



Deep Learning

Chapter 9: Autoencoders

Bernd Bischl

Department of Statistics – LMU Munich

Winter term 2020



Unsupervised learning

UNSUPERVISED LEARNING

- So far, we were dealing with different types of neural networks designed for classification and regression tasks.
- In these **supervised learning** scenarios, we exploit information of class memberships (or numeric values) to train our algorithm. That means in particular, that we have access to labeled data.
- There exists another learning paradigm, **unsupervised learning**, where:
 - training data consists of unlabeled input points $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(n)}$
 - and one aims at finding and describing intrinsic structure in the data.
- There is much more unlabeled data than labeled! But what can we learn from it?

UNSUPERVISED LEARNING - EXAMPLES

1. Clustering

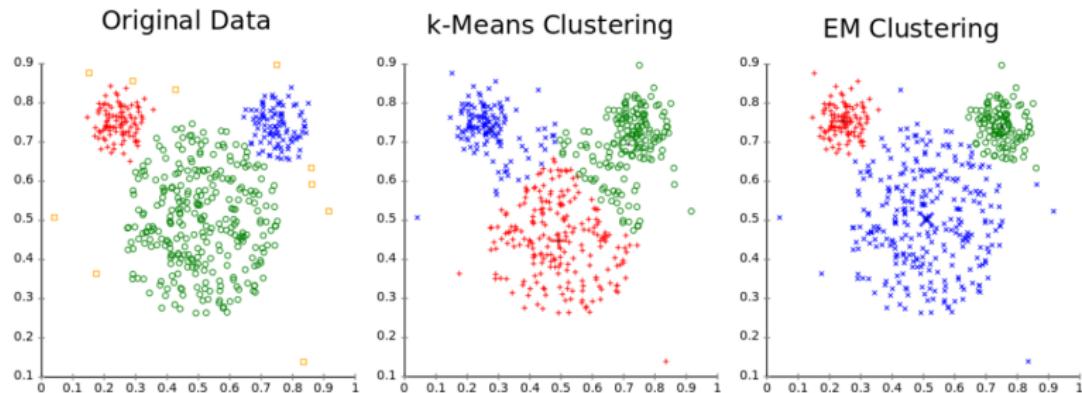


Figure: Different cluster analysis results on a dataset. True labels (the colors in the original data) are shown here but the algorithms only operate on unlabeled data.
(Source: Wikipedia)

UNSUPERVISED LEARNING - EXAMPLES

2. Dimensionality reduction/manifold learning

- E.g. for visualisation in a low dimensional space.

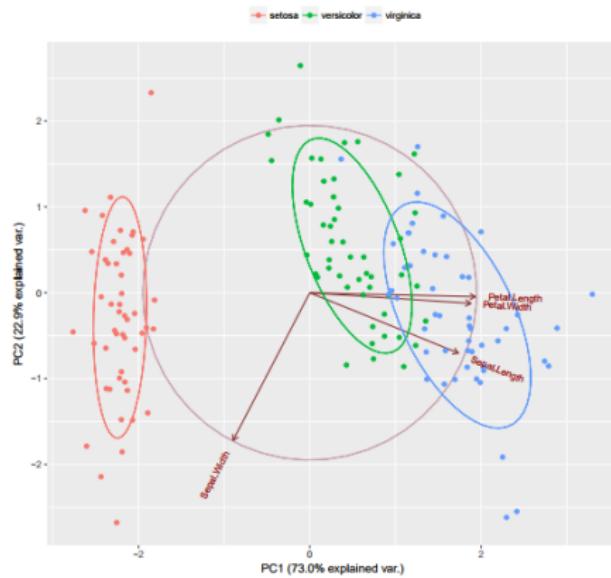


Figure: Principal Component Analysis (PCA)

UNSUPERVISED LEARNING - EXAMPLES

2. Dimensionality reduction/manifold learning
 - E.g. for image compression.

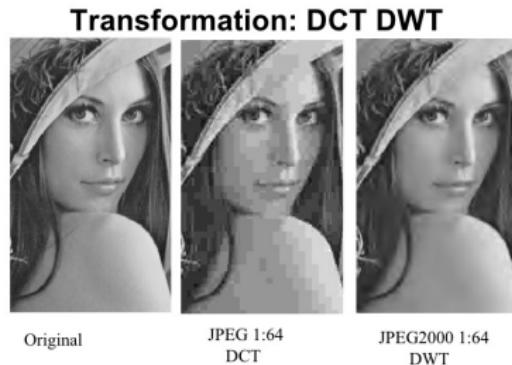


Figure: from <https://de.slideshare.net/hcycon/bildkompression>

UNSUPERVISED LEARNING - EXAMPLES

3. Feature extraction/representation learning

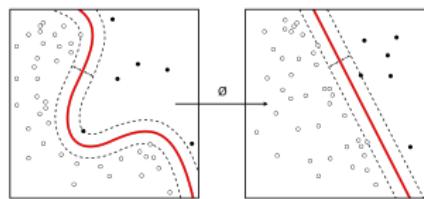


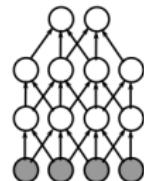
Figure: Source: Wikipedia

- E.g. for **semi-supervised learning**: features learned from an unlabeled dataset are employed to improve performance in a supervised setting.

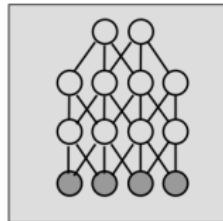
UNSUPERVISED LEARNING - EXAMPLES

4. Density fitting/learning a generative model

approx. inference
model $q(x, h)$



trained model
 $p^*(x, h)$



approx. inference
model $p(x, h)$

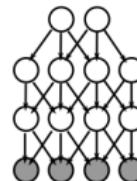


Figure: A generative model can reconstruct the missing portions of the images. (Bornschein, Shabanian, Fischer & Bengio, ICML, 2016)

UNSUPERVISED DEEP LEARNING

Given i.i.d. (unlabeled) data $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n \sim \mathbb{P}_x$, in unsupervised deep learning, one usually trains:

- an autoencoder (a special kind of neural network) for **representation learning** (feature extraction, dimensionality reduction, manifold learning, ...), or,

UNSUPERVISED DEEP LEARNING

Given i.i.d. (unlabeled) data $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n \sim \mathbb{P}_x$, in unsupervised deep learning, one usually trains:

- an autoencoder (a special kind of neural network) for **representation learning** (feature extraction, dimensionality reduction, manifold learning, ...), or,
- a **generative model**, i.e. a probabilistic model of the data generating distribution p_{data} (data generation, outlier detection, missing feature extraction, reconstruction, denoising or planning in reinforcement learning, ...).

Autoencoders

AUTOENCODER (AE)-TASK AND STRUCTURE

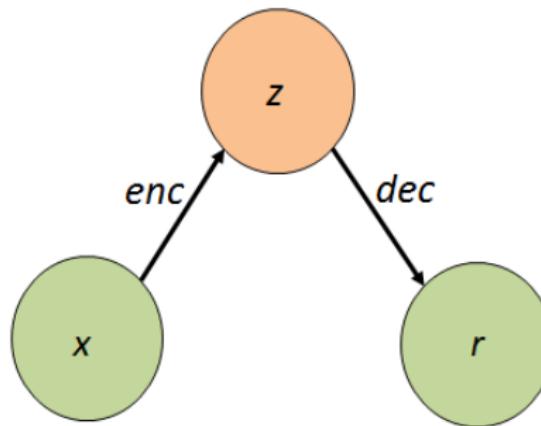
- Autoencoders (AEs) are a special kind of feedforward neural networks.
- Task: Compression and reconstruction of the input.
- They consist of two parts:
 - **encoder** function $\mathbf{z} = \text{enc}(\mathbf{x})$.
 - **decoder** that produces the reconstruction $\mathbf{r} = \text{dec}(\mathbf{z})$.
- Loss function:

$$L(\mathbf{x}, \text{dec}(\text{enc}(\mathbf{x})))$$

- Goal: Learn good **internal representations** \mathbf{z} (also called **code**).

AUTOENCODER (AE)- COMPUTATIONAL GRAPH

The general structure of an AE as a computational graph:



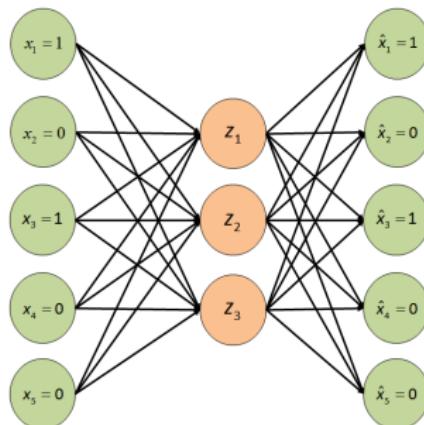
- An AE has two computational steps:
 - the encoder *enc*, mapping x to z .
 - the decoder *dec*, mapping z to r .

UNDERCOMPLETE AUTOENCODERS

- Learning simply the identity $\text{dec}(\text{enc}(\mathbf{x})) = \mathbf{x}$ is not useful.
- Therefore, we restrict the architecture, such that

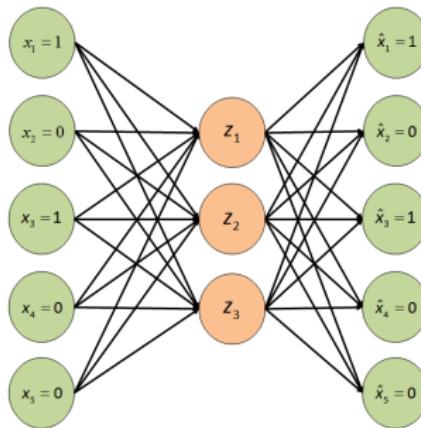
$$\dim(\mathbf{z}) < \dim(\mathbf{x})$$

- Such an AE is called **undercomplete**.



UNDERCOMPLETE AUTOENCODERS

- In other words: In an undercomplete AE, the hidden layer has fewer neurons than the input layer. \Rightarrow That will force the AE to
 - capture only the most salient features of the training data!
 - learn a “compressed” representation of the input.



UNDERCOMPLETE AUTOENCODERS

- Training an AE is done by minimizing the risk, where the loss function penalizes the reconstruction $dec(enc(\mathbf{x}))$ for differing from \mathbf{x} . The MSE

$$\|\mathbf{x} - dec(enc(\mathbf{x}))\|_2^2$$

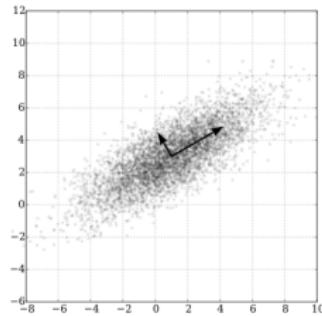
is a typical choice.

- For optimization, the very same optimization techniques as for standard feed-forward nets are applied (SGD, RMSProp, ADAM,...).

UNDERCOMPLETE AUTOENCODERS

- Example: For an undercomplete AE with
 - linear decoder $dec(\mathbf{z}) = \mathbf{Hz}$
 - squared error loss $L = ||\mathbf{x} - dec(enc(\mathbf{x}))||_2^2$
 - input normalized to have zero mean,
- the optimal encoder corresponds to **Principal Component Analysis (PCA)**.
- PCA is a statistical procedure that uses an orthogonal transformation to convert a set of observations of possibly correlated variables into a set of values of linearly uncorrelated variables called **principal components**.
- Transformation is defined in such a way that the first principal component has the largest possible variance (that is, accounts for as much of the variability in the data as possible)

UNDERCOMPLETE AUTOENCODERS



- An AE with a non-linear decoder/encoder can be seen as a non-linear generalization of PCA.
- Problem: If an AE is **overcomplete** ($\dim(\mathbf{z}) > \dim(\mathbf{x})$) or encoder and decoder are too powerful, the AE can learn to simply copy the input.

OVERCOMPLETE AE – PROBLEM

Example 1: Overcomplete AE (code dimension \geq input dimension).
⇒ even a linear AE can copy the input to the output without learning anything useful.

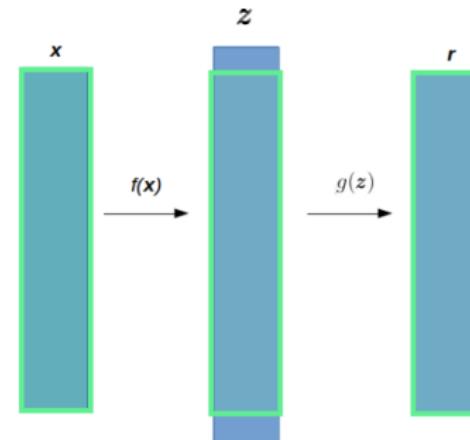


Figure: Overcomplete AE that learned to copy its inputs to the hidden layer and then to the output layer (Credits to A.-L. Popkes and P. Wenker).

TO POWERFUL AE – PROBLEM

Example 2: Very powerful nonlinear AE that learns a 1D code:

- Encoder: learns to map each training example $\mathbf{x}^{(i)}$ to the code i .
- Decoder: learns to map these integer indices back to the values of specific training examples.

EXPERIMENT: LEARN TO ENCODE MNIST

- Let us try to compress the MNIST data as good as possible.
- Therefore, we will fit an undercomplete autoencoder to learn the best possible representation,
→ as few as possible dimensions in the internal representation \mathbf{z} .

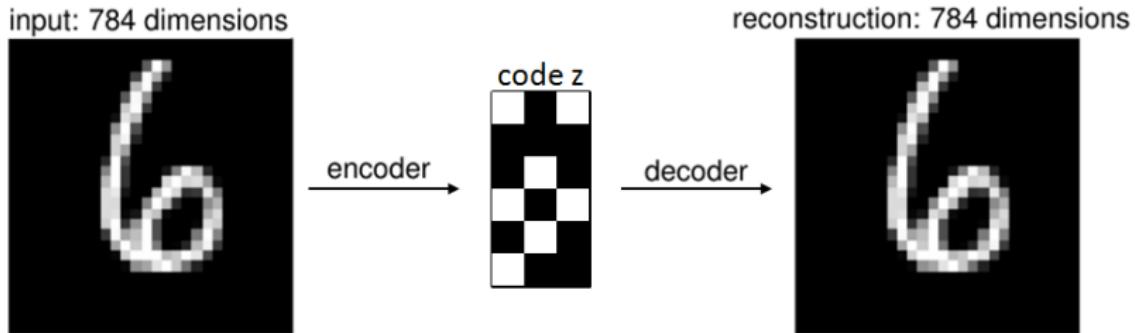


Figure: Flow chart of our autoencoder: reconstruct the input with fixed dimensions $dim(\mathbf{z}) << dim(\mathbf{x})$.

EXPERIMENT: LEARN TO ENCODE MNIST

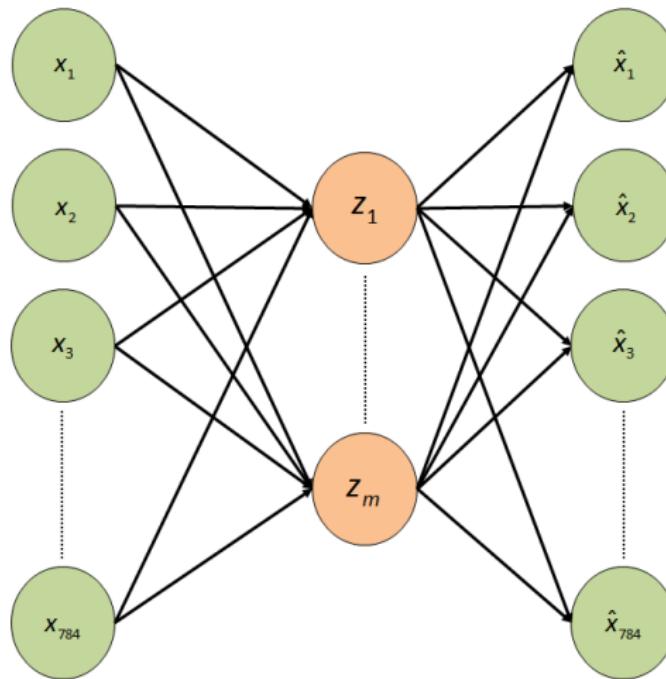


Figure: Architecture of the autoencoder.

EXPERIMENT: LEARN TO ENCODE MNIST

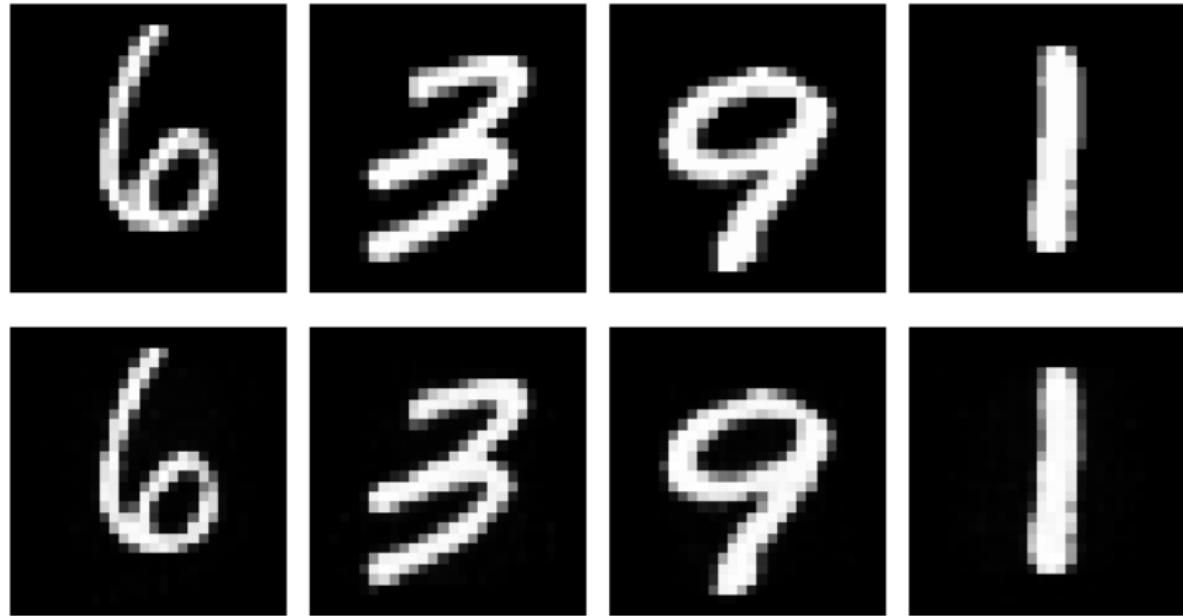


Figure: The top row shows the original digits, the bottom row the reconstructed ones.

- $\dim(\mathbf{z}) = 784 = \dim(\mathbf{x})$.

EXPERIMENT: LEARN TO ENCODE MNIST

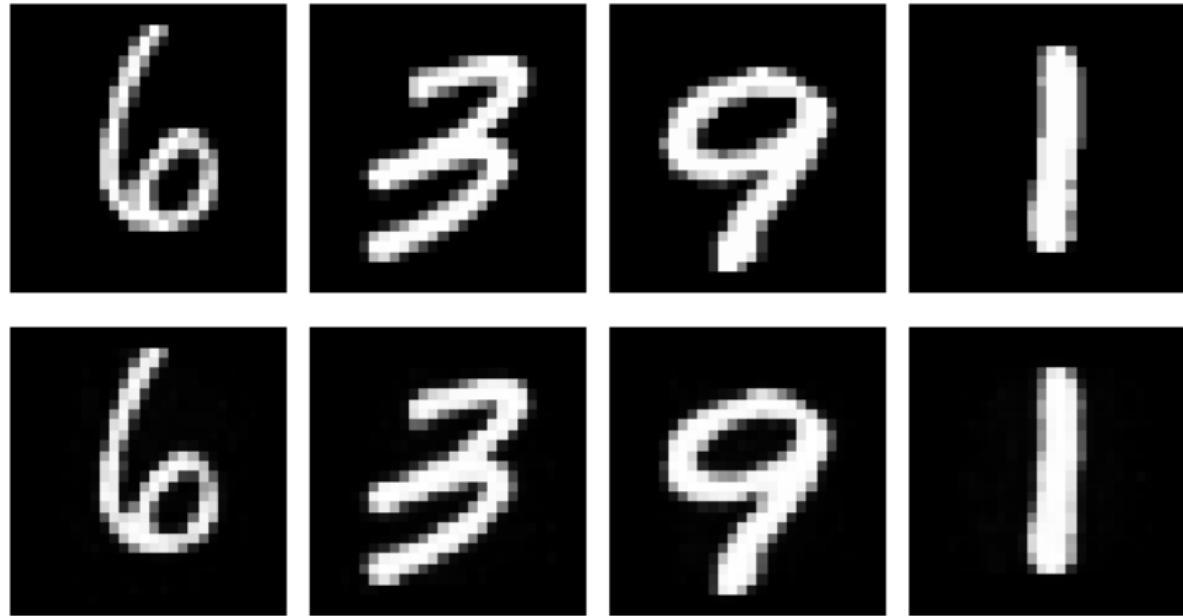


Figure: The top row shows the original digits, the bottom row the reconstructed ones.

- $\dim(\mathbf{z}) = 256$.

EXPERIMENT: LEARN TO ENCODE MNIST

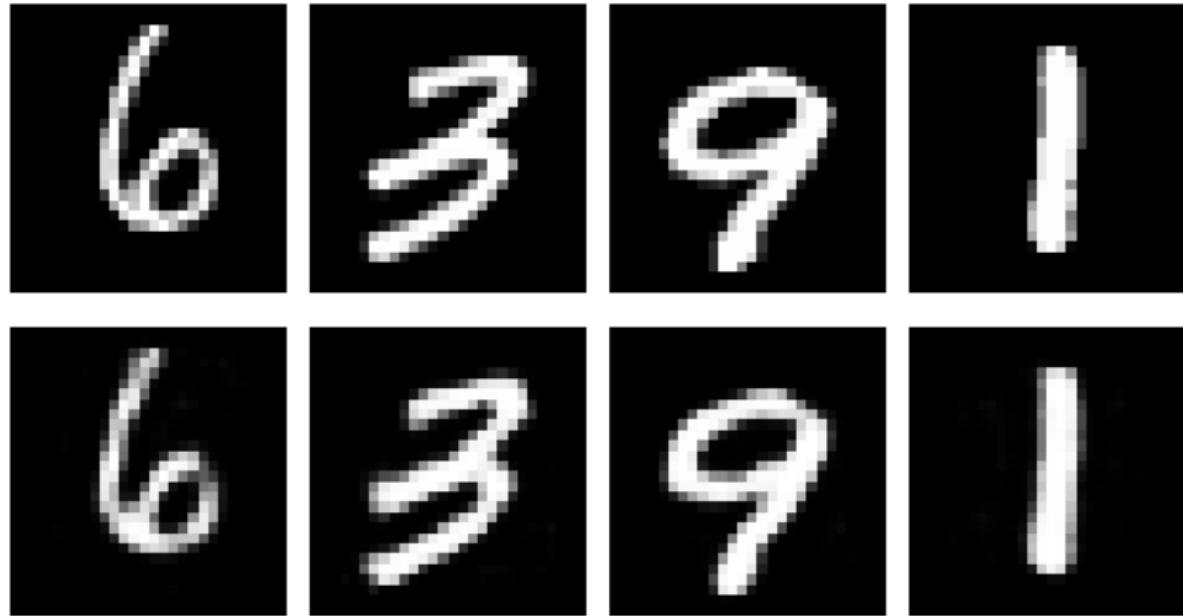


Figure: The top row shows the original digits, the bottom row the reconstructed ones.

- $\dim(\mathbf{z}) = 64$.

EXPERIMENT: LEARN TO ENCODE MNIST

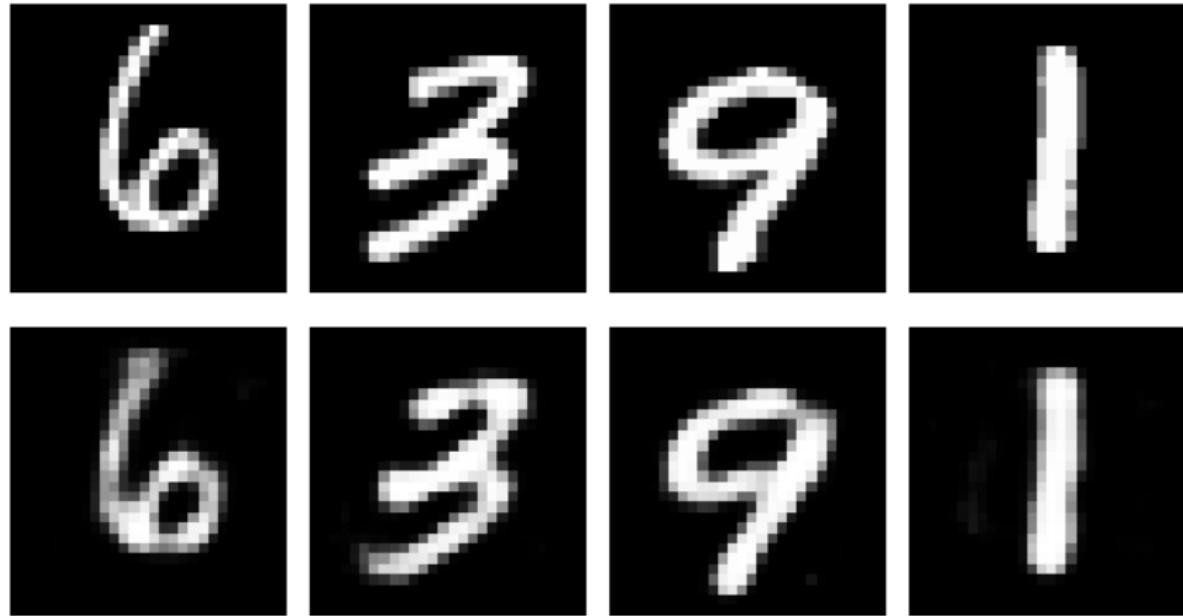


Figure: The top row shows the original digits, the bottom row the reconstructed ones.

- $\dim(\mathbf{z}) = 32$.

EXPERIMENT: LEARN TO ENCODE MNIST



Figure: The top row shows the original digits, the bottom row the reconstructed ones.

- $\dim(\mathbf{z}) = 16$.

EXPERIMENT: LEARN TO ENCODE MNIST

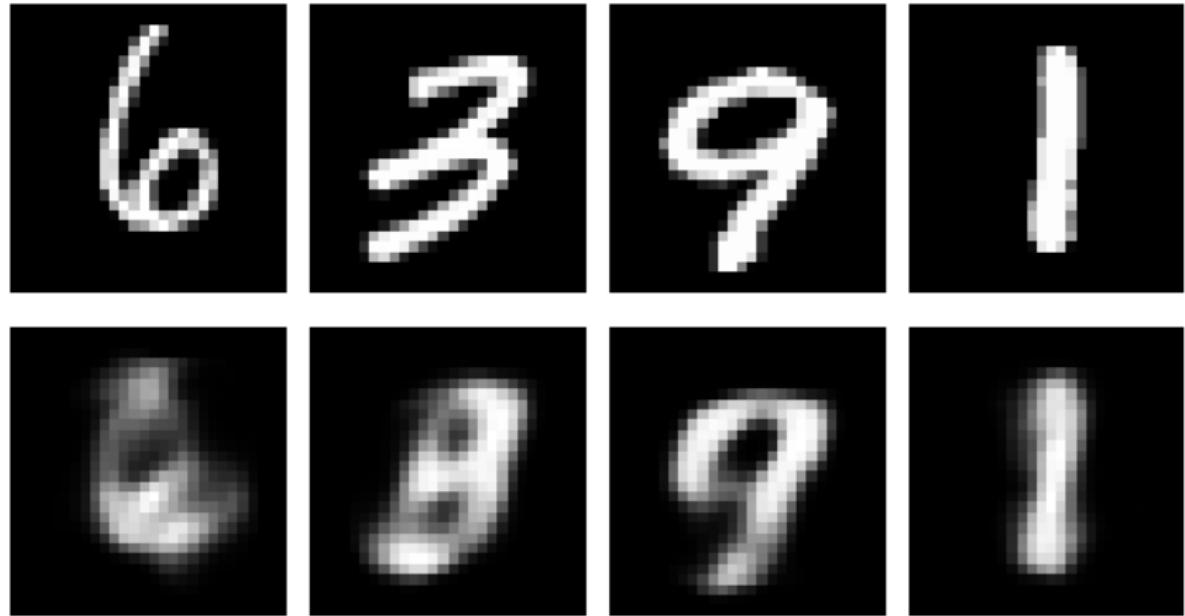


Figure: The top row shows the original digits, the bottom row the reconstructed ones.

- $\dim(\mathbf{z}) = 8$.

EXPERIMENT: LEARN TO ENCODE MNIST



Figure: The top row shows the original digits, the bottom row the reconstructed ones.

- $\dim(\mathbf{z}) = 4$.

EXPERIMENT: LEARN TO ENCODE MNIST

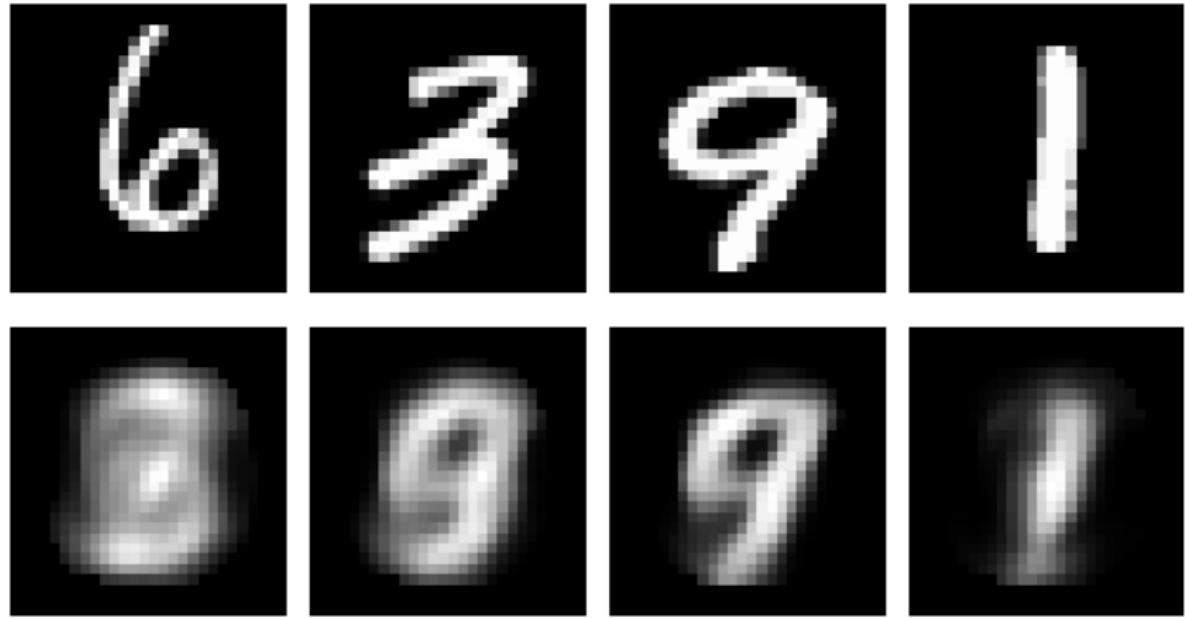


Figure: The top row shows the original digits, the bottom row the reconstructed ones.

- $\dim(\mathbf{z}) = 2$.

EXPERIMENT: LEARN TO ENCODE MNIST

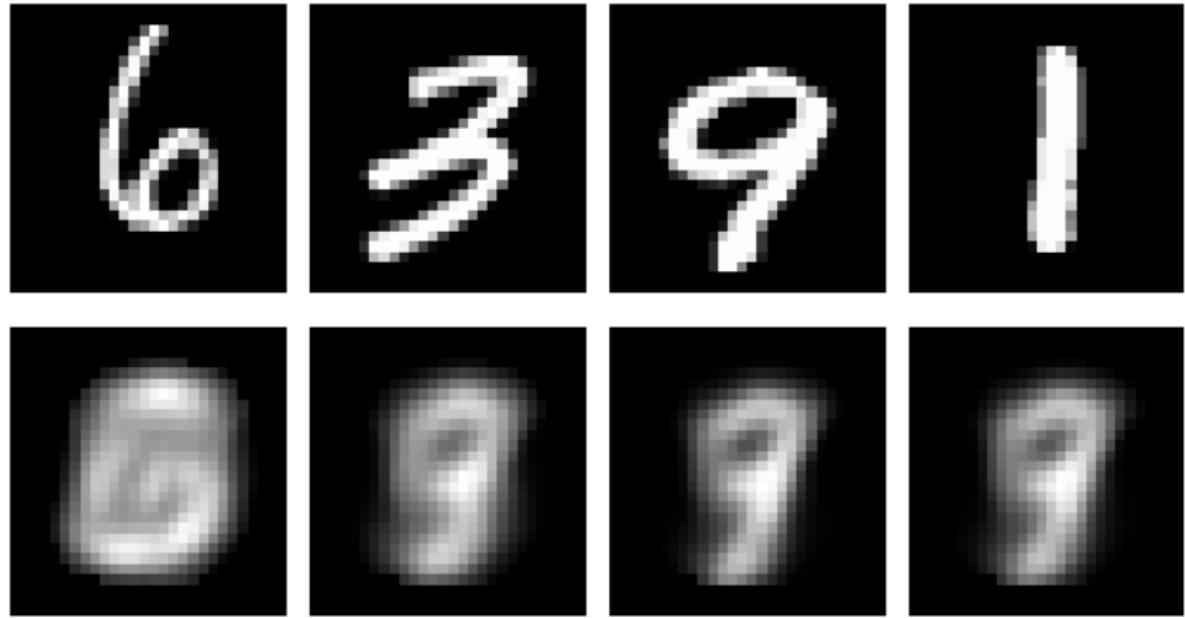


Figure: The top row shows the original digits, the bottom row the reconstructed ones.

- $\dim(\mathbf{z}) = 1$.

Regularized Autoencoder

REGULARIZED AUTOENCODER

- Goal: choose code dimension and capacity of encoder/decoder based on the problem.
- **Regularized AEs** modify the original loss function to:
 - prevent the network from trivially copying the inputs.
 - encourage additional properties.
- Examples:
 - **Sparse AE**: sparsity of the representation.
 - **Denoising AE**: robustness to noise.
 - **Contractive AE**: small derivatives of the representation w.r.t. input.

⇒ A regularized AE can be overcomplete and nonlinear but still learn something useful about the data distribution!

DENOISING AUTOENCODER (DAE)

- Idea: representation should be robust to introduction of noise.
- Produce corrupted version $\tilde{\mathbf{x}}$ of input \mathbf{x} , e.g. by
 - random assignment of subset of inputs to 0.
 - adding Gaussian noise.
- Modified reconstruction loss: $L(\mathbf{x}, \text{dec}(\text{enc}(\tilde{\mathbf{x}})))$
⇒ denoising AEs must learn to undo this corruption.

DAE-A PROBABILISTIC PERSPECTIVE

- With the corruption process, we induce stochasticity into the DAE.
- Formally: let $C(\tilde{\mathbf{x}}|\mathbf{x})$ present the conditional distribution of corrupted samples $\tilde{\mathbf{x}}$, given a data sample \mathbf{x} .
- Like feedforward NNs can model a distribution over targets $p(\mathbf{y}|\mathbf{x})$, output units and loss function of an AE can be chosen such that one gets a stochastic decoder $p_{decoder}(\mathbf{x}|\mathbf{z})$.
- E.g. linear output units to parametrize the mean of Gaussian distribution for real valued \mathbf{x} and negative log-likelihood loss (which is equal to MSE).
- The DAE then learns a reconstruction distribution $p_{reconstruct}(\mathbf{x}|\tilde{\mathbf{x}})$ from training pairs $(\mathbf{x}, \tilde{\mathbf{x}})$.
- (Note that the encoder could also be made stochastic, modelling $p_{encoder}(\mathbf{z}|\tilde{\mathbf{x}})$.)

DAE-A PROBABILISTIC PERSPECTIVE

The general structure of a DAE as a computational graph:

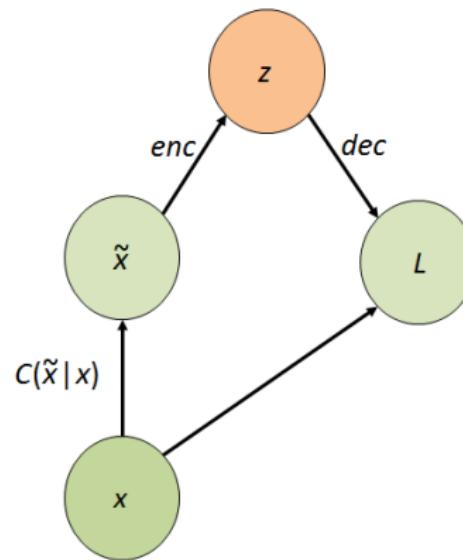


Figure: Denoising autoencoder: “making the learned representation robust to partial corruption of the input pattern.”

DAE-A PROBABILISTIC PERSPECTIVE

Algorithm Training denoising autoencoders

- 1: Sample a training example \mathbf{x} from the training data.
- 2: Sample a corrupted version $\tilde{\mathbf{x}}$ from $C(\tilde{\mathbf{x}}|\mathbf{x})$
- 3: Use $(\mathbf{x}, \tilde{\mathbf{x}})$ as a training example for estimating the AE reconstruction $p_{reconstruct}(\mathbf{x}|\tilde{\mathbf{x}}) = p_{decoder}(\mathbf{x}|\mathbf{z})$, where
 - \mathbf{z} is the output of the encoder $enc(\tilde{\mathbf{x}})$ and
 - $p_{decoder}$ defined by a decoder $dec(\mathbf{z})$

- All we have to do to transform an AE into a DAE is to add a stochastic corruption process on the input.
- The DAE still tries to preserve the information about the input (encode it), but also to undo the effect of a corruption process!
- Training a DAE by minimizing $\|dec(enc(\tilde{\mathbf{x}})) - \mathbf{x}\|^2$ corresponds to minimizing $\mathbb{E}_{\mathbf{x}, \tilde{\mathbf{x}} \sim p_{data}(\mathbf{x})C(\tilde{\mathbf{x}}|\mathbf{x})} [\log p_{decoder}(\mathbf{x}|enc(\tilde{\mathbf{x}}))]$.

DAE-A PROBABILISTIC PERSPECTIVE

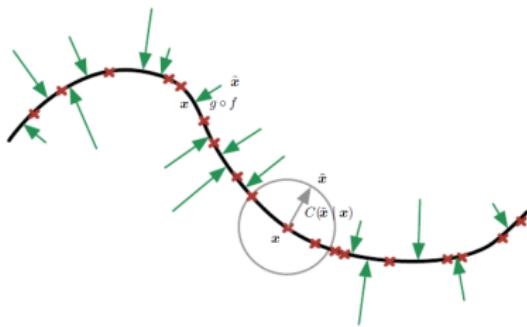


Figure: Denoising autoencoders - “manifold perspective” (Ian Goodfellow et al. (2016))

- A DAE is trained to map a corrupted data point $\tilde{\mathbf{x}}$ back to the original data point \mathbf{x} .
- The corruption process $C(\tilde{\mathbf{x}}|\mathbf{x})$ is displayed by the gray circle of equiprobable corruptions

DAE-A PROBABILISTIC PERSPECTIVE

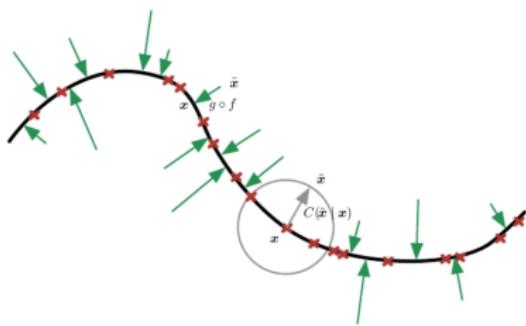


Figure: Denoising autoencoders - “manifold perspective” (Ian Goodfellow et al. (2016))

- The vector $dec(enc(\tilde{\mathbf{x}})) - \tilde{\mathbf{x}}$ points approximately towards the nearest point in the data manifold, since $dec(enc(\tilde{\mathbf{x}}))$ estimates the center of mass of clean points \mathbf{x} which could have given rise to $\tilde{\mathbf{x}}$.
- Thus, the DAE learns a vector field $dec(enc(\mathbf{x})) - \mathbf{x}$ indicated by the green arrows.

DAE-A PROBABILISTIC PERSPECTIVE

An example of a vector field learned by a DAE.

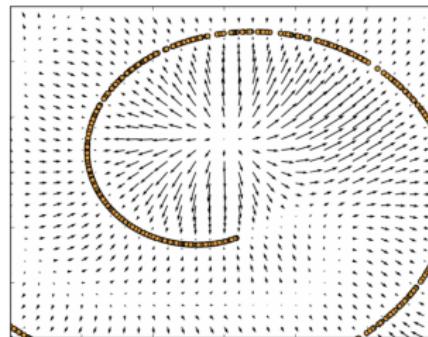


Figure: source : Ian Goodfellow et al. (2016)

EXPERIMENT: ENCODE MNIST WITH A DAE

- We will now corrupt the MNIST data with Gaussian noise and then try to denoise it as good as possible.

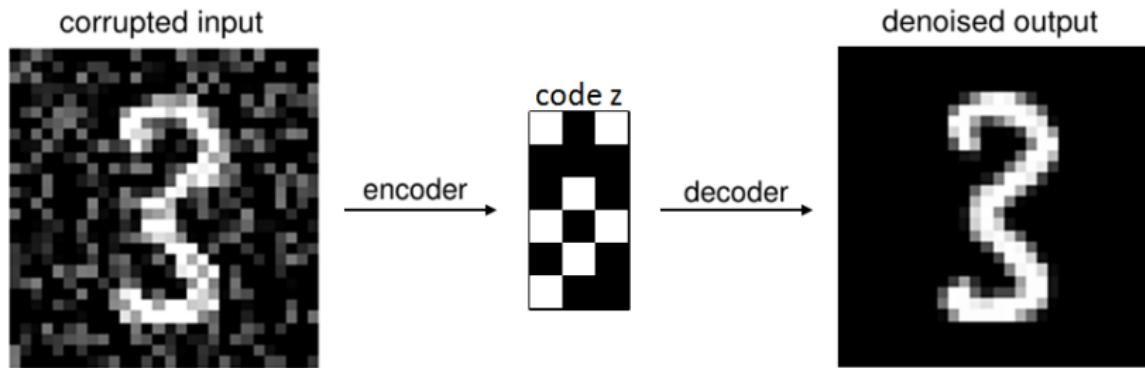


Figure: Flow chart of our autoencoder: denoise the corrupted input.

EXPERIMENT: ENCODE MNIST WITH A DAE

- To corrupt the input, we randomly add or subtract values from a uniform distribution to each of the image entries.

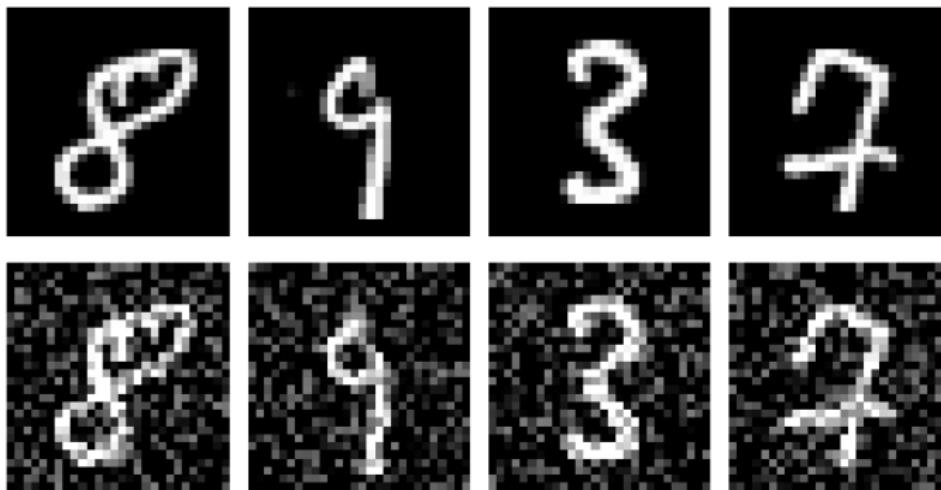


Figure: Top row: original data, bottom row: corrupted mnist data.

EXPERIMENT: ENCODE MNIST WITH A DAE

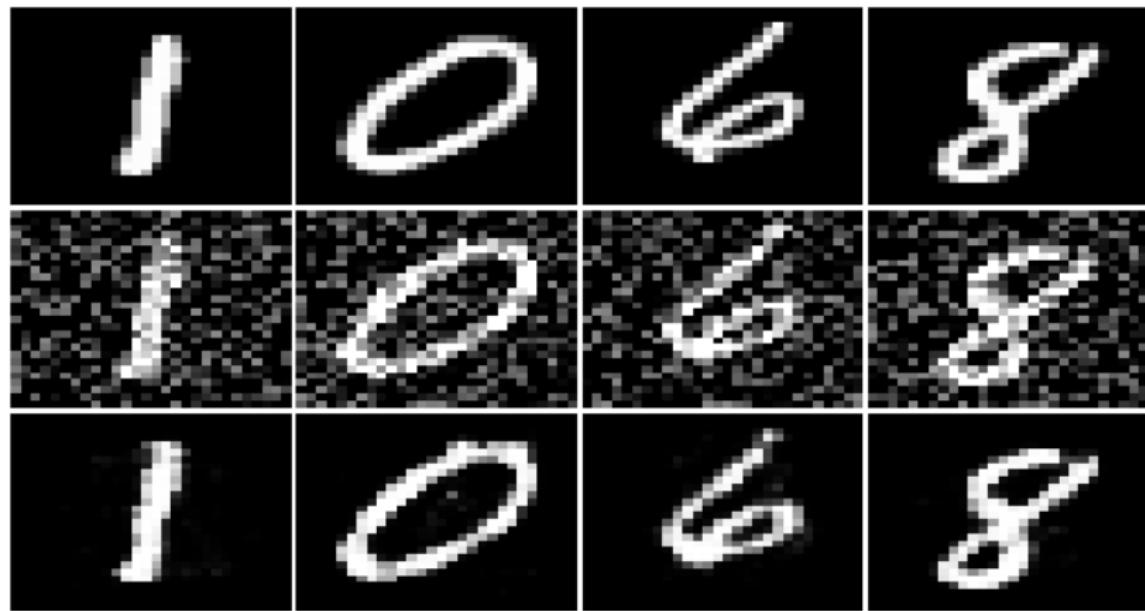


Figure: The top row shows the original digits, the intermediate one the corrupted and the bottom row the denoised/reconstructed digits (prediction).

- $\dim(\mathbf{z}) = 1568$ (overcomplete).

EXPERIMENT: ENCODE MNIST WITH A DAE

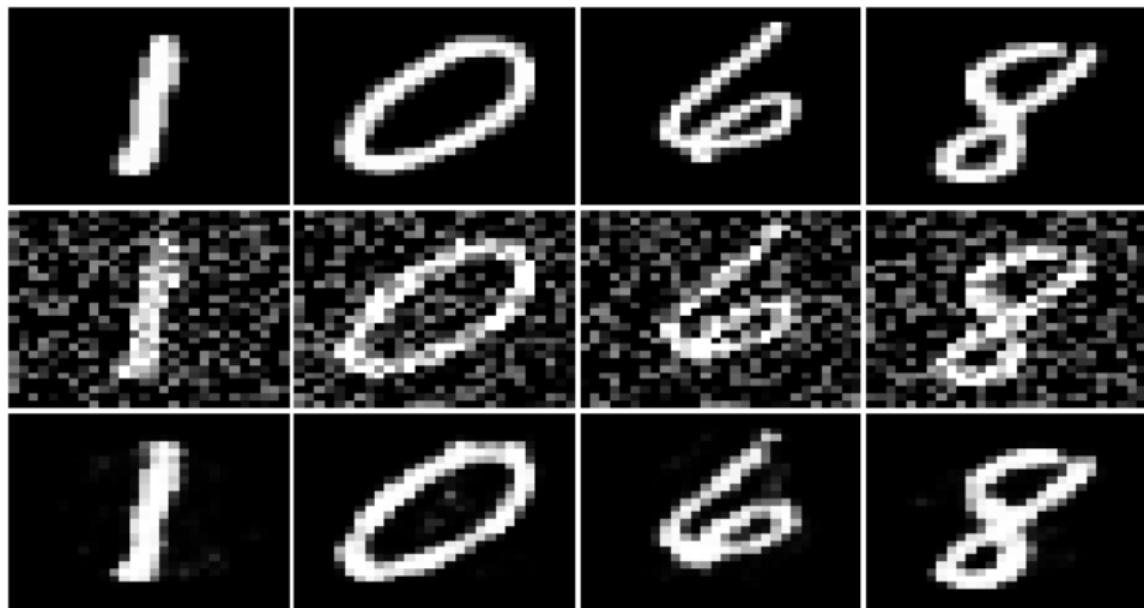


Figure: The top row shows the original digits, the intermediate one the corrupted and the bottom row the denoised/reconstructed digits (prediction).

- $\dim(\mathbf{z}) = 784 (= \dim(\mathbf{x}))$.

EXPERIMENT: ENCODE MNIST WITH A DAE

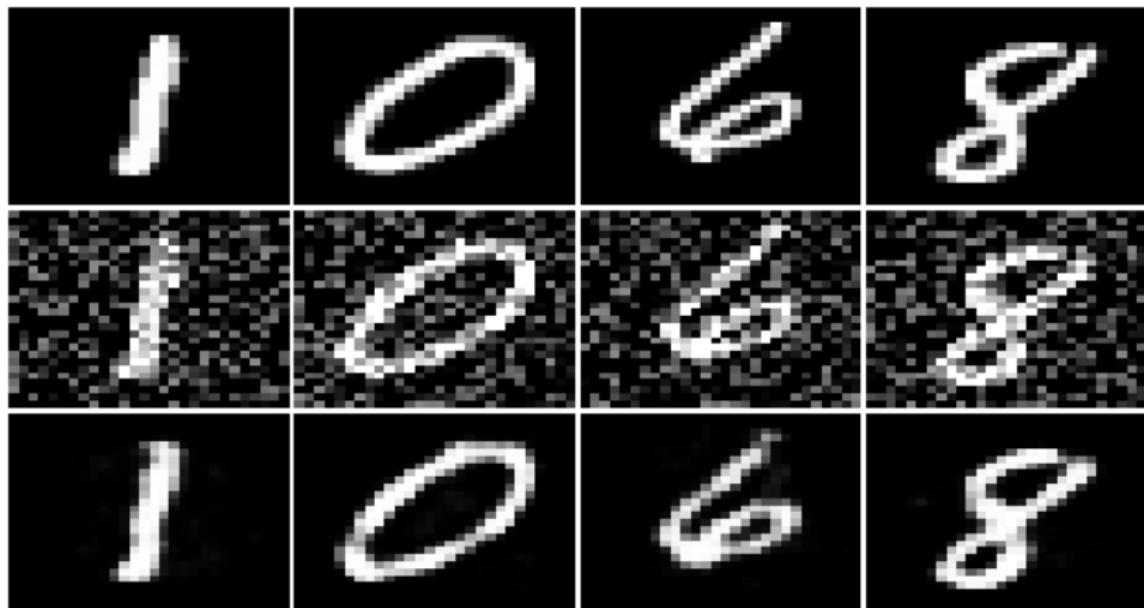


Figure: The top row shows the original digits, the intermediate one the corrupted and the bottom row the denoised/reconstructed digits (prediction).

- $\dim(\mathbf{z}) = 256$.

EXPERIMENT: ENCODE MNIST WITH A DAE

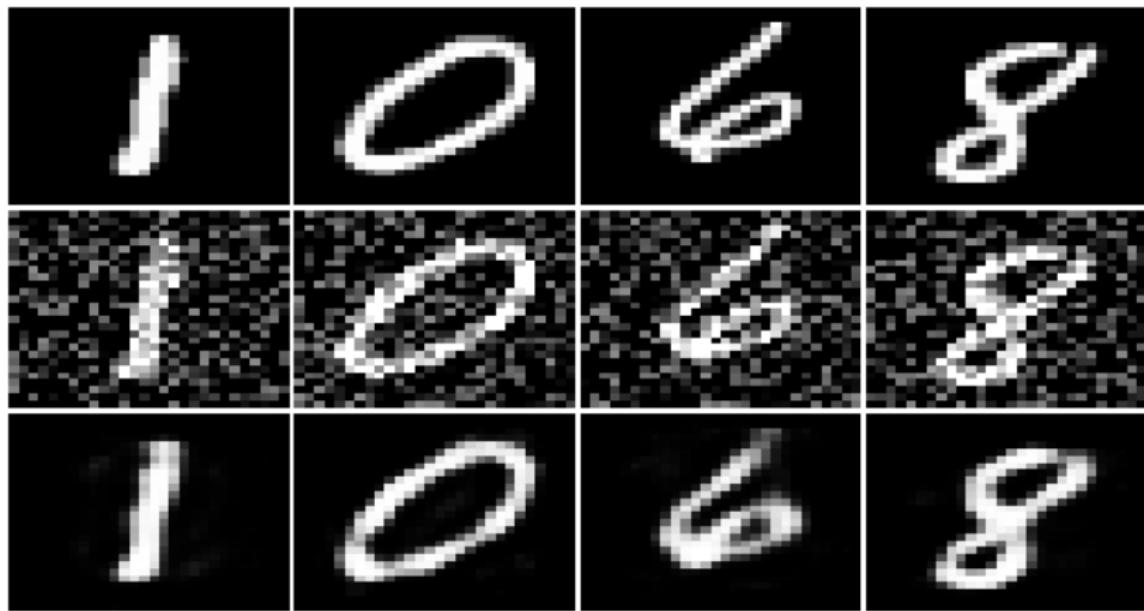


Figure: The top row shows the original digits, the intermediate one the corrupted and the bottom row the denoised/reconstructed digits (prediction).

- $\dim(\mathbf{z}) = 64$.

EXPERIMENT: ENCODE MNIST WITH A DAE

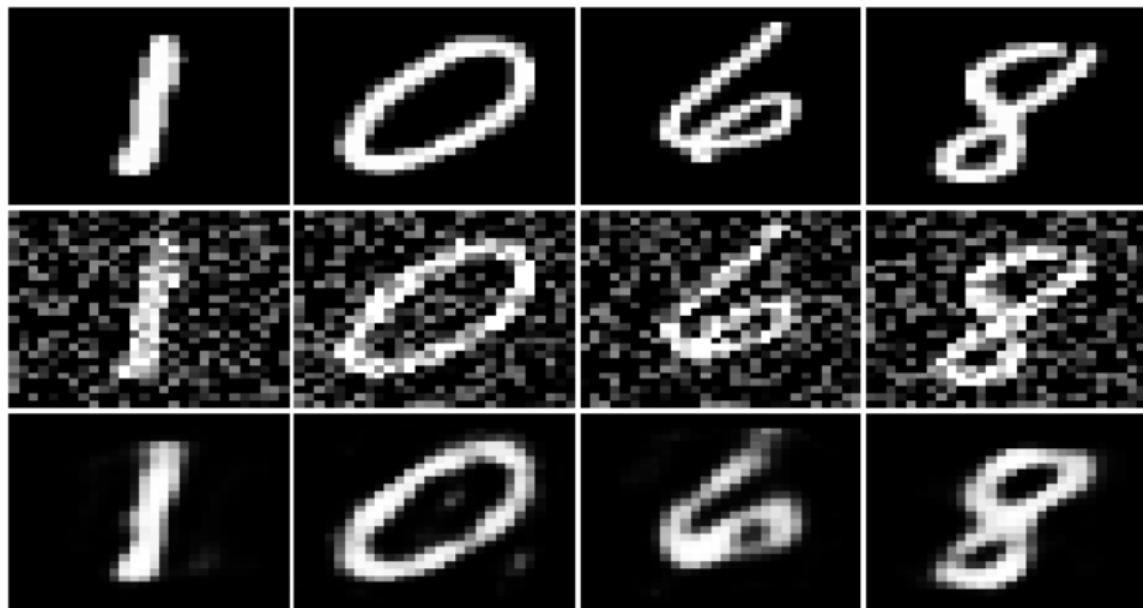


Figure: The top row shows the original digits, the intermediate one the corrupted and the bottom row the denoised/reconstructed digits (prediction).

- $\dim(\mathbf{z}) = 32$.

EXPERIMENT: ENCODE MNIST WITH A DAE

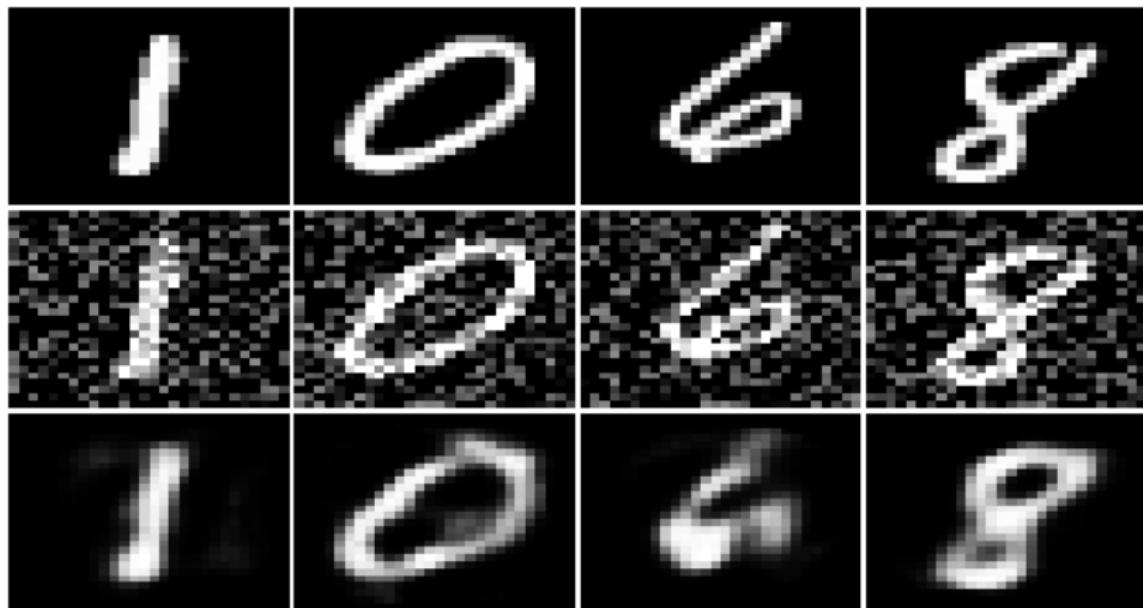


Figure: The top row shows the original digits, the intermediate one the corrupted and the bottom row the denoised/reconstructed digits (prediction).

- $\dim(\mathbf{z}) = 16$.

EXPERIMENT: ENCODE MNIST WITH A DAE

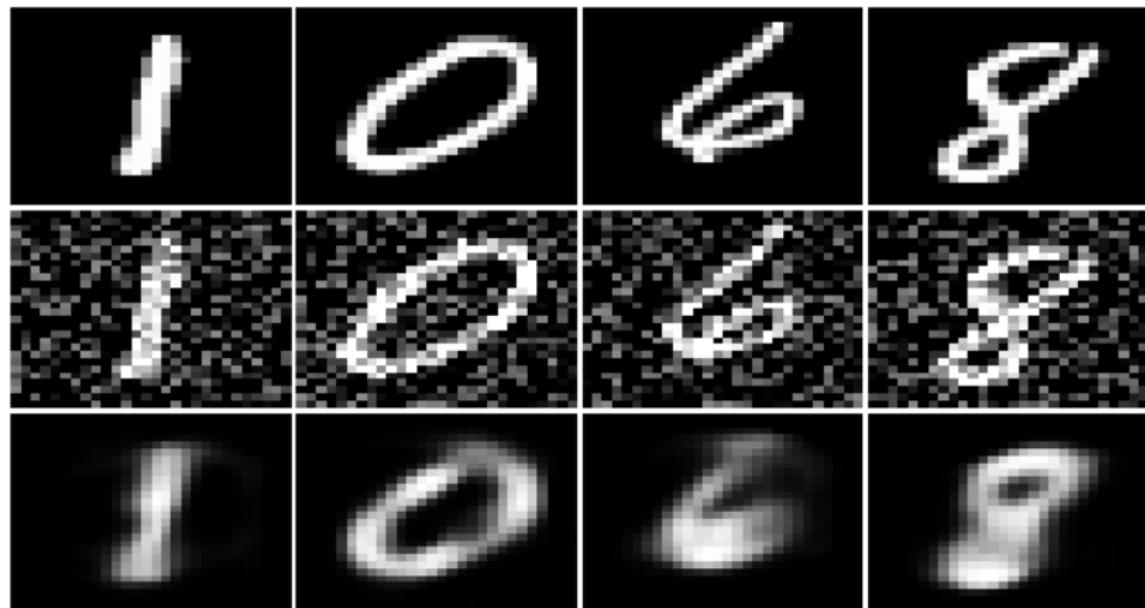


Figure: The top row shows the original digits, the intermediate one the corrupted and the bottom row the denoised/reconstructed digits (prediction).

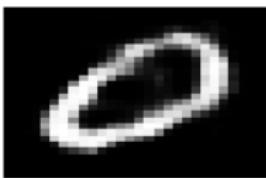
- $\dim(\mathbf{z}) = 8$.

EXPERIMENT: ENCODE MNIST WITH A DAE

- Let us increase the amount of noise and see how the autoencoder with $\dim(z) = 64$ deals with it (for science!).

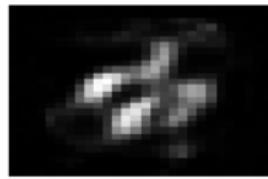
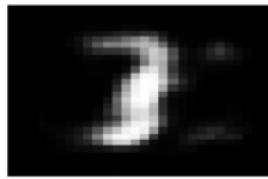
EXPERIMENT: ENCODE MNIST WITH A DAE

- A lot of noise.



EXPERIMENT: ENCODE MNIST WITH A DAE

- A lot of noise.



CONTRACTIVE AUTOENCODER (CAE)

- Idea: extract features that only reflect variations found in the training set.
- Add explicit regularization term to the reconstruction loss:

$$L(\mathbf{x}, dec(enc(\mathbf{x})) + \lambda \left\| \frac{\partial enc(\mathbf{x})}{\partial \mathbf{x}} \right\|_F^2$$

- ⇒ Derivatives of the encoder function w.r.t. the input are encouraged to be small.
- ⇒ Only a small number of input directions will have significant derivatives.
- ⇒ The encoder function is encouraged to resist infinitesimal perturbations of the input.

DAE VS. CAE

DAE

the *decoder* function is trained to resist infinitesimal perturbations of the input.

CAE

the *encoder* function is trained to resist infinitesimal perturbations of the input.

WHICH AUTOENCODER?

- Both the denoising and contractive autoencoders perform well.
- Advantage of denoising autoencoder: simpler to implement
 - requires adding one or two lines of code to regular AE.
 - no need to compute Jacobian of hidden layer.
- Advantage of contractive autoencoder: gradient is deterministic
 - can use second order optimizers (conjugate gradient, LBFGS, etc.).
 - might be more stable than the denoising autoencoder, which uses a sampled gradient.

Manifold learning

MANIFOLD LEARNING

- **Manifold hypothesis:** Data of interest lies on an embedded non-linear manifold within the higher-dimensional space.
- **A manifold:**
 - is a topological space that locally resembles the Euclidean space.
 - in ML more loosely refers to a connected set of points that can be approximated well by considering only a small number of dimensions.

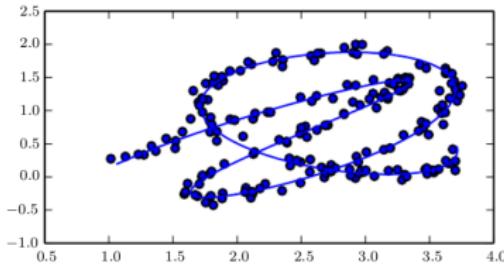


Figure: from Goodfellow et. al

MANIFOLD LEARNING

- An important characterization of a manifold is the set of its tangent planes.
- **Definition:** At a point \mathbf{x} on a d -dimensional manifold, the **tangent plane** is given by d basis vectors that span the local directions of variation allowed on the manifold.

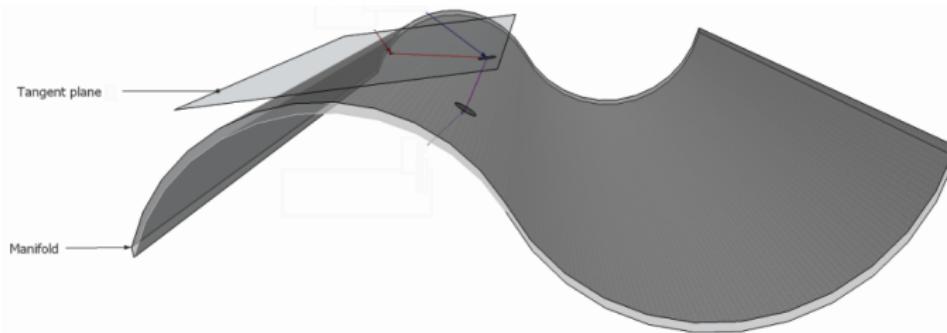


Figure: A pictorial representation of the tangent space of a single point, \mathbf{x} , on a manifold (Goodfellow et al. (2016)).

MANIFOLD LEARNING

- Manifold hypothesis does not need to hold true.
- In the context of AI tasks (e.g. processing images, sound, or text) it seems to be at least approximately correct, since :
 - probability distributions over images, text strings, and sounds that occur in real life are highly concentrated (randomly sampled pixel values do not look like images, randomly sampling letters is unlikely to result in a meaningful sentence).
 - samples are connected to each other by other samples, with each sample surrounded by other highly similar samples that can be reached by applying transformations (E.g. for images: Dim or brighten the lights, move or rotate objects, change the colors of objects, etc).

LEARNING MANIFOLDS WITH AES

- AEs training procedures involve a compromise between two forces:
 - ➊ Learning a representation \mathbf{z} of a training example \mathbf{x} such that \mathbf{x} can be approximately recovered from \mathbf{z} through a decoder.
 - ➋ Satisfying an architectural constraint or regularization penalty.
- Together, they force the hidden units to capture information about the structure of the data generating distribution
- important principle: AEs can afford to represent only the variations that are needed to reconstruct training examples.
- If the data-generating distribution concentrates near a low-dimensional manifold, this yields representations that implicitly capture a local coordinate system for the manifold.

LEARNING MANIFOLDS WITH AES

- Only the variations tangent to the manifold around \mathbf{x} need to correspond to changes in $\mathbf{z} = \text{enc}(\mathbf{x})$. Hence the encoder learns a mapping from the input space to a representation space that is only sensitive to changes along the manifold directions, but that is insensitive to changes orthogonal to the manifold.

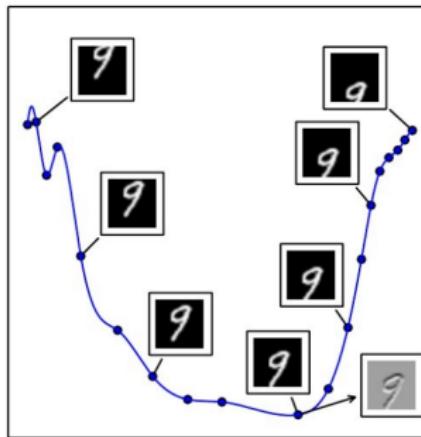


Figure: from Goodfellow et al. (2016)

LEARNING MANIFOLDS WITH AES

- Common setting: a representation (embedding) for the points on the manifold is learned.
- Two different approaches
 - ① Non-parametric methods: learn an embedding for each training example.
 - ② Learning a more general mapping for *any* point in the input space.
- AI problems can have very complicated structures that can be difficult to capture from only local interpolation.
⇒ Motivates use of **distributed representations** and deep learning!

Specific AEs and applications

CONVOLUTIONAL AUTOENCODER (CONVAE)

- For the image domain, using convolutions is advantageous. Can we also make use of them in AEs?
- In a ConvAE, the encoder consists of convolutional layers. The decoder, on the other hand, consists of transpose convolution layers or simple upsampling operations.

CONVOLUTIONAL AUTOENCODER (CONVAE)

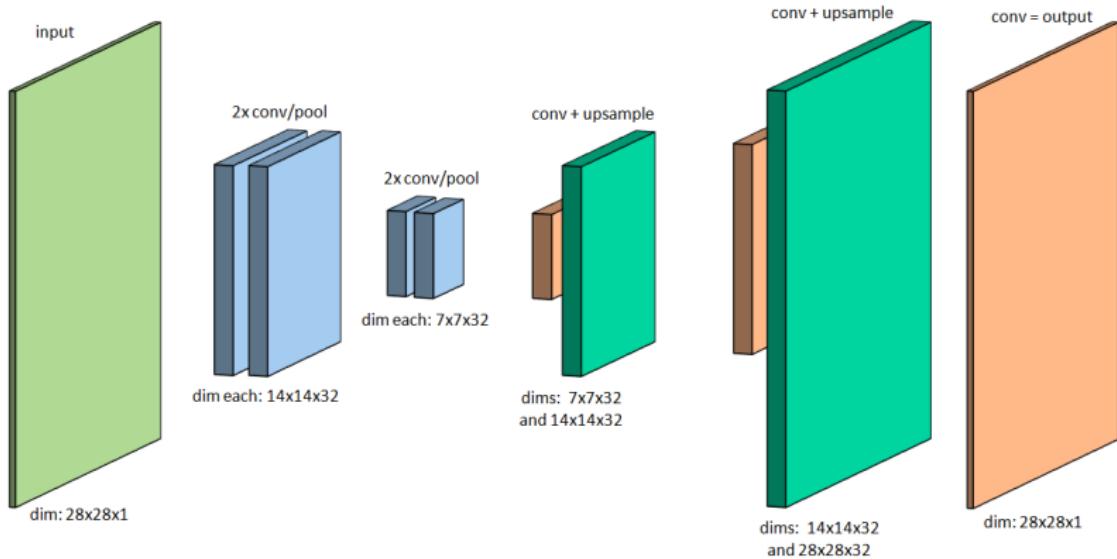


Figure: Potential architecture of a convolutional autoencoder.

We now apply this architecture to denoise MNIST.

CONVOLUTIONAL AUTOENCODER (CONVAE)

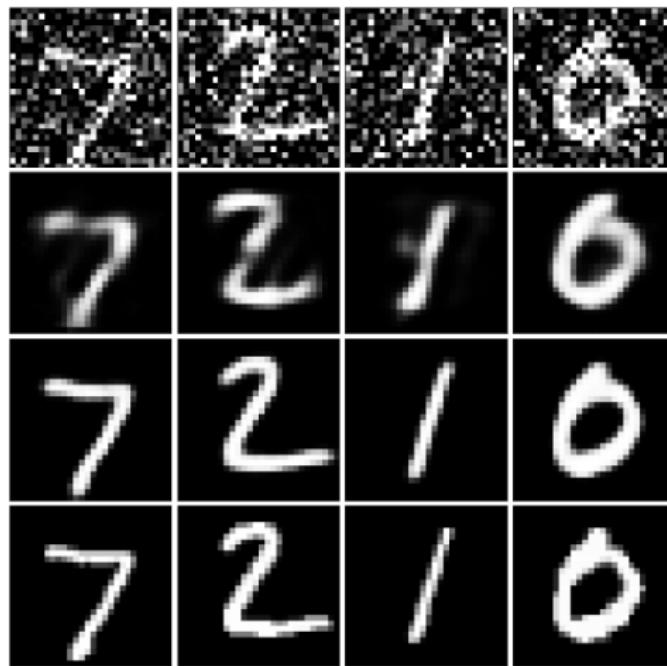


Figure: Top row: noised data, second row: AE with $\dim(\mathbf{z}) = 32$ (roughly 50k params), third row: ConvAE (roughly 25k params), fourth row: ground truth.

AES FOR UNSUPERVISED PRETRAINING

- Stacked AEs can be used for layer-wise unsupervised pretraining of deep neural networks.
- This corresponds to subsequently training each layer as an AE.
- It aims at yielding better weight initializations for the actual supervised training.
- This usually eliminates the risk of vanishing gradients in feed forward nets.
- It played an important role in the past before general techniques for stabilizing optimization were invented (e.g. ReLUs, batch normalization, dropout, etc.)

REAL-WORLD APPLICATIONS

Today, autoencoders are still used for tasks such as:

- data de-noising,
- compression,
- and dimensionality reduction for the purpose of visualization.

REAL-WORLD APPLICATIONS

Medical image denoising using convolutional denoising autoencoders

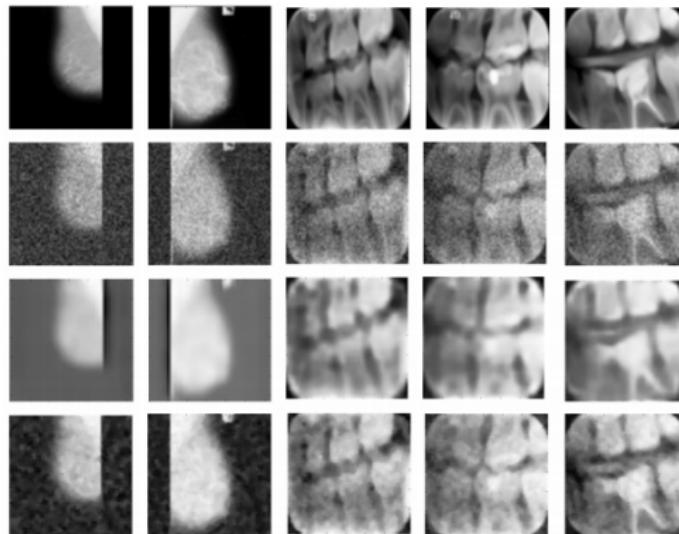


Figure: Top row : real image, second row : noisy version, third row : results of a (convolutional) denoising autoencoder and fourth row : results of a median filter (Lovedeep Gondara (2016))

REAL-WORLD APPLICATIONS

AE-based image compression.

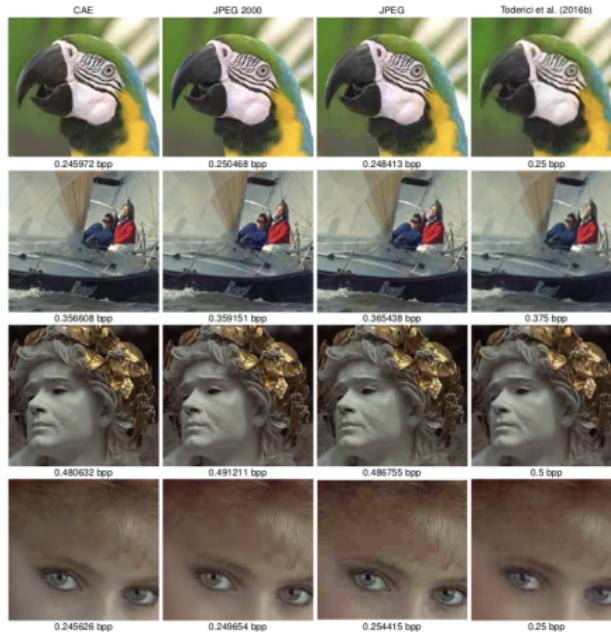


Figure: from Theis et al.

REFERENCES

-  Ian Goodfellow, Yoshua Bengio and Aaron Courville (2016)
Deep Learning
<http://www.deeplearningbook.org/>
-  Lovedeep Gondara (2016)
Medical image denoising using convolutional denoising autoencoders
<https://arxiv.org/abs/1608.04667>
-  Vijay Badrinarayanan, Alex Kendall and Roberto Cipolla (2016)
SegNet: A Deep Convolutional Encoder-Decoder Architecture for Image Segmentation
<https://arxiv.org/abs/1511.00561>
-  Theis, Lucas; Shi, Wenzhe; Cunningham, Andrew; Huszár, Ferenc (2017)
Lossy Image Compression with Compressive Autoencoders
<https://arxiv.org/abs/1703.00395>