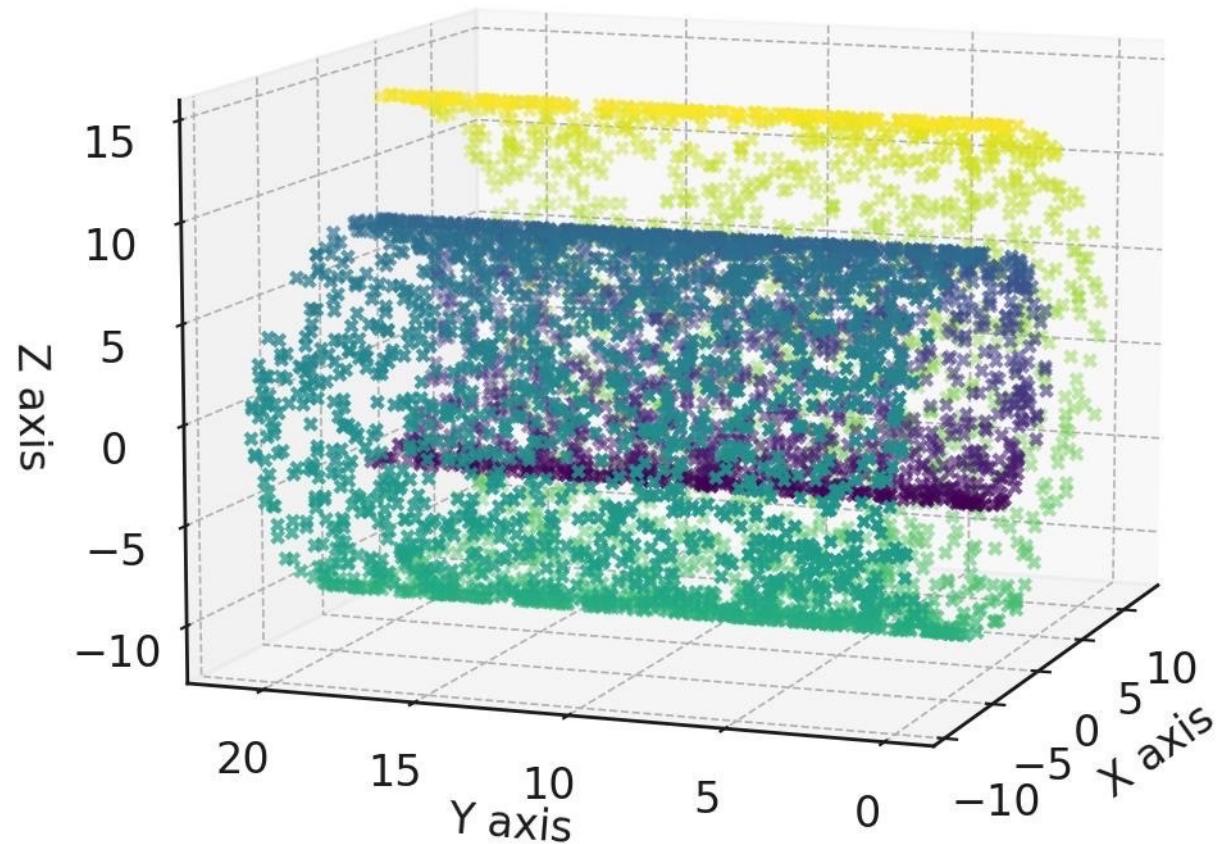
... how is this different from PCA ?

NNs can learn **non-linear** relationships

$\Rightarrow AE \equiv$ more powerful (nonlinear)
generalization of PCA

PCA
↓ learn
lower dimensional
hyperplane

AE
↓ learn
lower dimensional
non-linear manifold





PROBLEM: "CAPACITY" OF THE AUTOENCODER

If the encoder and decoder are complex enough ...

... even with a very small number of latent dimensions ...

... input data can be memorized !

SPARSE AUTOENCODERS

IDEA

~~Reduction in nb. of nodes of
hidden layers~~



Penalize activations
within a layer

REMARKS

- { Usually: Regularizer $\xrightarrow{\text{penalize}}$ weights
- (i) e.g. $\sum_{ij} |w_{ij}|$
- { Sparse AEs: Regularizer $\xrightarrow{\text{penalize}}$ activations
- (ii) The neurons that activate are data-dependent
- (iii) latent dimension indep. of regularization

REGULARIZATION

Given: sparsity parameter ρ
(average activation of a neuron over a collection of samples)

Define:

$$\hat{\rho}_j^{(h)} = \frac{1}{m} \sum_{i=1}^m a_j^{(h)}(x^{(i)})$$

activation on layer h
neuron j
i-th training obs.
nb. of training obs.

layer
↑
neuron

REGULARIZATION

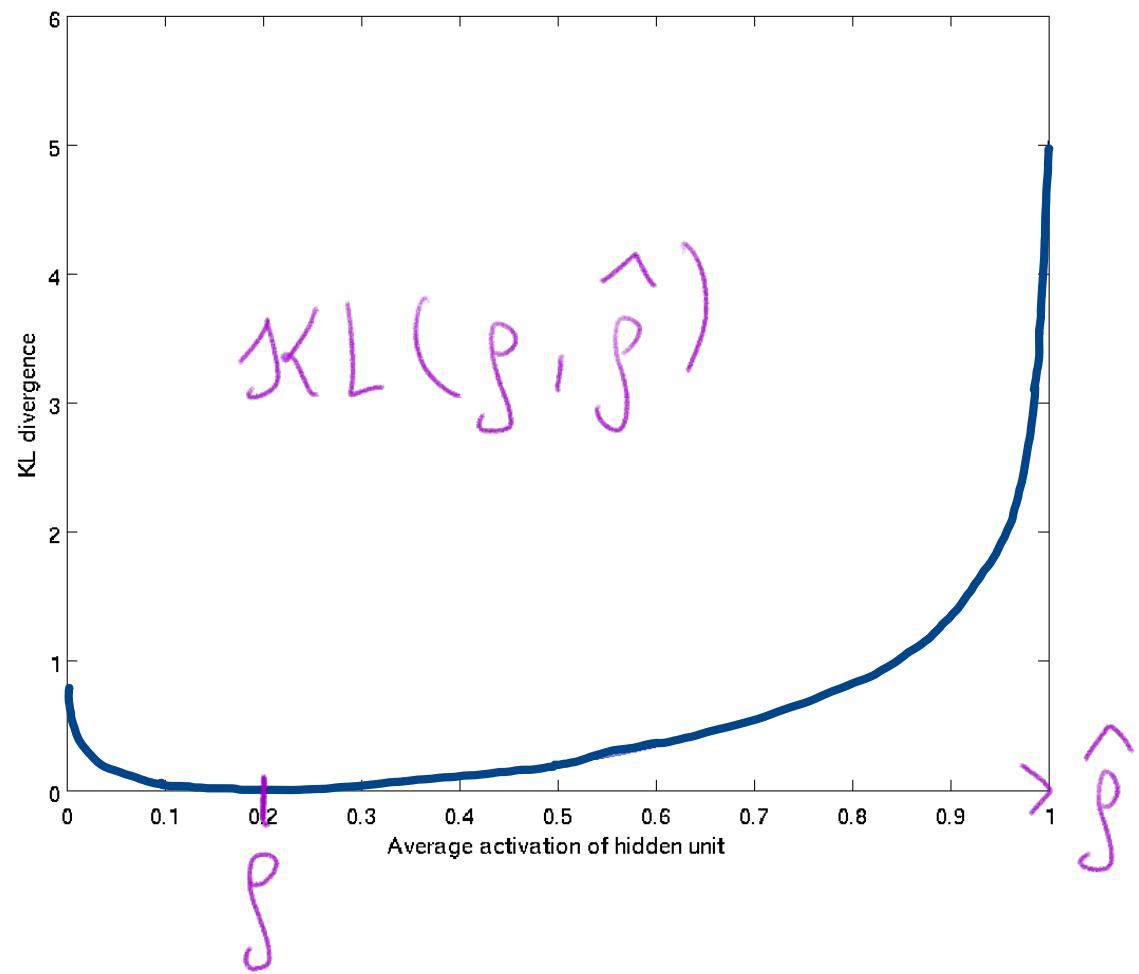
Given: sparsity parameter ρ
(average activation of a neuron over a collection of samples)

Define:

$$\hat{p}_j^{(h)} = \frac{1}{m} \sum_{i=1}^m a_j^{(h)}(x^{(i)})$$

$$\text{Loss} = L(x, \hat{x}) + \lambda \sum_{h,j} \text{KL}(\rho, \hat{p}_j^{(h)})$$

KL-DIVERGENCE



$KL(p, \hat{p})$ has a minimum at $\hat{p} = p$ and is NOT symmetrical!

BACKPROPAGATION

Given that:

$$\hat{p}_j^{(h)} = \frac{1}{m} \sum_{i=1}^m a_j^{(h)}(x^{(i)})$$

$$\text{Loss} = L(x, \hat{x}) + \lambda \sum_{h,j} KL(p, \hat{p}_j)$$

How to perform backpropagation?

BACKPROPAGATION

Given that:

$$\hat{p}_j^{(h)} = \frac{1}{m} \sum_{i=1}^m a_j^{(h)}(x^{(i)})$$

$$\text{Loss} = L(x, \hat{x}) + \lambda \sum_{h,j} KL(p, \hat{p}_j)$$

How to perform backpropagation?

Sol: perform multiple forward passes per backward pass!

EXERCISE

Consider an AE with an input dim of n .

Call the weights of the first (FC) layer $W^{(1)}$ so that activations after the 1st layer are

$$a_i^{(1)}(x) = \sum_{1 \leq j \leq n} w_{ij}^{(1)} x_j + b_i^{(1)}, \text{ for } i=1, \dots, n_1$$

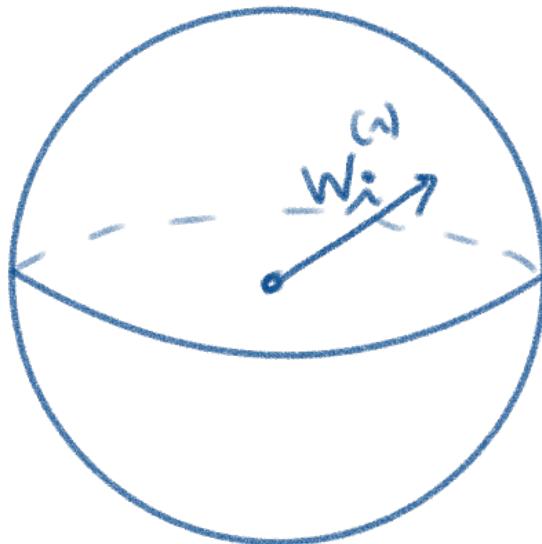
Assuming that the input satisfies

$\|x\|_2^2 \leq 1$. Find the input x that maximizes $a_i^{(1)}(x)$.

SOLUTION

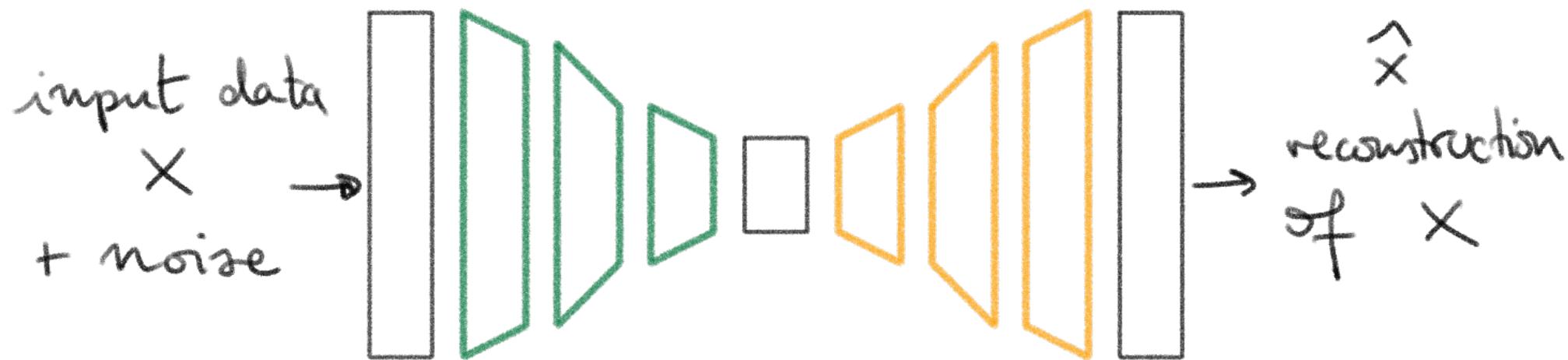
$$a_i(x) = \sum_{j=1}^n w_{ij}^{(1)} x_j + b_i = \langle w_i^{(1)}, x \rangle + b_i$$

is maximal when $x = \frac{w_i^{(1)}}{\|w_i^{(1)}\|_2}$



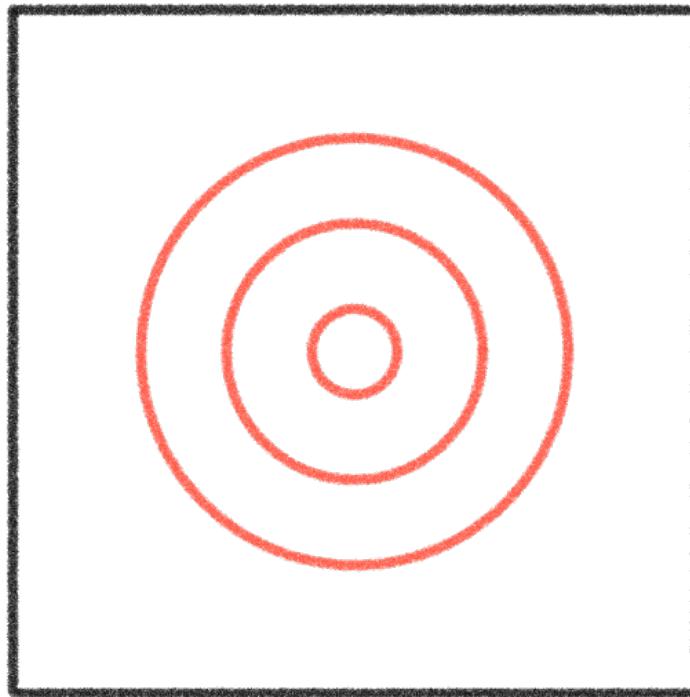
DENOISING AUTOENCODERS

IDEA

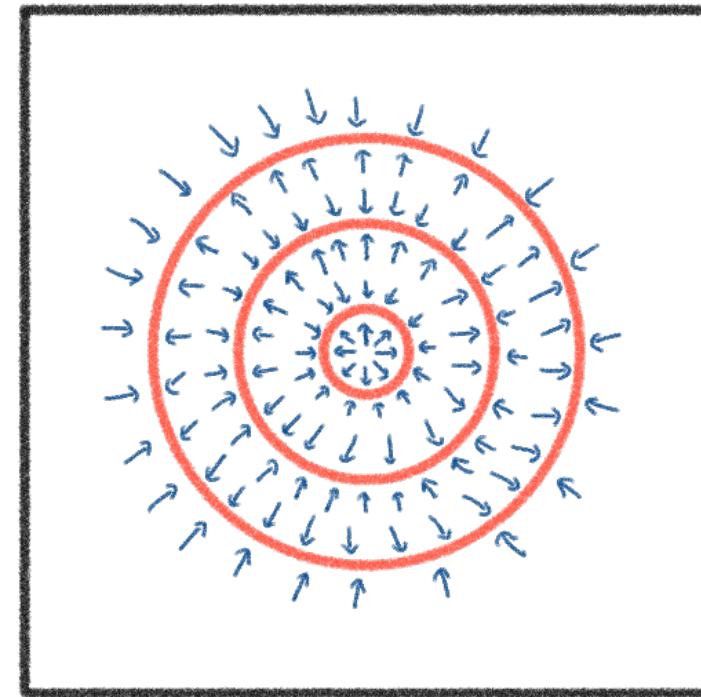


\Rightarrow Prevent AE from memorizing
(why?)

LEARNED VECTOR FIELD



original data



learned vector field

What vector field do we learn in regions where no training data is available?

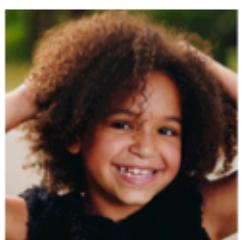
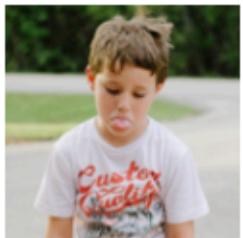
VARIATIONAL AUTOENCODERS

AE: input $\xrightarrow{\text{encoder}}$ representation
(single value)

VAE: input $\xrightarrow{\text{encoder}}$ representation
(probability distribution)

\Rightarrow VAEs can easily generate new,
unseen data.

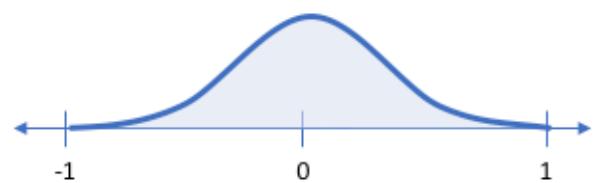
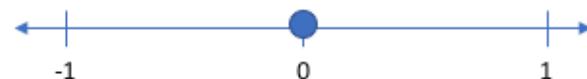
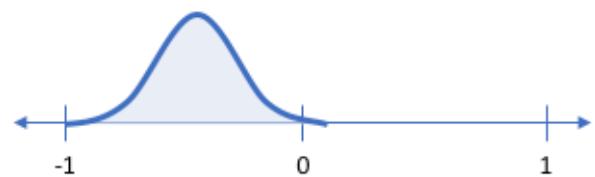
EXAMPLE



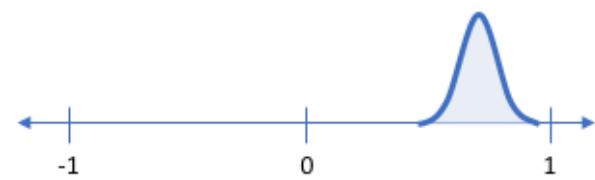
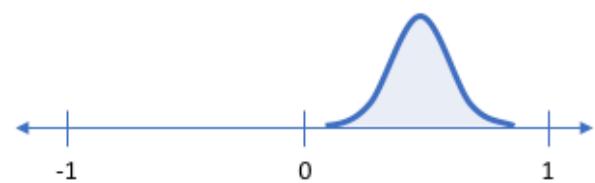
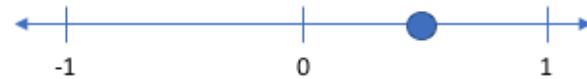
Smile (discrete value)



Smile (probability distribution)

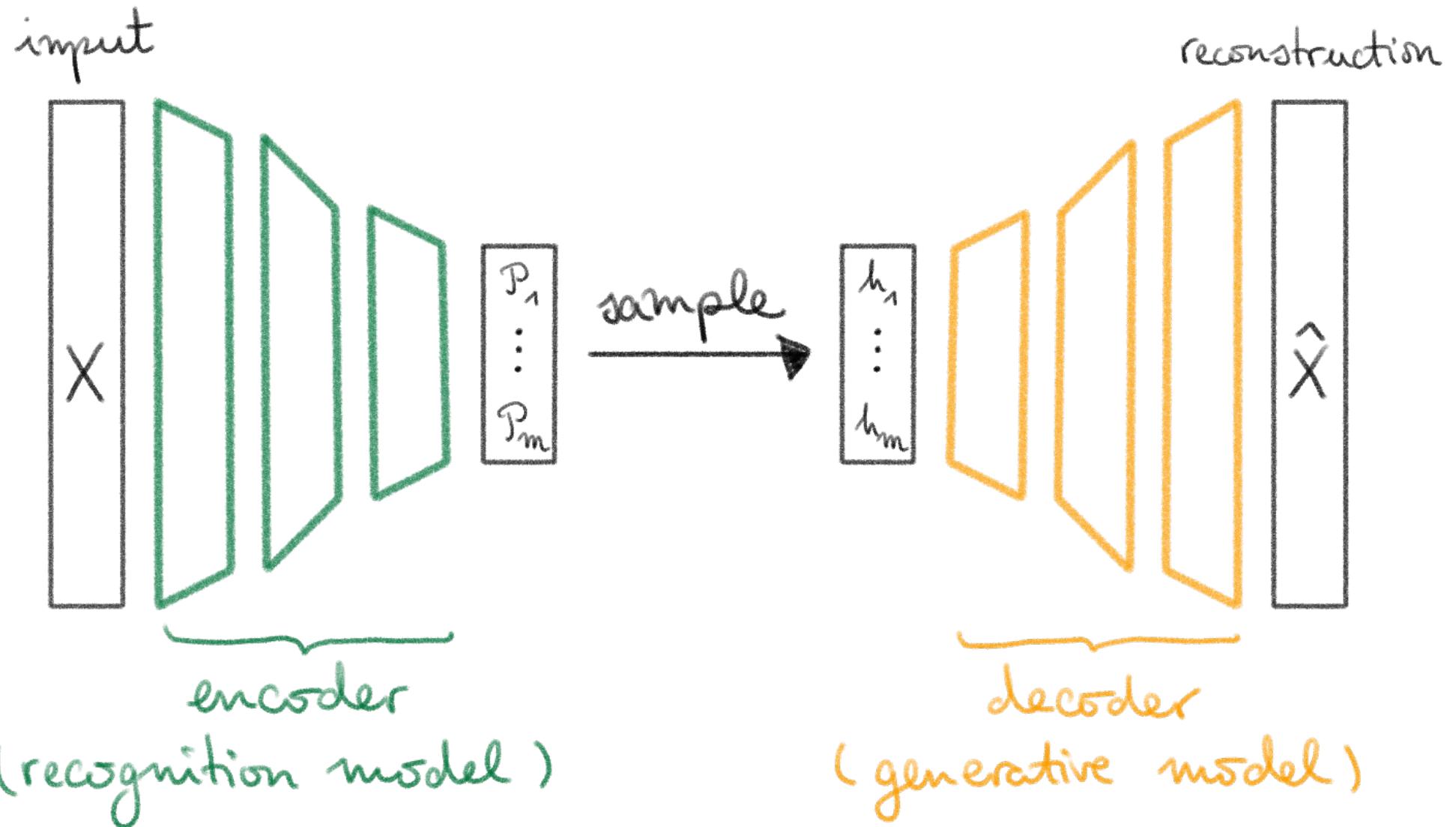


vs.



from Jeremy Jordan's blog

VAE ARCHITECTURE

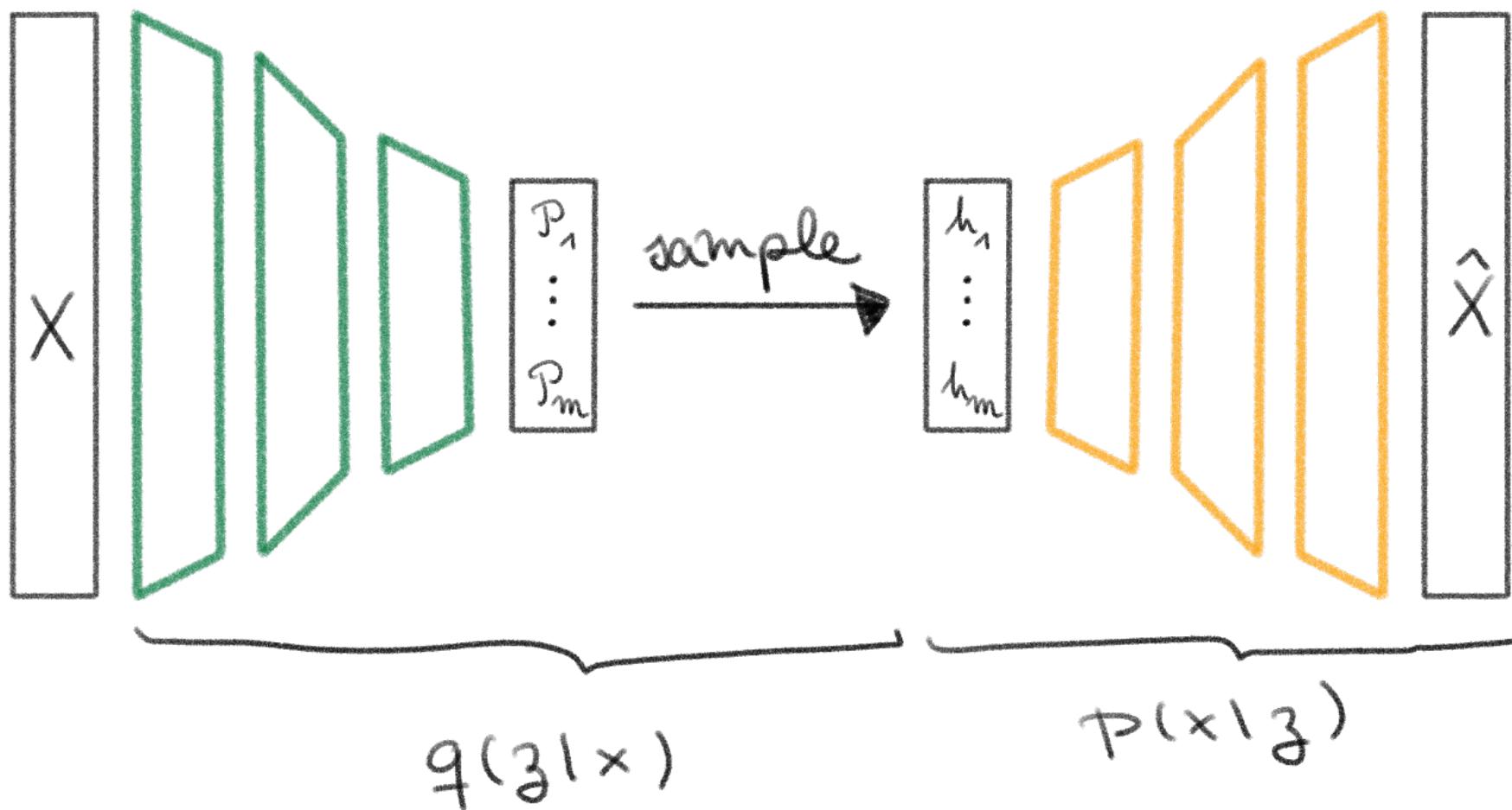


REMARKS

- (i) We are enforcing a continuous, smooth latent representation.
- (ii) If the decoder reconstructs well for any sampling, necessarily:

close in latent space \Rightarrow similar reconstructions

STATISTICAL MOTIVATION



LOSS FUNCTION

$$\text{Loss} = L(x, \hat{x}) + \sum_{j=1}^m \text{KL}(q(z_j|x) || p(z_j))$$

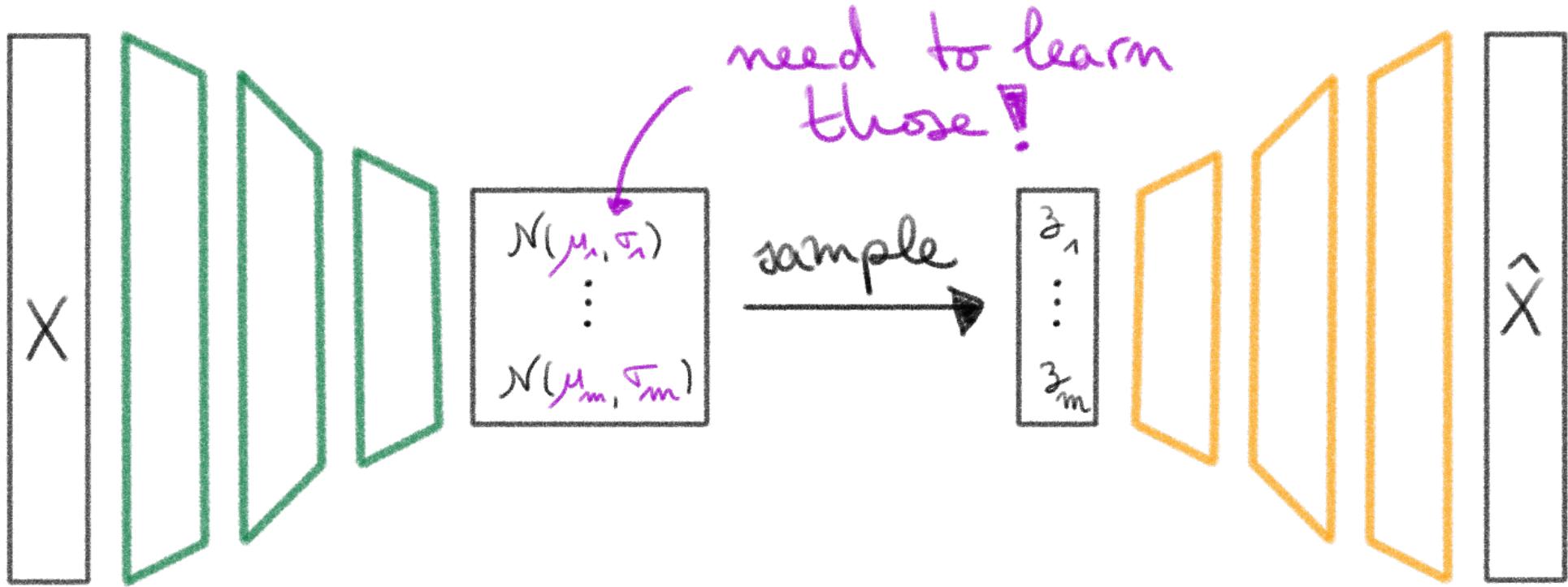
↑ ↑
 MSE or BCE true prior
 usually assumed to be
 $N(0, 1)$

IMPLEMENTATION

We have to learn the mean vector and covariance matrix of the latent Gaussian distribution.

Simplifying Assumption .- The covariance matrix is diagonal, i.e. the latent features are uncorrelated.

BACKPROPAGATION



Since $z_i = N(\mu_i, \Sigma_i) \dots$

$$\dots \frac{\partial L(x, \hat{x})}{\partial \mu_i} = ?$$

REPARAMETERIZATION TRICK

~~$z \sim N(\mu, \sigma)$~~ instead $z = \mu + \sigma \varepsilon, \varepsilon \sim N(0, 1)$

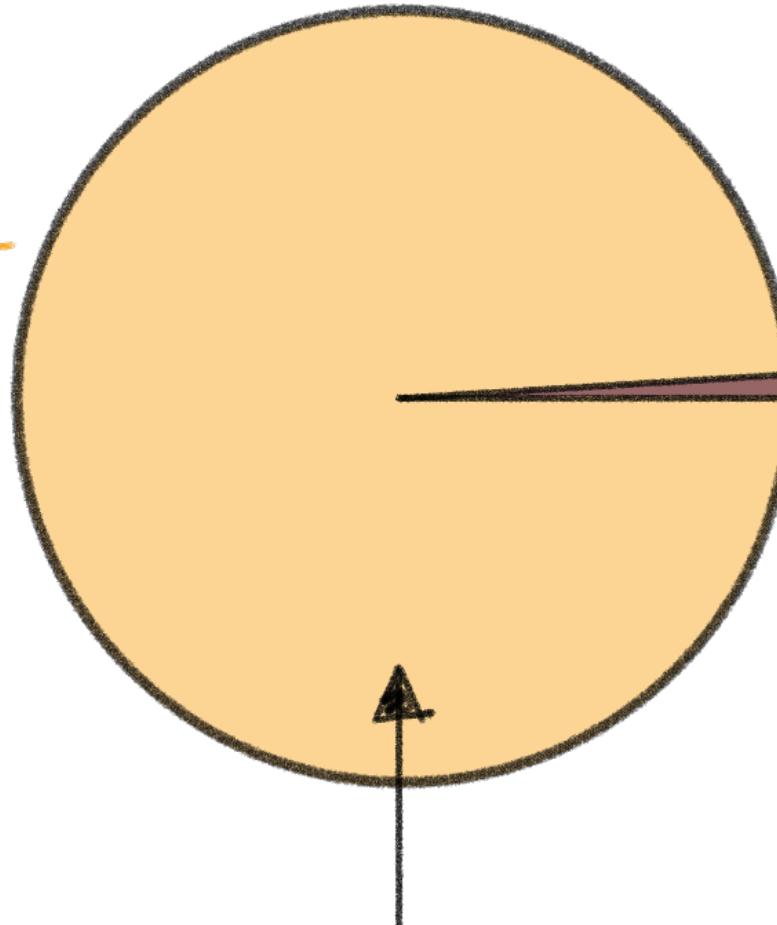
$$\Rightarrow \begin{cases} \frac{\partial L}{\partial \mu} = \frac{\partial L}{\partial z} \cdot \frac{\partial z}{\partial \mu} \\ \frac{\partial L}{\partial \sigma} = \frac{\partial L}{\partial z} \cdot \frac{\partial z}{\partial \sigma} \end{cases}$$



We could learn negative σ !
Sol: learn $\log(\sigma)$ then take \exp to get σ .

DATA

unlabeled

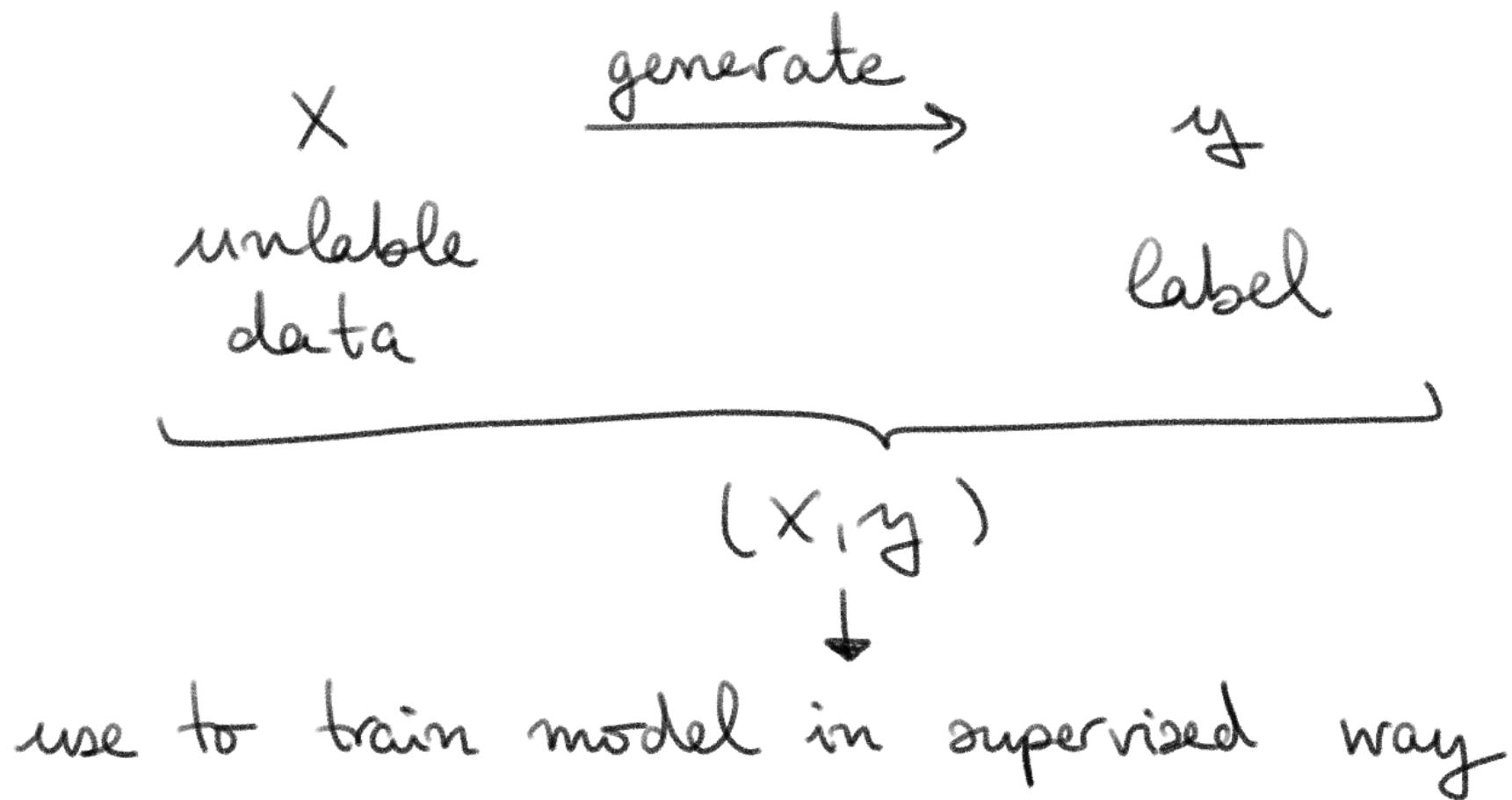


labeled

how do we leverage all this data ?

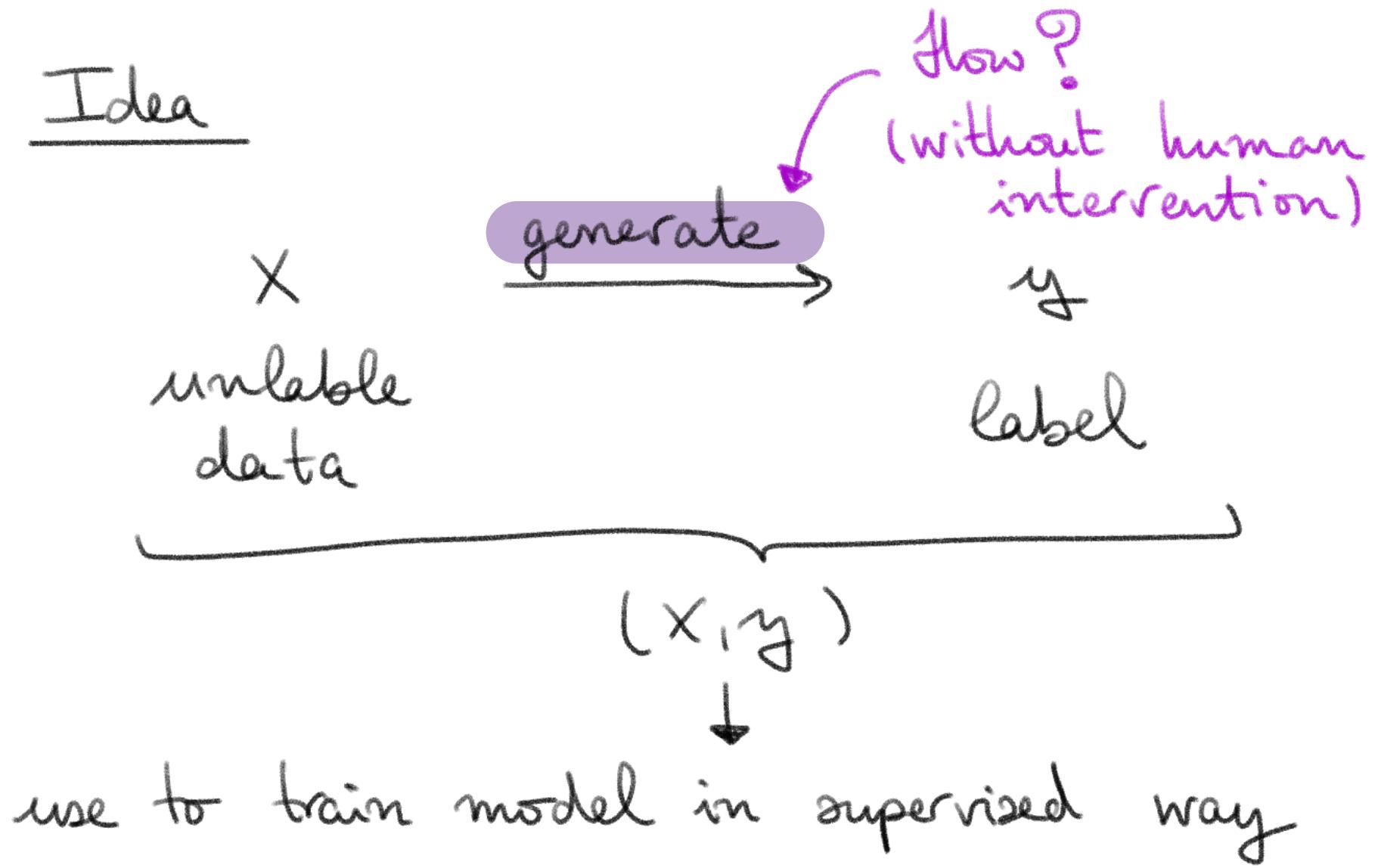
SELF-SUPERVISED LEARNING

Idea



SELF-SUPERVISED LEARNING

Idea



PRETEXT vs DOWNSTREAM TASKS

Pretext Task.- Designed so that NN learns visual features.

Downstream Task.- Has real world applications and human-labeled data.

EXAMPLES : DOWNSTREAM TASKS

Classification

Regression

Object Detection

Segmentation

Anomaly Detection

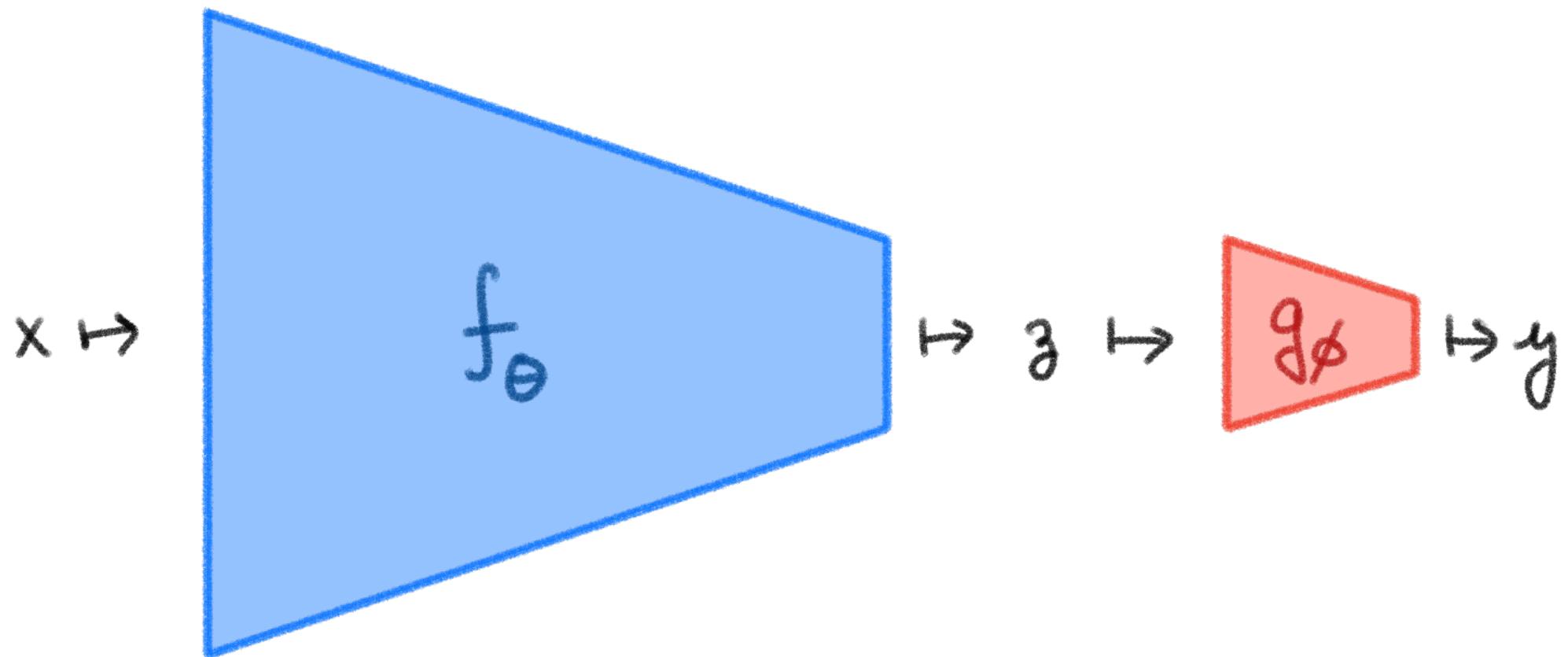
...

EXAMPLES: PRETEXT TASKS

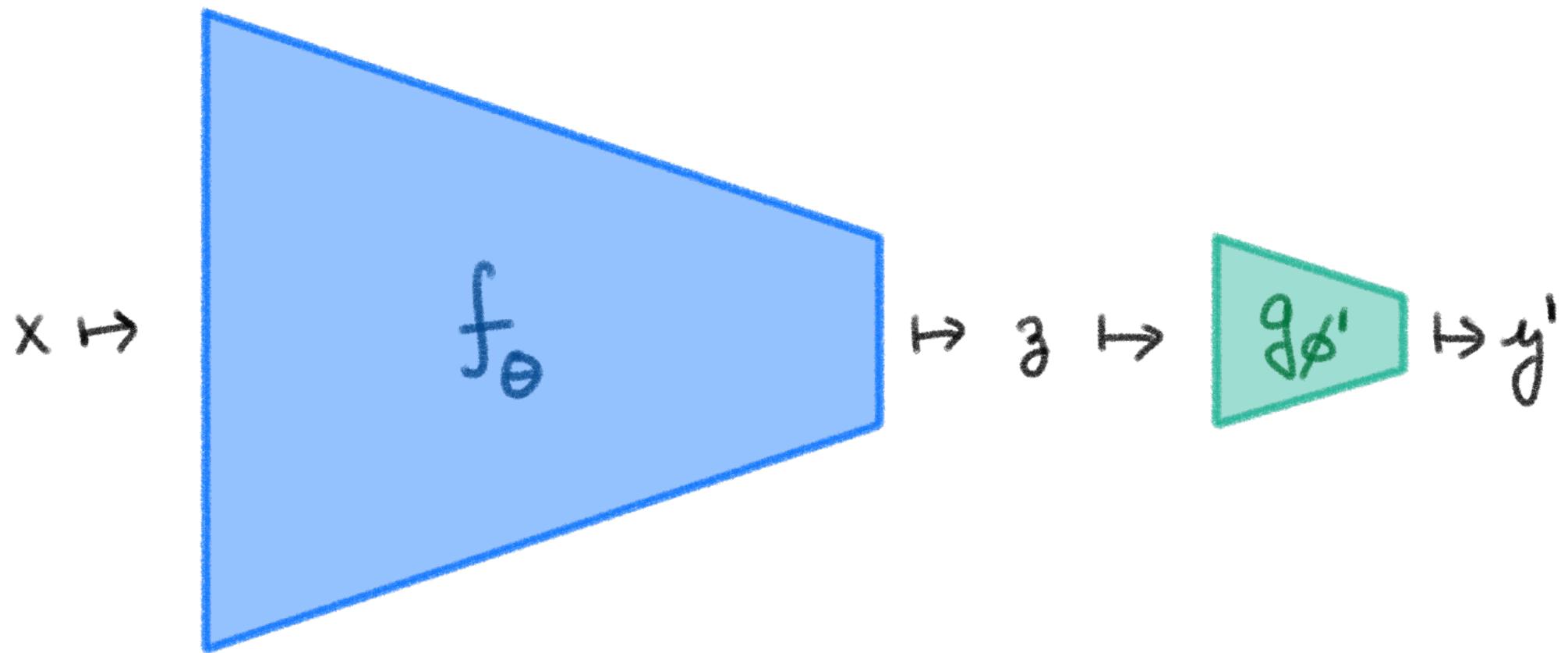
- Identifying augmentations
- Colorizing
- Masking
- Inpainting
- Contrastive Learning

How can a model
trained on a (pretext) task
be useful for a
different (downstream) task ?

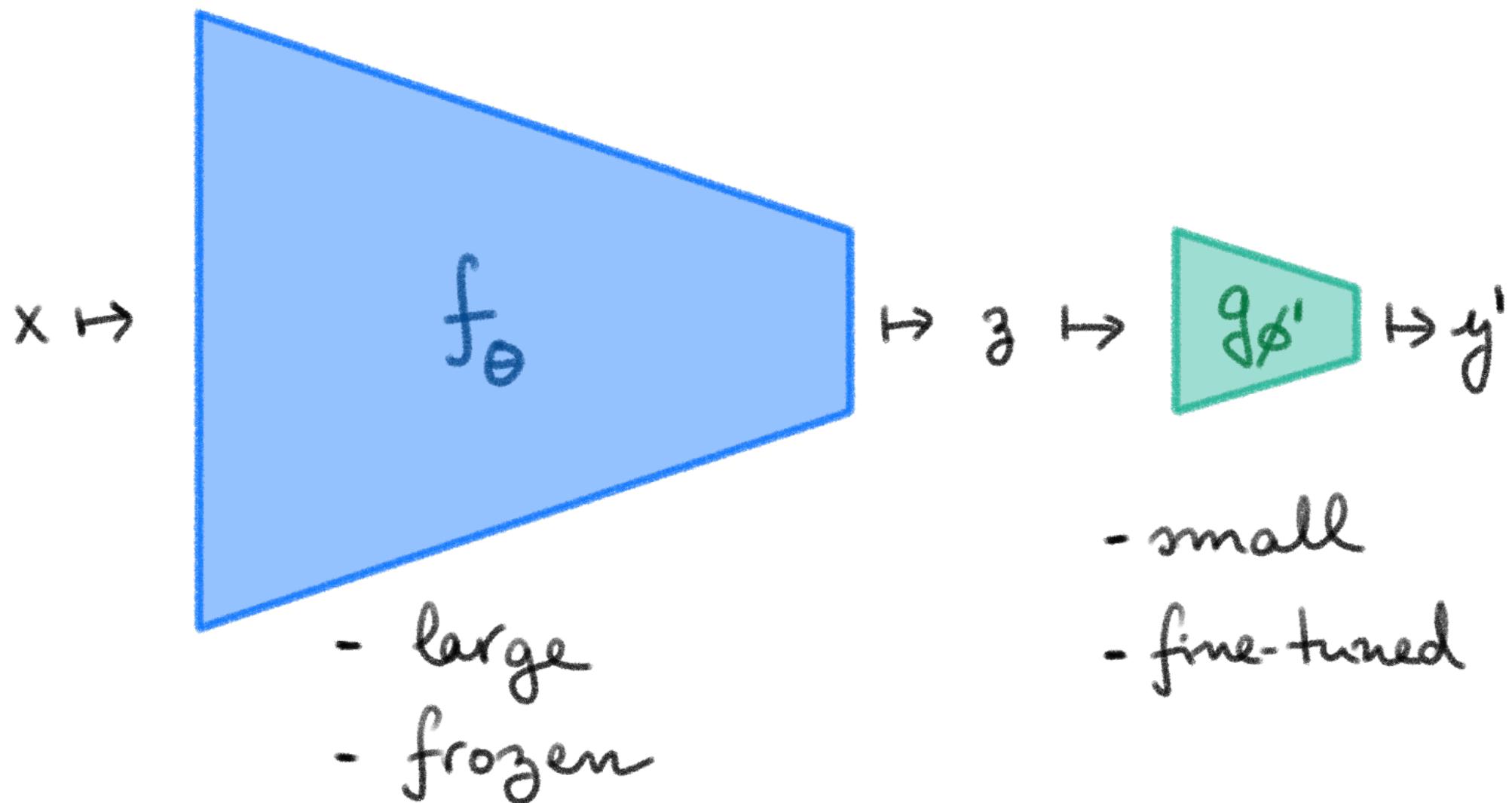
BACKBONE & HEAD



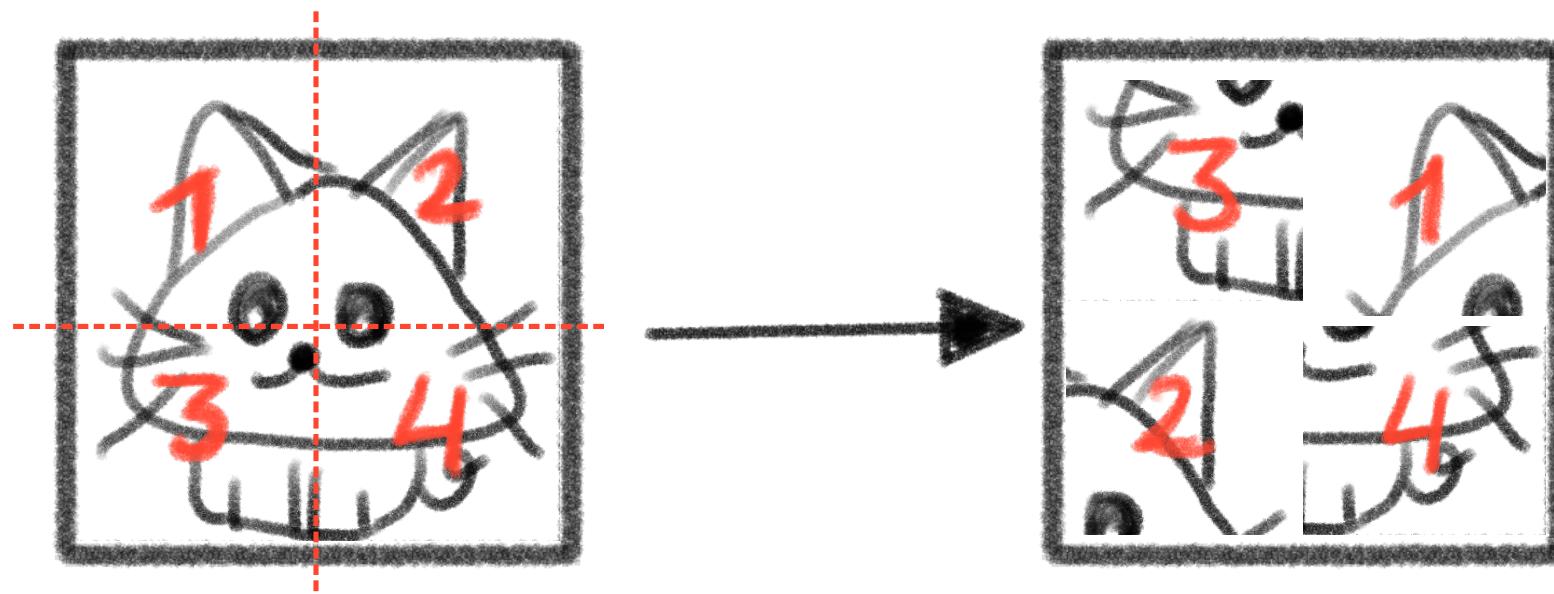
BACKBONE & HEAD



BACKBONE & HEAD



IDENTIFYING AUGMENTATIONS: — JIGSAW PUZZLE —

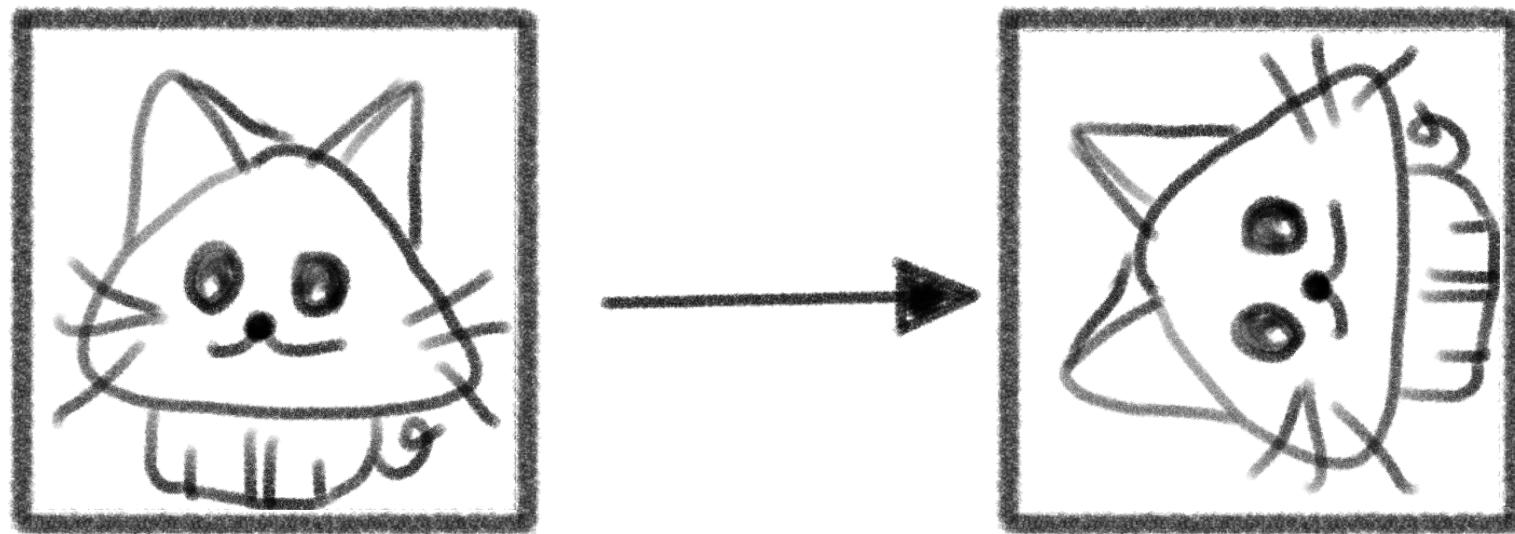


Loss?

$$y = \begin{pmatrix} 1 & 2 & 3 & 4 \\ 2 & 3 & 1 & 4 \end{pmatrix}$$

IDENTIFYING AUGMENTATIONS:

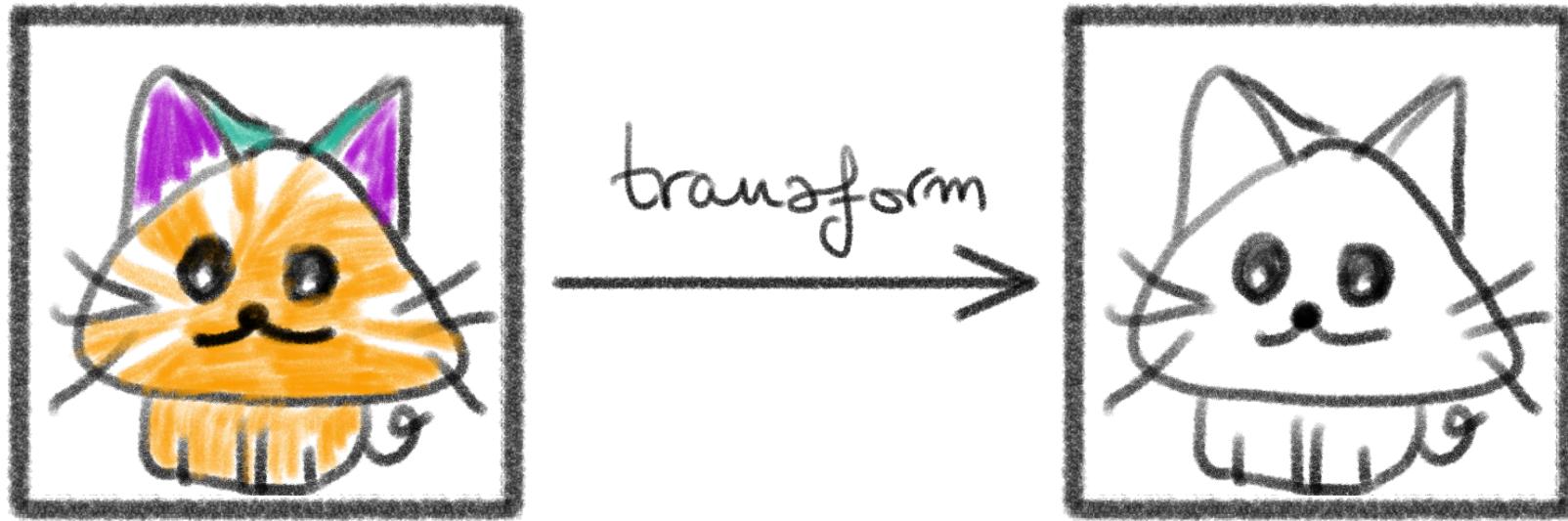
— ROTATIONS —



LOSS?

$y = 90^\circ$

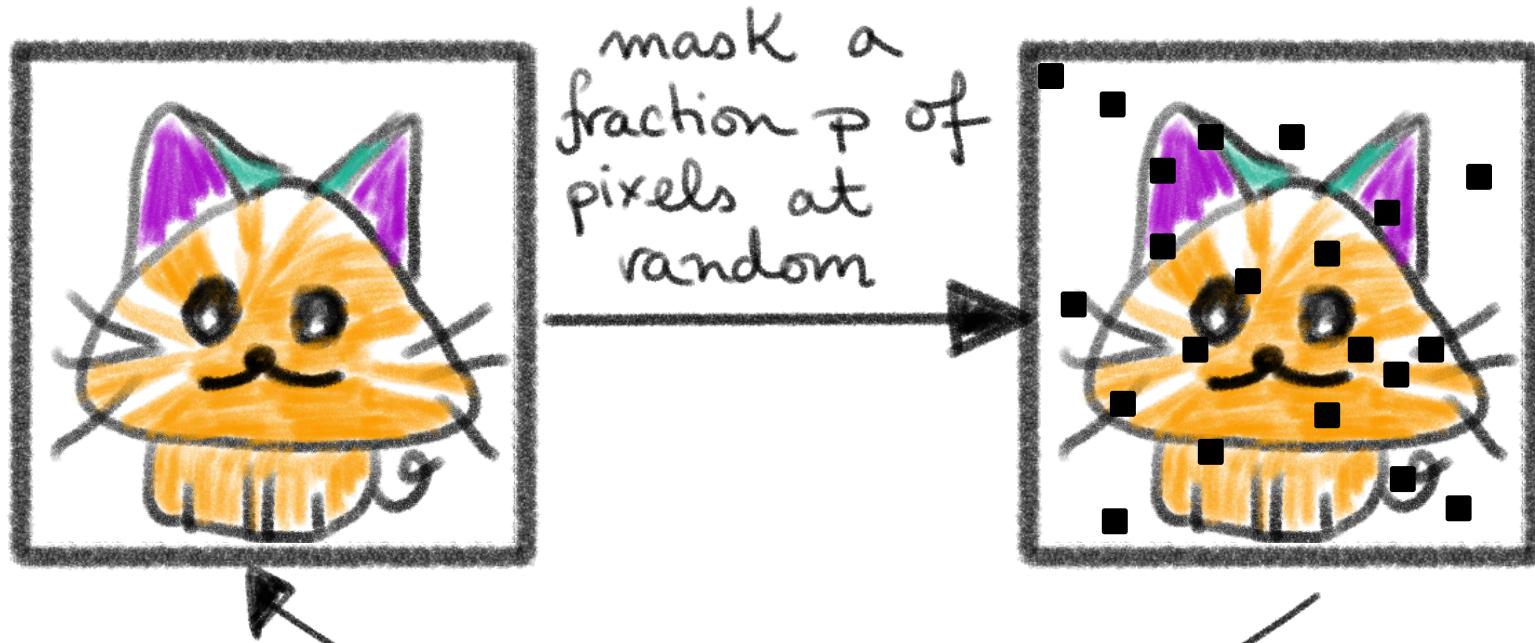
COLORIZING



train to predict

LOSS?

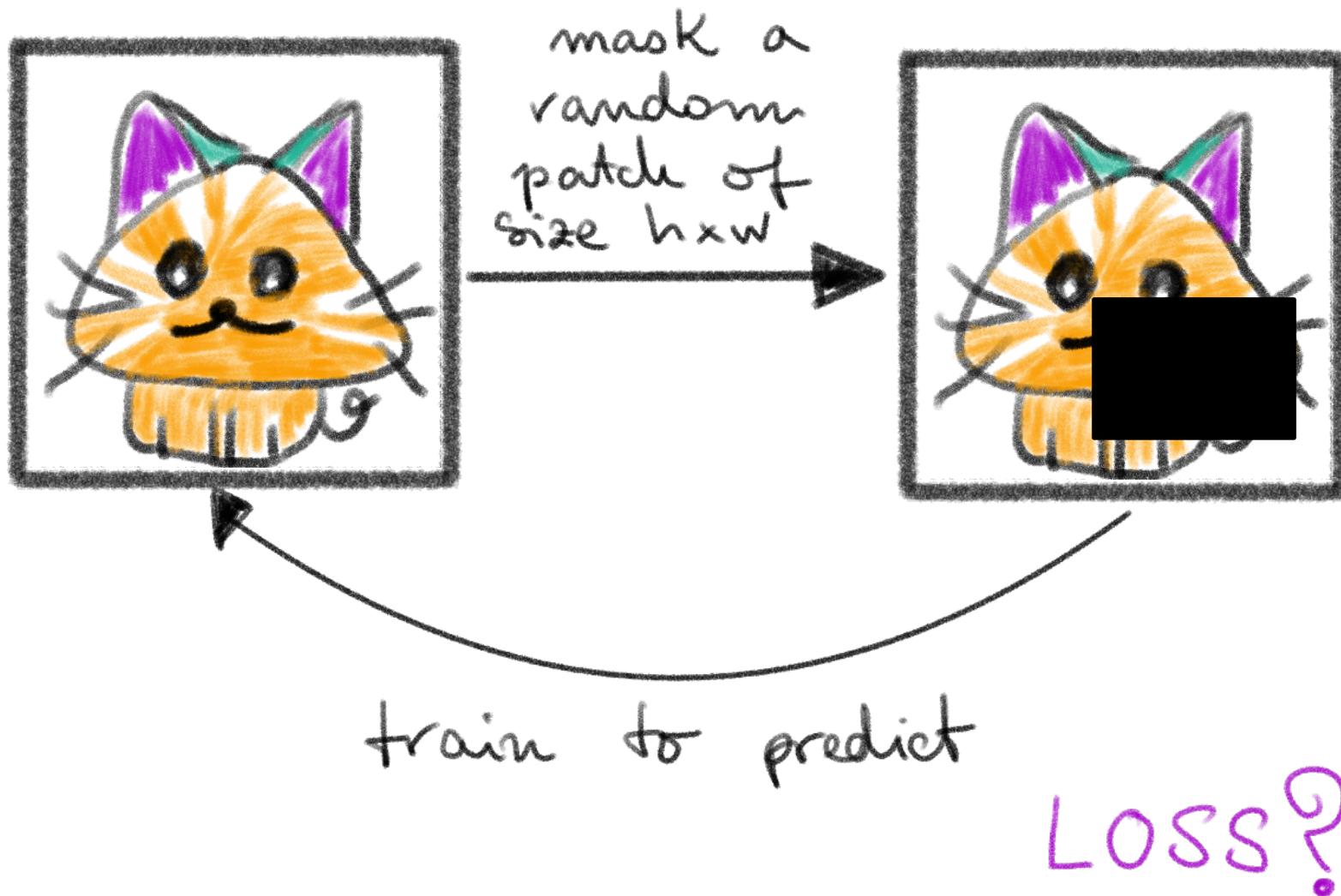
MASKING



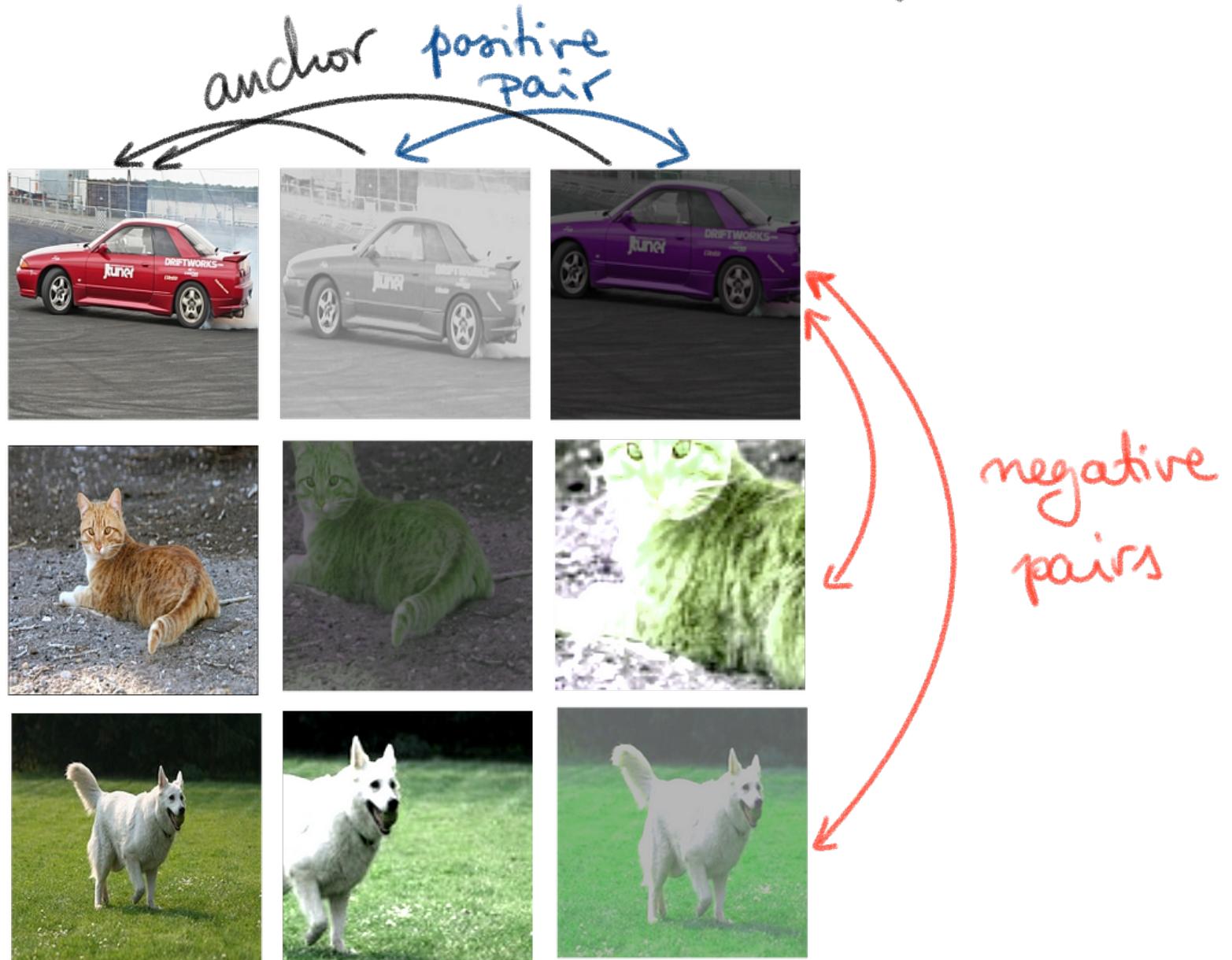
train to predict

LOSS?

INPAINTING



CONTRASTIVE LEARNING



CONTRASTIVE LEARNING

Idea

Contrastive loss :

"pull together" positive pairs

"push apart" negative pairs

CONTRASTIVE Loss

$$\text{Loss}(x, x') = \underbrace{\mathbb{1}_{x \sim x'}}_{\text{positive pair}} \|E(x) - E(x')\|_2^2 + \underbrace{\mathbb{1}_{x \neq x'}}_{\text{negative pair}} \max(0, m - \|E(x) - E(x')\|_2)^2$$

↑
margin

CONTRASTIVE Loss

$$\text{Loss}(x, x') = \underbrace{\mathbb{1}_{x \sim x'} \|E(x) - E(x')\|_2^2}_{\text{positive pair}} + \underbrace{\max(0, m - \|E(x) - E(x')\|_2)^2}_{\begin{array}{l} \text{margin} \\ \uparrow \end{array}}$$

negative pair

$\text{Loss} \approx 0 \Rightarrow$ if $x \sim x' \Rightarrow E(x) \simeq E(x')$

if $x \not\sim x' \Rightarrow \|E(x) - E(x')\|_2 \geq m$

TRIPLET Loss

$$\text{Loss} = \max(0, d(A, P) - d(A, N) + \alpha)$$

similarity metric: \uparrow anchor \downarrow
euclidean/cosine/... positive input \uparrow
negative input \downarrow margin \uparrow

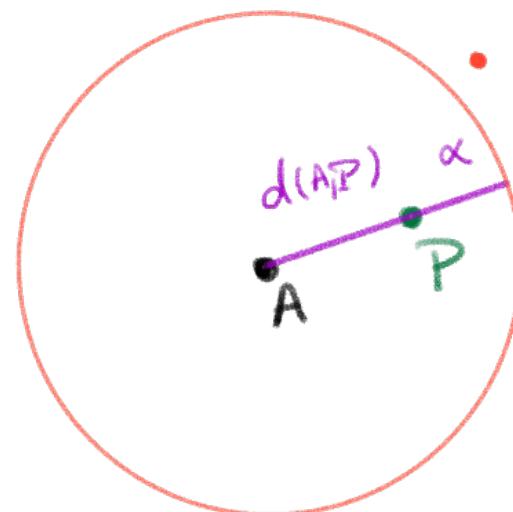
TRIPLET Loss

$$\text{Loss} = \max(0, d(A, P) - d(A, N) + \alpha)$$

↓ anchor ↓ negative input
↑ similarity metric: ↑ positive input ↑ margin
euclidean/cosine/...

$$\text{Loss} = 0 \Rightarrow$$

$$d(A, P) + \alpha \leq d(A, N)$$



Kaloot

Time!

Is SSL REALLY USEFUL ?

- Really useful for text data
- Really useful for multi-modal data
- For Computer Vision, struggles to beat purely supervised learning in most downstream tasks (e.g. classification) *
- Really useful for Anomaly Detection

J E P A (JOINT EMBEDDING PREDICTIVE ARCHITECHTURE)

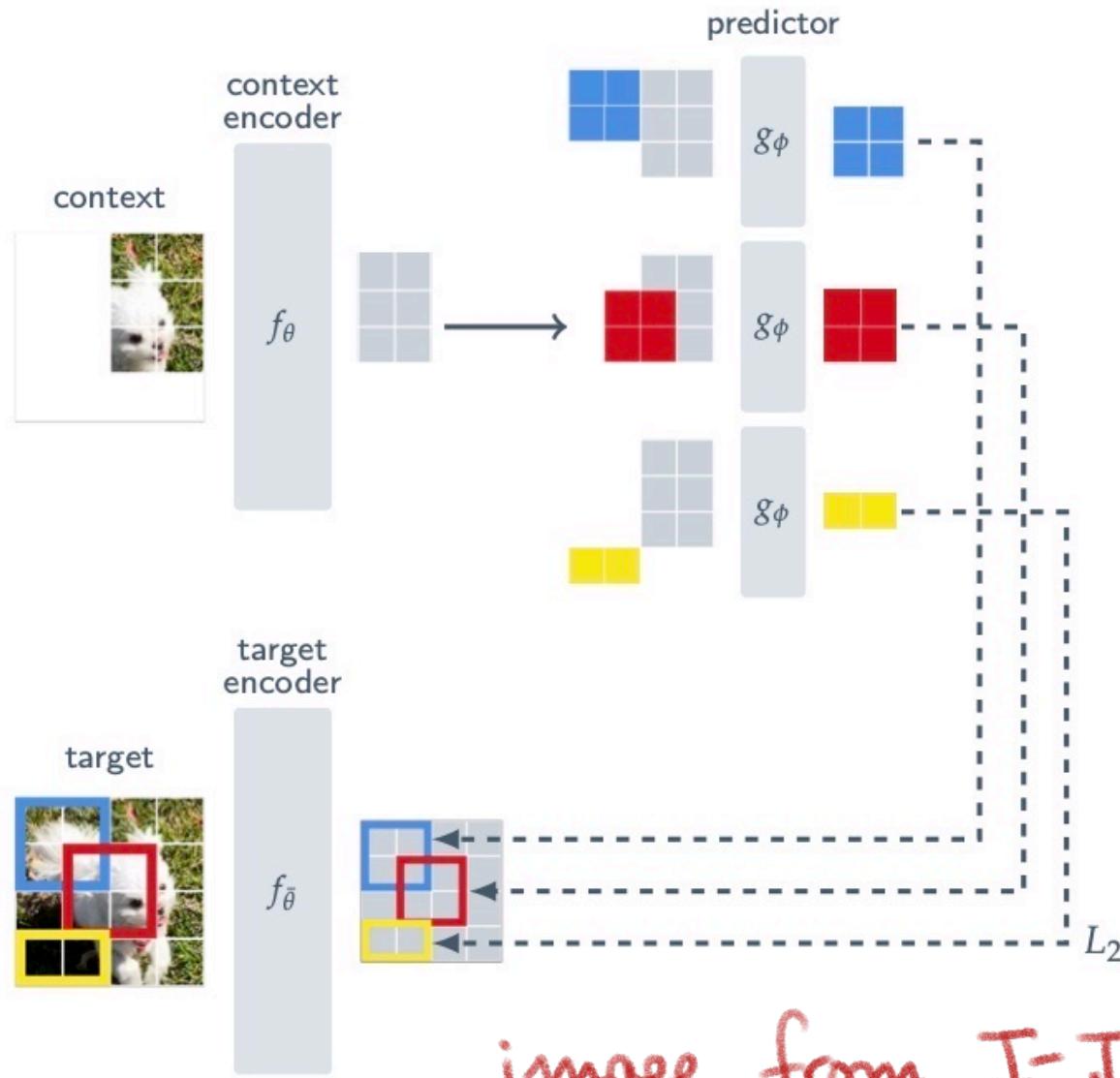


image from I-JEPA paper

DINO - v3

- Very good at Dense Prediction Tasks
- Good at semantic tasks
- Universal, frozen backbone
 - ⇒ robust across domains

Promising for medical imaging,
satellite imaging ...