

Deep Learning

17 Autoencoders

Dr. Konda Reddy Mopuri
Dept. of AI, IIT Hyderabad
Jan-May 2023

Beyond Classification and Regression

- ① Applications such as image synthesis, image-to-image transformations model high-dim signals

Beyond Classification and Regression

- ① Applications such as image synthesis, image-to-image transformations model high-dim signals
- ② These applications require to learn the meaningful degrees of freedom that constitute the signal

- ① Applications such as image synthesis, image-to-image transformations model high-dim signals
- ② These applications require to learn the meaningful degrees of freedom that constitute the signal
- ③ Typically, these degrees of freedom are of lesser dimensions than the signal

Example: Synthesizing Human faces

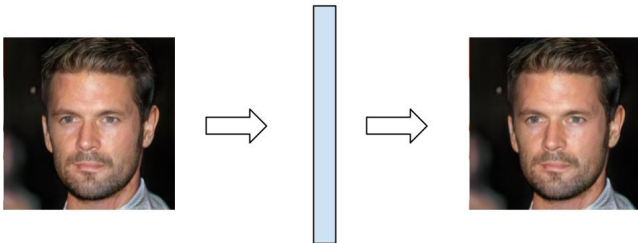
- ① For generating new faces, it makes sense to capture a small number of degrees of freedom such as
 - skull size and shape
 - color of skin and eyes
 - features of nose and lips, etc.

Example: Synthesizing Human faces

- ① For generating new faces, it makes sense to capture a small number of degrees of freedom such as
 - skull size and shape
 - color of skin and eyes
 - features of nose and lips, etc.
- ② Even a comprehensive list of such things will be less than the number of pixels in the image (i.e. resolution)

Example: Synthesizing Human faces

- ① If we can model these relatively small number of dimensions, we can synthesize a face with thousands of dimensions



Autoencoder: architecture

- ① Feed-forward Neural network that maps a space to itself

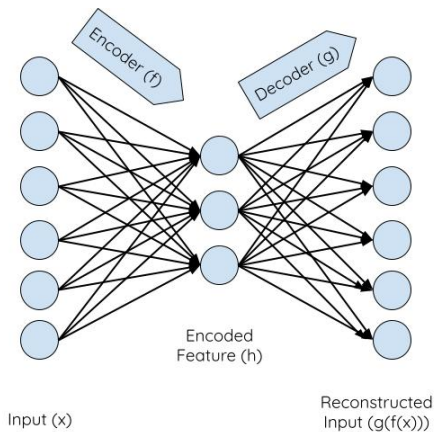
Autoencoder: architecture

- ① Feed-forward Neural network that maps a space to itself
- ② Trained to map its input to itself (but not an identity function)

Autoencoder: architecture

- ① Feed-forward Neural network that maps a space to itself
- ② Trained to map its input to itself (but not an identity function)
- ③ Network consists of two parts: encoder (f) and decoder (g)

Autoencoder: architecture



Autoencoder: principle

- ① Original (input) space is of higher dimensions but the data lies in a manifold of smaller dimension

Autoencoder: principle

- ① Original (input) space is of higher dimensions but the data lies in a manifold of smaller dimension
- ② Dimension of the latent space is a hyper-parameter chosen from prior knowledge, and/or through heuristics

Autoencoder: principle

- ① Original (input) space is of higher dimensions but the data lies in a manifold of smaller dimension
- ② Dimension of the latent space is a hyper-parameter chosen from prior knowledge, and/or through heuristics

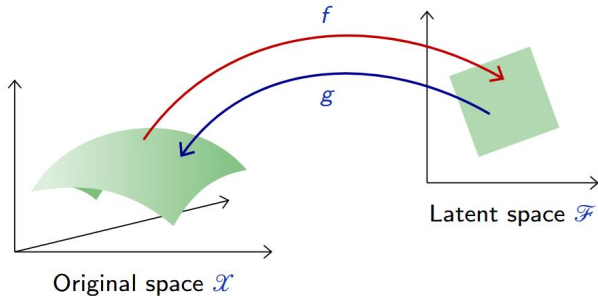


Figure credits: Francois Fluoret

Autoencoder: nonlinearity

- ① For a binary i/p vector, what could be an appropriate nonlinearity for g ?

Autoencoder: nonlinearity

- ① For a binary i/p vector, what could be an appropriate nonlinearity for g ?
- ② For a real i/p vector?

Autoencoder: nonlinearity

- ① For a binary i/p vector, what could be an appropriate nonlinearity for g ?
- ② For a real i/p vector?
- ③ Nonlinearity for f ?

Autoencoder: Objective

- ① Enforces the reconstructed o/p to be very similar to i/p

Autoencoder: Objective

- ① Enforces the reconstructed o/p to be very similar to i/p
- ② Loss function takes care of this via training

Autoencoder: Objective

- ① Let p be the data distribution in the input space, autoencoder is characterized with the following loss

$$\mathbb{E}_{x \sim p} \|x - g \circ f(x)\|^2 \approx 0$$

Autoencoder: Objective

- ① Let p be the data distribution in the input space, autoencoder is characterized with the following loss

$$\mathbb{E}_{x \sim p} \|x - g \circ f(x)\|^2 \approx 0$$

- ② Training: finding the parameters for the encoder ($f(\cdot; w_f)$) and decoder ($g(\cdot; w_g)$) optimizing the empirical loss

$$\hat{w}_f, \hat{w}_g = \operatorname{argmin}_{w_f, w_g} \frac{1}{N} \sum_n \|x_n - g(f(x_n; w_f); w_g)\|^2$$

Autoencoder: Objective

- ① For binary i/p , we may interpret the reconstructions as probabilities (with a sigmoid nonlinearity)

Autoencoder: Objective

- ① For binary i/p , we may interpret the reconstructions as probabilities (with a sigmoid nonlinearity)
- ② Hence, we may use BCE loss for training

Autoencoder: Connection to PCA

- ① f and g are linear functions (data is normalized $x_i = \frac{1}{\sqrt{|X|}}(x_i - \mu)$)
→ optimal solution is PCA

Autoencoder: Connection to PCA

- ① f and g are linear functions (data is normalized $x_i = \frac{1}{\sqrt{|X|}}(x_i - \mu)$)
→ optimal solution is PCA
- ② Better results can be made possible with sophisticated transformations such as deep neural networks → Deep Autoencoders

Deep Autoencoders



Top row: original data samples

Bottom row: corresponding reconstructed samples (single ReLU layer of dimension 32)

Figure credits: [Keras blog](#)

Autoencoder: Regularization

- ① idea is to help generalization

Autoencoder: Regularization

- ① idea is to help generalization
- ② Over-complete autoencoders (more hidden units than the i/p dimension) \rightarrow parameter explosion and prone to overfitting

Autoencoder: Regularization

- ① idea is to help generalization
- ② Over-complete autoencoders (more hidden units than the i/p dimension) \rightarrow parameter explosion and prone to overfitting
- ③ Even the under-complete configurations benefit from regularization

Autoencoder: Regularization

- ① idea is to help generalization
- ② Over-complete autoencoders (more hidden units than the i/p dimension) → parameter explosion and prone to overfitting
- ③ Even the under-complete configurations benefit from regularization
- ④ Simplest is to add l_2 regularization term to the objective

Autoencoder: Regularization

- ① idea is to help generalization
- ② Over-complete autoencoders (more hidden units than the i/p dimension) \rightarrow parameter explosion and prone to overfitting
- ③ Even the under-complete configurations benefit from regularization
- ④ Simplest is to add l_2 regularization term to the objective
- ⑤ Tie the weights, i.e., $w_g = w_f^T$

Besides dimensionality reduction

- ① Autoencoders can capture the dependencies across signal components

Besides dimensionality reduction

- ① Autoencoders can capture the dependencies across signal components
- ② This can help to restore the missing components from an input

Besides dimensionality reduction

- ① In this scenario, we may ignore the encoder/decoder architecture

Besides dimensionality reduction

- ① In this scenario, we may ignore the encoder/decoder architecture
- ② Goal in this case is not to learn a ϕ such that $\phi(X) \approx X$

Besides dimensionality reduction

- ① In this scenario, we may ignore the encoder/decoder architecture
- ② Goal in this case is not to learn a ϕ such that $\phi(X) \approx X$
- ③ It is to learn a ϕ such that $\phi(\tilde{X}) \approx X$, where \tilde{X} is a perturbed version of X

Besides dimensionality reduction

- ① In this scenario, we may ignore the encoder/decoder architecture
- ② Goal in this case is not to learn a ϕ such that $\phi(X) \approx X$
- ③ It is to learn a ϕ such that $\phi(\tilde{X}) \approx X$, where \tilde{X} is a perturbed version of X
- ④ This is referred to as a **Denoising Autoencoder**

- ① This can be illustrated with an additive Gaussian noise in case of a 2D signal and MSE

$$\hat{w} = \underset{w}{\operatorname{argmin}} \frac{1}{N} \sum_{n=1}^N \|x_n - \phi(x_n + \epsilon_n; w)\|^2,$$

where x_n are data samples and ϵ_n are Gaussian random noise

Denoising Autoencoder

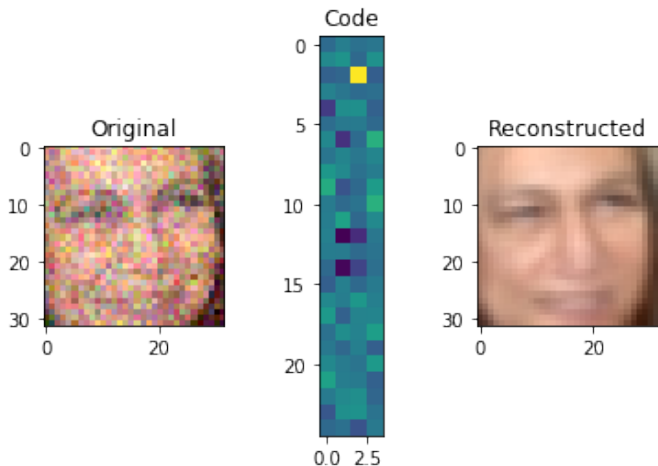


Figure credits: Ali Abdelal, <https://stackabuse.com/>

Sparse Autoencoders

- ① Tries to enforce the hidden neurons to be inactive mostly

Sparse Autoencoders

- ① Tries to enforce the hidden neurons to be inactive mostly
- ② Restricts the freedom of the parameters by forcing them to fire sparsely

Sparse Autoencoders

- ① Tries to enforce the hidden neurons to be inactive mostly
- ② Restricts the freedom of the parameters by forcing them to fire sparsely
- ③ Uses a sparsity parameter (ρ) (typically close to 0, say 0.01)
- ④ Enforces the mean neuron activation ($\hat{\rho}_l$) to be close to ρ

Sparse Autoencoders

① Mean activation: $\hat{\rho}_l = \frac{1}{m} \sum_{i=1}^m f(x_i)_l$

Sparse Autoencoders

- ① Mean activation: $\hat{\rho}_l = \frac{1}{m} \sum_{i=1}^m f(x_i)_l$
- ② $R(w) = \sum_{l=1}^k \rho \log \frac{\rho}{\hat{\rho}_l} + (1 - \rho) \log \frac{1-\rho}{1-\hat{\rho}_l}$

Sparse Autoencoders

- ① Mean activation: $\hat{\rho}_l = \frac{1}{m} \sum_{i=1}^m f(x_i)_l$
- ② $R(w) = \sum_{l=1}^k \rho \log \frac{\rho}{\hat{\rho}_l} + (1 - \rho) \log \frac{1-\rho}{1-\hat{\rho}_l}$
- ③ k - dimension of hidden layer
m -size of training dataset

Contractive Autoencoders

- ① Prevents an autoencoder from learning an identity function

Contractive Autoencoders

- ① Prevents an autoencoder from learning an identity function
- ② $R(w) = \left| \frac{\partial f}{\partial x} \right|_F$

Contractive Autoencoders

- ① Prevents an autoencoder from learning an identity function
- ② $R(w) = \left| \frac{\partial f}{\partial x} \right|_F$
- ③ Competition (in the latent/hidden layer) b/w 'being sensitive' and 'not sensitive' to the i/p variations

- ① Prevents an autoencoder from learning an identity function
- ② $R(w) = \left| \frac{\partial f}{\partial x} \right|_F$
- ③ Competition (in the latent/hidden layer) b/w 'being sensitive' and 'not sensitive' to the i/p variations
- ④ Ends up capturing only the important variations in the i/p (something like PCA)

Latent Representations

- 1 Consider two samples in the latent space and reconstruct the samples along the line joining these

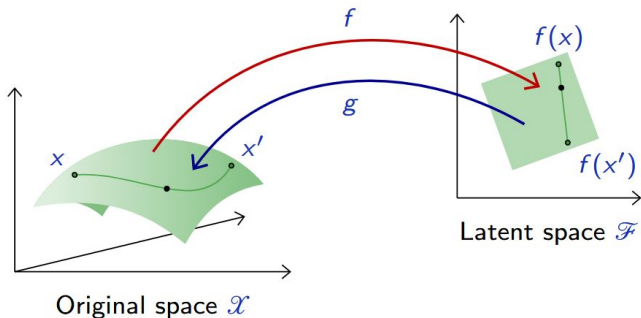


Figure credits: Francois Fleuret

Latent Representations

- 1 Consider two samples in the latent space and reconstruct the samples along the line joining these
- 2 $g(\alpha x + (1 - \alpha)x')$

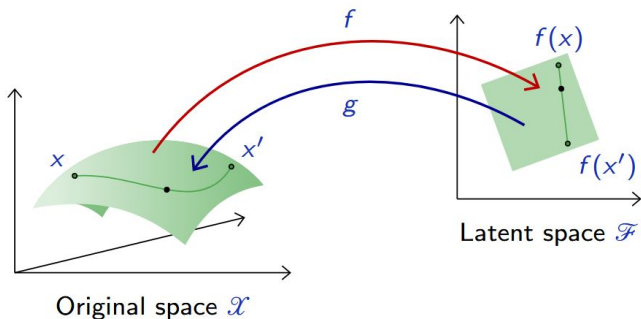
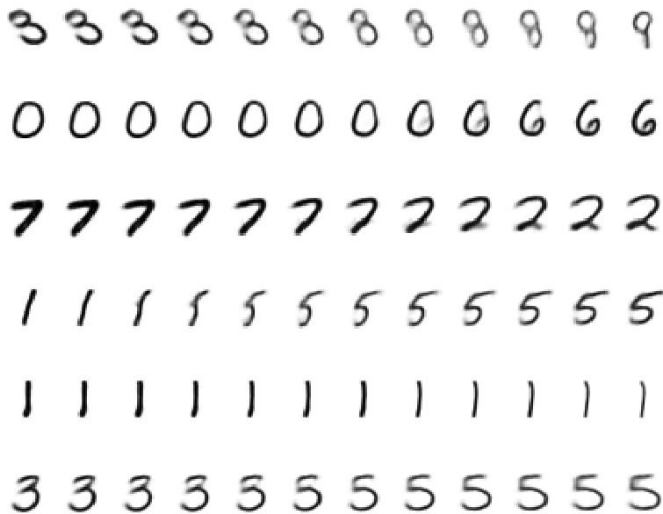


Figure credits: Francois Fleuret

Latent Representations





- 1 Introduce a density model over the latent space

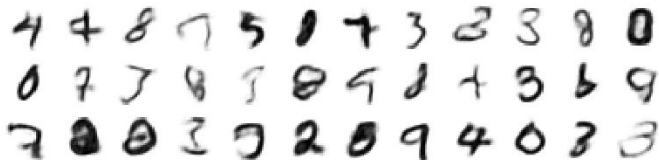
- ① Introduce a density model over the latent space
- ② Sample there and reconstruct using the decoder g

- ① Introduce a density model over the latent space
- ② Sample there and reconstruct using the decoder g
- ③ For instance, use a Gaussian density for modeling the latent space from the training data (estimate mean and a diagonal covariance matrix)

Generative Modeling by Autoencoder



Autoencoder sampling ($d = 8$)



Autoencoder sampling ($d = 16$)

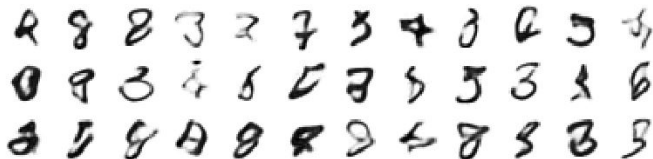


Figure credits: Francois Fleuret



- ① Reconstructions are not convincing

- ① Reconstructions are not convincing
- ② Because the density model is too simple



- ① Reconstructions are not convincing
- ② Because the density model is too simple
- ③ Good model still needs to capture the empirical distribution on the data although in a lower dimensional space