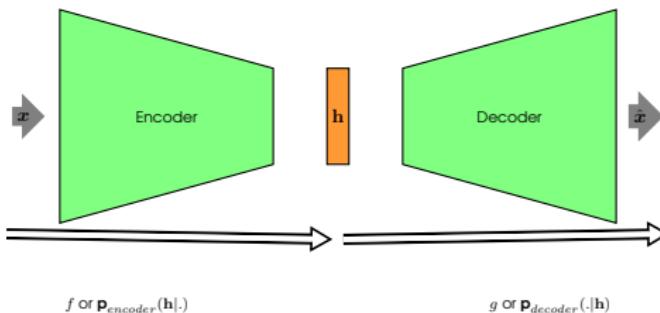


AUTOENCODERS

Vincent Barra
LIMOS, UMR 6158 CNRS, Université Clermont Auvergne

WHAT IS AN AUTOENCODER ?

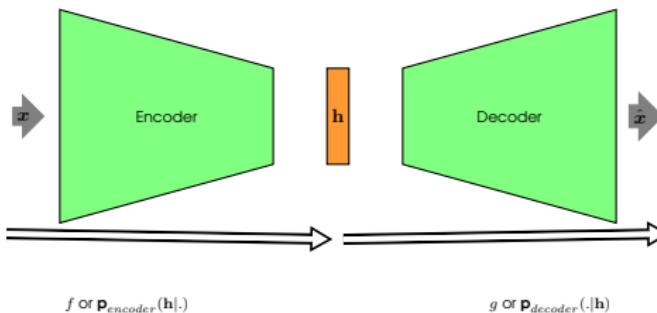


Key ideas

- ▶ A neural network trained using unsupervised learning
- ▶ Trained to copy its input to its output
- ▶ Learns an embedding h

$$\hat{\mathbf{x}} = g[f(\mathbf{x})] \quad \mathbf{h} = f(\mathbf{x})$$

WHAT IS AN AUTOENCODER ?



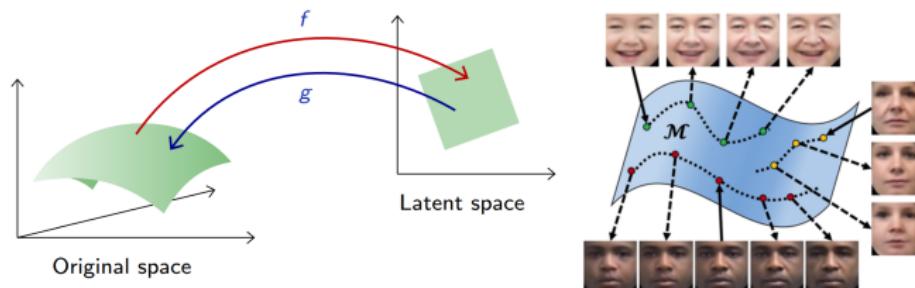
Key ideas

- ▶ $f : \mathbb{R}^n \rightarrow \mathbb{R}^q$: Encoder (MLP, CNN, RNN,..)
- ▶ $g : \mathbb{R}^q \rightarrow \mathbb{R}^n$: Decoder (MLP, CNN, RNN,..)
- ▶ $h = f(x)$: Encoded representation of x into a latent space of dimension q
- ▶ $q > n$: overcomplete / $q < n$: undercomplete

LATENT SPACE

Undercomplete autoencoder

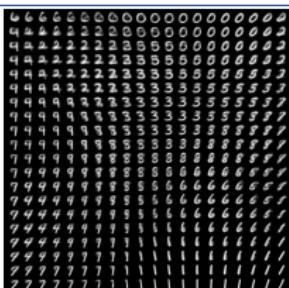
$(\forall \mathbf{x}) f(\mathbf{x})$ lies onto a manifold in \mathbb{R}^q



WHAT IS LEARNED ?

Learn relevant features

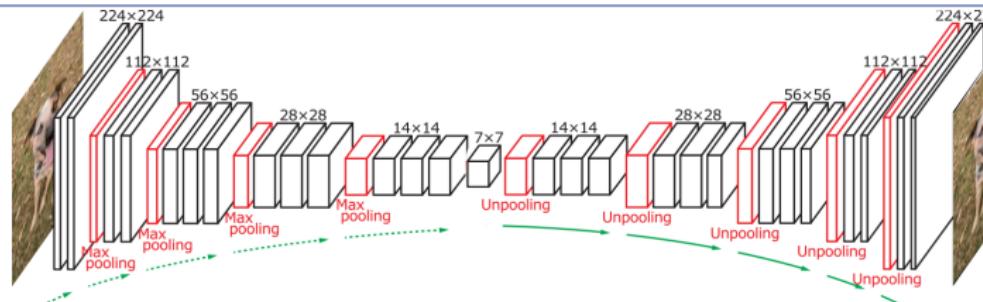
- ▶ Learning $\hat{\mathbf{x}} = \mathbf{x}$ is not useful ...
 - ▶ Autoencoders are designed to be unable to copy perfectly
 - ▶ Forced to capture most salient features of training data
 - ▶ So... what to retain ?
 - ▶ f and g can be probabilistic mappings $\mathbf{p}_{encoder}(\mathbf{h}|\cdot)$ and $\mathbf{p}_{decoder}(\cdot|\mathbf{h})$



WHAT IS LEARNED ?

What/how to learn ?

- ▶ Autoencoders are composed of two neural nets
- ▶ $\Rightarrow f = f_{\mathbf{W}_1, \mathbf{b}_1}$ and $g = g_{\mathbf{W}_2, \mathbf{b}_2}$
 - As usual, definition of a loss function to be minimized
 - Gradients computation using backpropagation
- ▶ Can also be trained using recirculation
 - Compare activations on \hat{x} to activations of \hat{x}
 - Biologically plausible than backpropagation, rarely used



LOSS FUNCTION

Loss function

$$\min_{W_1, \mathbf{b}_1, W_2, \mathbf{b}_2} L(\mathbf{x}, g[f(\mathbf{x})]) + \lambda \Omega(\mathbf{h})$$

- ▶ L : penalizing $g[f(\mathbf{x})]$ for being dissimilar from \mathbf{x}
- ▶ $\Omega(\mathbf{h})$: Regularization term

Regularization

Allows

- ▶ Sparsity of representation
- ▶ Robustness to missing inputs
- ▶ Robustness to noise
- ▶ ...

A SIMPLE EXAMPLE

Encoder = Decoder = 1 hidden layer MLP

$$\mathbf{x} \in \mathbb{R}^n \xrightarrow[f]{\quad} \mathbf{h} \in \mathbb{R}^q \xrightarrow[g]{\quad} \hat{\mathbf{x}} \in \mathbb{R}^n$$

σ_1, σ_2 : element-wise real activation functions, $W \in \mathcal{M}_{qn}(\mathbb{R})$, $b_1 \in \mathbb{R}^q$, $b_2 \in \mathbb{R}^n$

$$f(\mathbf{x}) = \sigma_1(W\mathbf{x} + b_1) \quad \hat{\mathbf{x}} = \sigma_2(W^T\mathbf{h} + b_2)$$

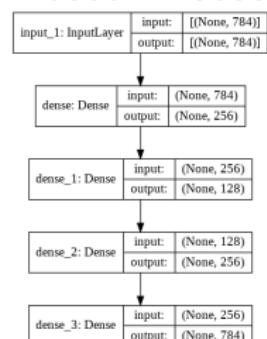
An autoencoder without regularization is trained to minimize

$$L(\mathbf{x}, \hat{\mathbf{x}}) = \|\mathbf{x} - \hat{\mathbf{x}}\|^2 = \|\mathbf{x} - \sigma_2(W^T(\sigma_1(W\mathbf{x} + b_1) + b_2))\|^2$$

If σ_2 is linear and $q < n$, this autoencoder is equivalent to PCA.

ADDING A LAYER...

Encoder = Decoder = 2 hidden layers MLP.. what you will have to do...

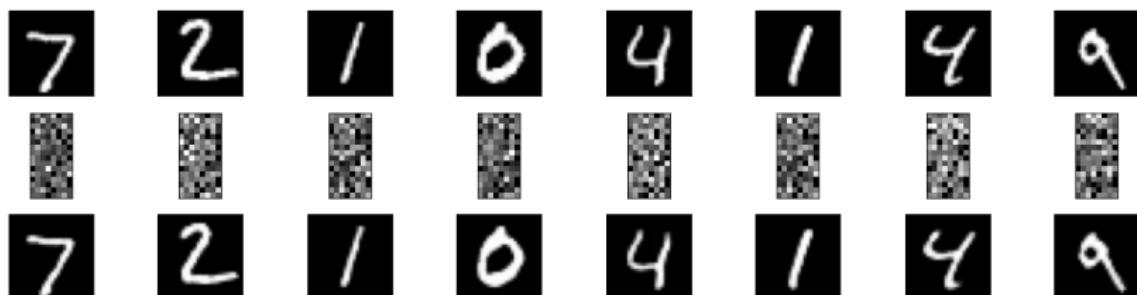


$$\sigma_{Encoder} = \{\text{Relu}, \text{Relu}\}$$

$$\sigma_{Decoder} = \{\text{Relu}, \text{sigmoid}\}$$

$$L = \text{MSE}$$

$$\text{Optimizer} = \text{Adam}.$$



ENCODER/DECODER CAPACITY

Not too much capacity

- ▶ too much capacity for f and $g \Rightarrow$ autoencoder can learn identity without learning any useful information about distribution of data
- ▶ if the encoder is very powerful, then $q = 1$ is possible and the decoder can learn to map back to the values of specific training examples

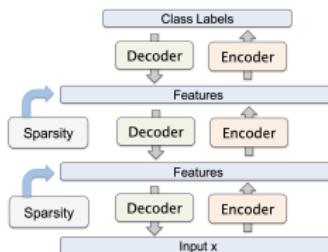
Failure if...

- ▶ capacity too high (controlled by depth)
- ▶ $q \geq n$ (if no strong constraint is applied on the loss)
⇒ Need for regularization.

REPRESENTATION POWER

Deep Encoder/Decoder can provide many advantages.
Common strategy: greedily pretrain a stack of shallow autoencoders

- ▶ The first layer is trained on input data
→ W_1, b_1
- ▶ Input → vector of the activation values of the hidden layer
- ▶ The second layer is trained on this vector
→ W_2, b_2
- ▶ Repeat the process



REGULARIZATION

Regularized autoencoders can help

- ▶ Rather than limiting model capacity by keeping encoder/decoder shallow and q small
- ▶ Use of $\Omega(h)$ in the loss function
- ▶ Forces the autoencoder to have more properties than the simple identity

Some properties

- ▶ Sparsity of representation
- ▶ Robustness to missing inputs
- ▶ Robustness to noise
- ▶ Even in $q \geq n$ can learn useful information
- ▶ ...

REGULARIZATION

Sparse autoencoder

- ▶ Only a few neurons are prone to activate when a single sample is input into the network
 - ▶ Can be done using $\Omega(\mathbf{h}) = \|\mathbf{W}\|_1$ (L_1 norm)
 - ▶ Can be done using the Kullback-Leibler divergence

Contractive autoencoder

- ▶ $\Omega(\mathbf{h}) = \sum_{i=1}^q \|\nabla_x h_i\|^2$
 - ▶ Learns a function that doesn't change much when \mathbf{x} slightly changes
 - ▶ Warps space: resists to perturbations of its input
 - ▶ Encourages to map a neighborhood of input \mathbf{x} to a smaller neighborhood of output points

DENOISING AUTOENCODERS

Receive a corrupted data as input and trained to predict the original, uncorrupted data as its output

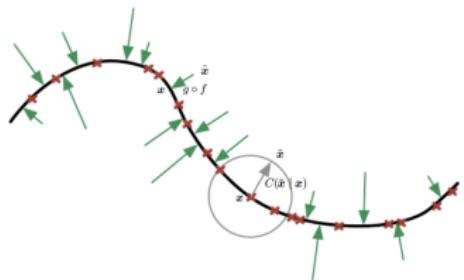
$$\text{Min } L(\mathbf{x}, \tilde{\mathbf{x}}) = -\log \mathbf{p}_{\text{encoder}}(\mathbf{x}|\mathbf{h} = f(\mathbf{x})) \quad \tilde{\mathbf{x}} \text{ corrupted version of } \mathbf{x}$$

Performs SGD on $\mathbb{E}_{\tilde{\mathbf{x}}} \log \mathbf{p}_{\text{decoder}}(\mathbf{x}|\mathbf{h} = f(\tilde{\mathbf{x}}))$



DENOISING AUTOENCODERS

- ▶ Can learn a corruption process $C(\tilde{x}|x)$ and the corresponding denoising process.
- ▶ Learns a vector field: the training samples lie on a low-dimensional manifold. The vector field estimates the slope of the density of data



- ▶ red crosses: training examples
- ▶ grey circle: equiprobable corruption of a training example
- ▶ green arrows: vector field

DENOISING AUTOENCODERS

Vector field

Encouraging the model to have the same score $S = \nabla_x \log \mathbf{p}(\mathbf{x})$ as the data distribution at every training point \mathbf{x}

Denoising autoencoder with Gaussian $\mathbf{p}(\mathbf{x} | \mathbf{h})$ estimates the score as

$$S \approx g[f(\mathbf{x})] - \mathbf{x}$$

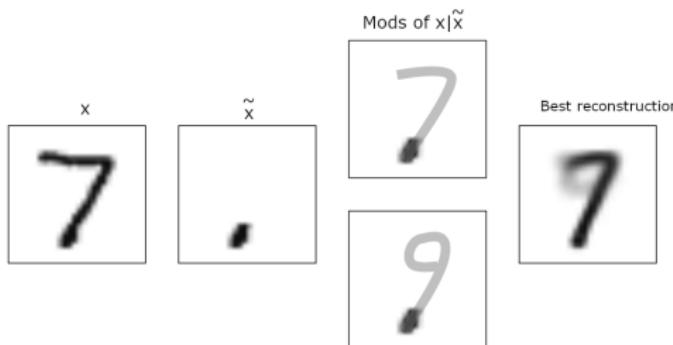
and is trained to minimize

$$\|g[f(\tilde{\mathbf{x}})] - \mathbf{x}\|^2$$

DENOISING AUTOENCODERS

The posterior $\mathbf{P}_{\mathbf{x}|\tilde{\mathbf{x}}}$ can be non-deterministic.

If $L(\mathbf{x}, \tilde{\mathbf{x}}) = \|\mathbf{x} - \tilde{\mathbf{x}}\|^2$: best reconstruction is $\mathbb{E}(\mathbf{x} | \tilde{\mathbf{x}})$: very unlikely under $\mathbf{P}_{\mathbf{x}|\tilde{\mathbf{x}}}$



Adversarial networks

Use in place of loss a second network that assesses if the output is realistic

Additional ressource

See lecture "Generative models (GAN)".

SUPER RESOLUTION

Special case of denoising autoencoders, where the encoder's input is smaller than the decoder output.

Original
7 2 1 0 4 1 4 9 5 9 0 6
9 0 1 5 9 7 8 4 9 6 6 5
4 0 7 4 0 1 3 1 3 4 7 2

Input
7 2 1 0 4 1 4 9 5 9 0 6
9 0 1 5 9 7 8 4 9 6 6 5
4 0 7 4 0 1 3 1 3 4 7 2

Bilinear interpolation
7 2 1 0 4 1 4 9 5 9 0 6
9 0 1 5 9 7 8 4 9 6 6 5
4 0 7 4 0 1 3 1 3 4 7 2

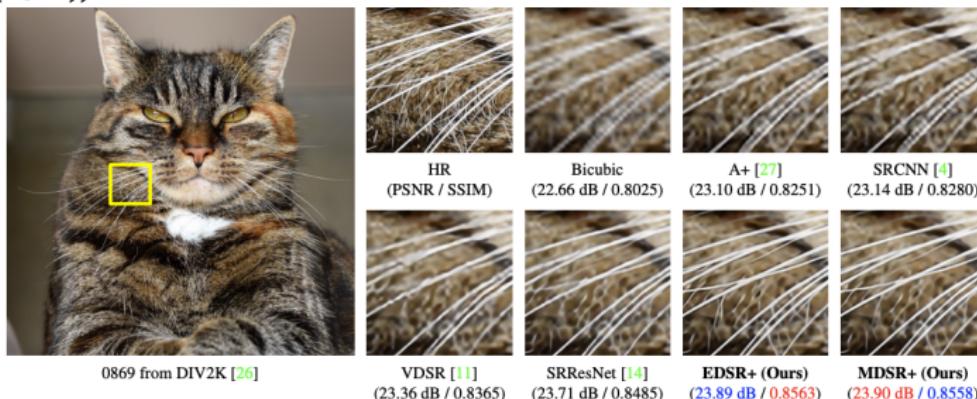
Original
7 2 1 0 4 1 4 9 5 9 0 6
9 0 1 5 9 7 8 4 9 6 6 5
4 0 7 4 0 1 3 1 3 4 7 2

Input
7 2 1 0 4 1 4 9 5 9 0 6
9 0 1 5 9 7 8 4 9 6 6 5
4 0 7 4 0 1 3 1 3 4 7 2

Autoencoder output
7 2 1 0 4 1 4 9 5 9 0 6
9 0 1 5 9 7 8 4 9 6 6 5
4 0 7 4 0 1 3 1 3 4 7 2

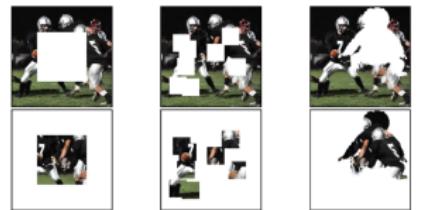
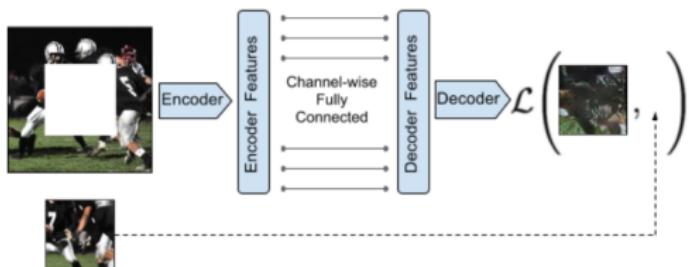
SUPER RESOLUTION

For super resolution on real images, use deep nets (here 2 different ResNets, Lim et al (2017))

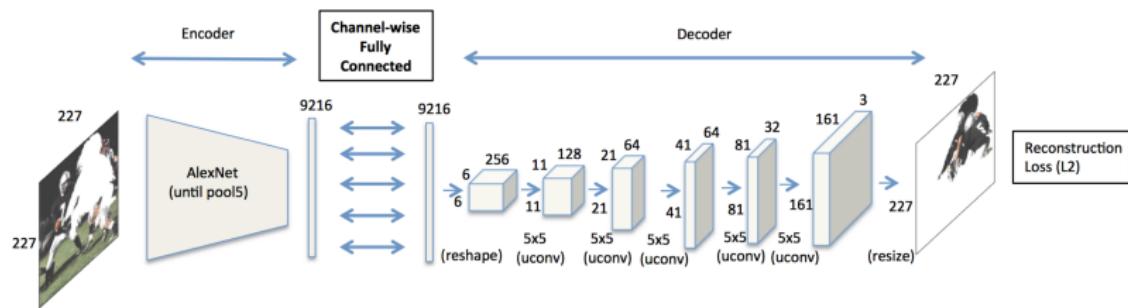


CONTEXT ENCODERS

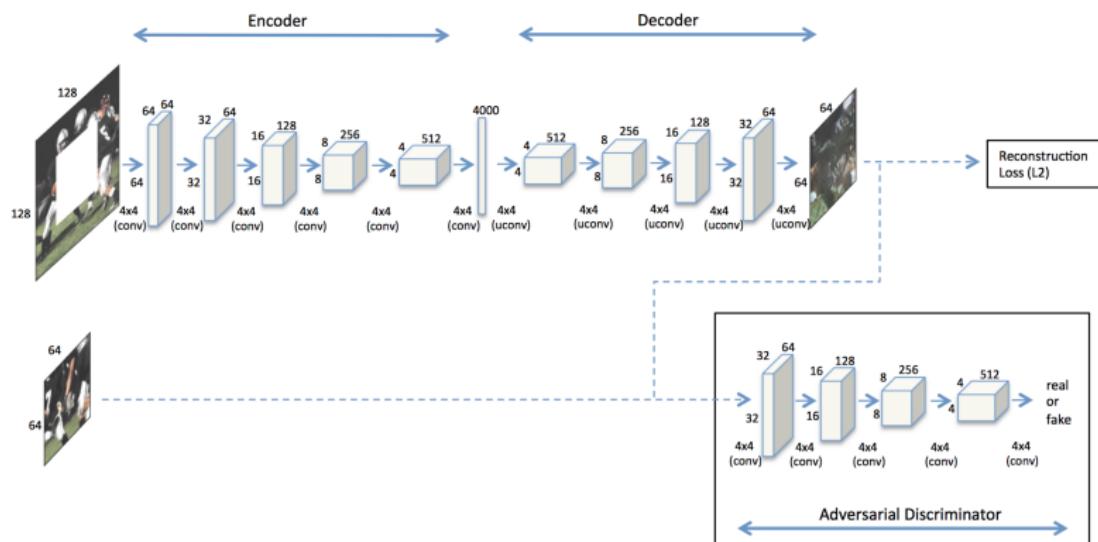
An interesting work of Pathak et al., 2016



CONTEXT ENCODERS



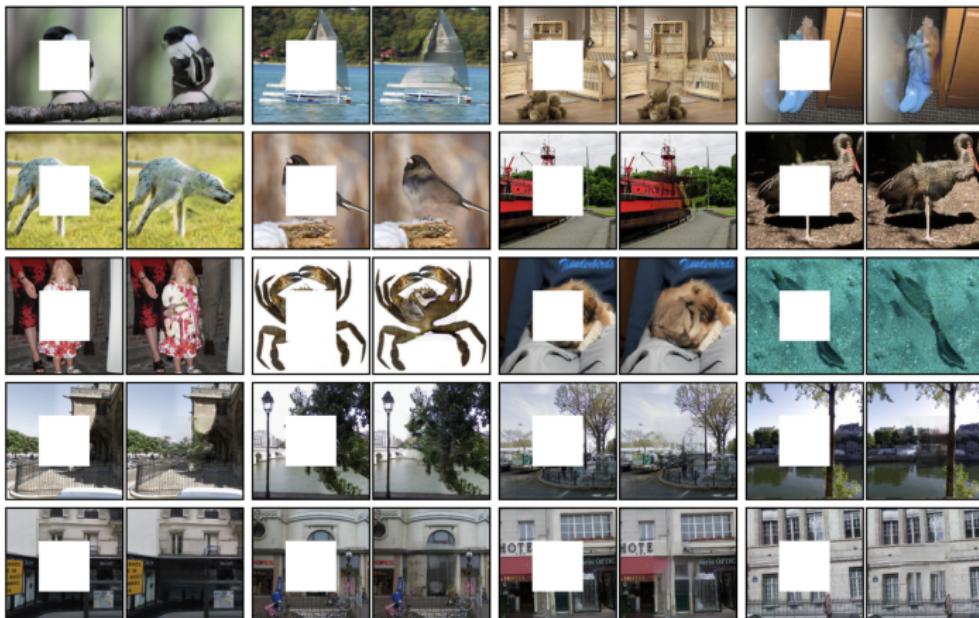
CONTEXT ENCODERS



Additional ressource

See lecture "Generative models (GAN)".

CONTEXT ENCODERS



LATENT SPACE

Can be used in several ways:

- ▶ Once learned, the Encoder provides a concise way to compress the data, while retaining relevant information
- ⇒ input to classic machine learning methods
- ▶ Directly working in the latent space : analysis, generation with the decoder...

Additional ressource

See lecture “Generative models (VAE)”.

MANIFOLD IN LATENT SPACE

Hypothesis

$f(\forall \mathbf{x}) (\mathbf{x})$ concentrates around a low dimensional manifold \mathcal{M} in \mathbb{R}^q .

"Justification"

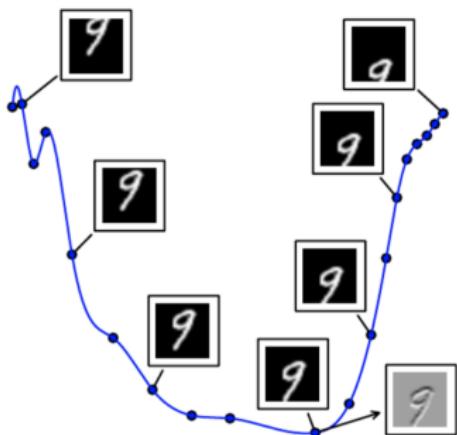
- ▶ \mathbf{x} = image of size $m \times n$, each pixel being coded between 0 and 255.
- ▶ Set of all images : $256^{m \times n}$ possible images.
- ▶ "Cat" images" impose constraints on grey level distribution
- ⇒ Less degrees of freedom in the high dimensional space



Autoencoders aim to learn the structure of \mathcal{M} .

MANIFOLD IN LATENT SPACE

Manifolds are defined by tangent planes: for $x \in \mathcal{M}$, specifies how x can change while staying on \mathcal{M} .

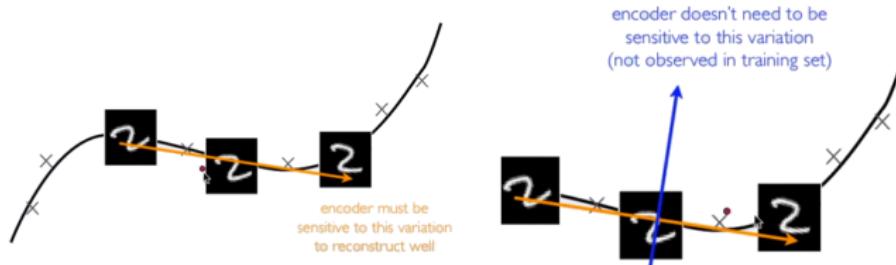


Moving along tangent

- ▶ gray pixels: don't change
- ▶ white pixels brighten
- ▶ black pixels darken

MANIFOLD IN LATENT SPACE

- ▶ Encoder captures the information needed to reconstruct x
→ representation h
- ▶ If data generating distribution concentrates near a low-dimensional \mathcal{M} , h implicitly captures a local coordinate system for \mathcal{M}
 - Only variations tangential to \mathcal{M} at x need to correspond to changes in h
 - The learned f is only sensitive to changes in the tangent plane, not to orthogonal changes.



MANIFOLD IN LATENT SPACE

An example of application: interpolation between images

