# CSCE 496/896 Lecture 5: Autoencoders

Stephen Scott

(Adapted from Paul Quint and Ian Goodfellow)
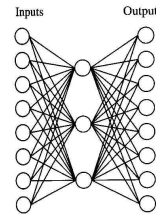
sscott@cse.unl.edu

---

# Introduction

- **Autoencoding** is training a network to replicate its input to its output
- Applications:
  - Unlabeled pre-training for semi-supervised learning
  - Learning **embeddings** to support information retrieval
  - Generation of new instances similar to those in the training set
  - Data compression

---

# Outline

- Basic idea
- Stacking
- Types of autoencoders
  - Denoising
  - Sparse
  - Contractive
  - Variational
  - Generative adversarial networks
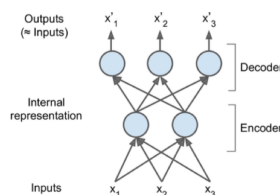
---

# Basic Idea

| Input | | Hidden Values | | | | Output |
|---|---|---|---|---|---|---|
| 10000000 | → | .89 | .04 | .08 | → | 10000000 |
| 01000000 | → | .15 | .99 | .99 | → | 01000000 |
| 00100000 | → | .01 | .97 | .27 | → | 00100000 |
| 00010000 | → | .99 | .97 | .71 | → | 00010000 |
| 00001000 | → | .03 | .05 | .02 | → | 00001000 |
| 00000100 | → | .01 | .11 | .88 | → | 00000100 |
| 00000010 | → | .80 | .01 | .98 | → | 00000010 |
| 00000001 | → | .60 | .94 | .01 | → | 00000001 |

Inputs   Outputs

- Sigmoid activation functions, 5000 training epochs, square loss, no regularization
- What's special about the hidden layer outputs?

---

# Basic Idea

- An **autoencoder** is a network trained to learn the **identity function:** output = input

Outputs (≈ Inputs)   $x'_1$   $x'_2$   $x'_3$

Decoder

Internal representation

Encoder

Inputs   $x_1$   $x_2$   $x_3$

- Subnetwork called **encoder** $f(\cdot)$ maps input to an **embedded representation**
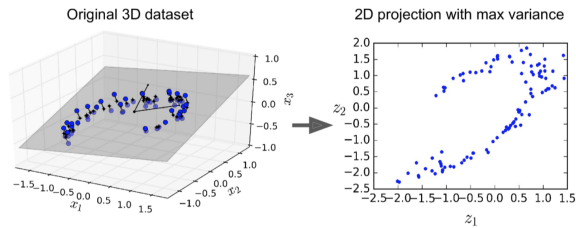- Subnetwork called **decoder** $g(\cdot)$ maps back to input space

- Can be thought of as **lossy compression** of input
- Need to identify the important attributes of inputs to reproduce faithfully

---

# Basic Idea

- General types of autoencoders based on size of hidden layer
  - **Undercomplete** autoencoders have hidden layer size smaller than input layer size
    - ⇒ Dimension of embedded space lower than that of input space
    - ⇒ Cannot simply memorize training instances
  - **Overcomplete** autoencoders have much larger hidden layer sizes
    - ⇒ Regularize to avoid overfitting, e.g., enforce a **sparsity** constraint
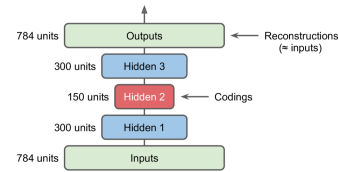
# Basic Idea
## Example: Principal Component Analysis

Original 3D dataset

2D projection with max variance

- A 3-2-3 autoencoder with linear units and square loss performs **principal component analysis**: Find linear transformation of data to maximize variance

---

# Stacked Autoencoders

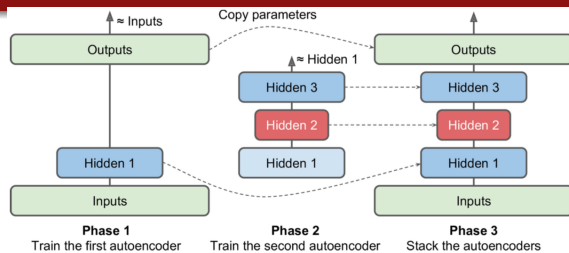| | |
|---|---|
| 784 units — Outputs ← Reconstructions (≈ inputs) | |
| 300 units — Hidden 3 | |
| 150 units — Hidden 2 ← Codings | |
| 300 units — Hidden 1 | |
| 784 units — Inputs | |

- A **stacked autoencoder** has multiple hidden layers

- Can share parameters to reduce their number by exploiting symmetry: $W_4 = W_1^\top$ and $W_3 = W_2^\top$

```
weights1 = tf.Variable(weights1_init, dtype=tf.float32, name="weights1")
weights2 = tf.Variable(weights2_init, dtype=tf.float32, name="weights2")
weights3 = tf.transpose(weights2, name="weights3")    # shared weights
weights4 = tf.transpose(weights1, name="weights4")    # shared weights
```
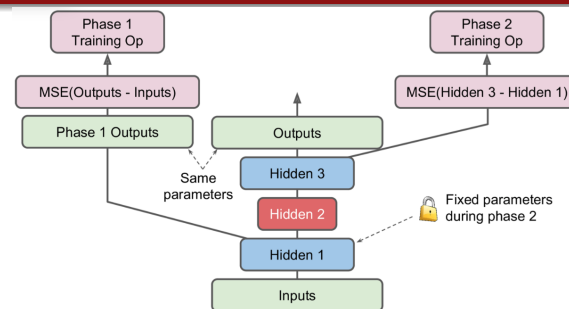
---

# Stacked Autoencoders
## Incremental Training

**Phase 1**
Train the first autoencoder

**Phase 2**
Train the second autoencoder

**Phase 3**
Stack the autoencoders

- Can simplify training by starting with single hidden layer $H_1$
- Then, train a second AE to mimic the output of $H_1$
- Insert this into first network
- Can build by using $H_1$'s output as training set for Phase 2

---

# Stacked Autoencoders
## Incremental Training (Single TF Graph)

Phase 1 Training Op

Phase 2 Training Op

MSE(Outputs - Inputs)

MSE(Hidden 3 - Hidden 1)

Phase 1 Outputs

Outputs

Same parameters

Hidden 3

Hidden 2 — Fixed parameters during phase 2
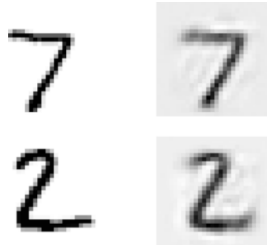
Hidden 1

Inputs

- Previous approach requires multiple TensorFlow graphs
- Can instead train both phases in a single graph: First left side, then right

---

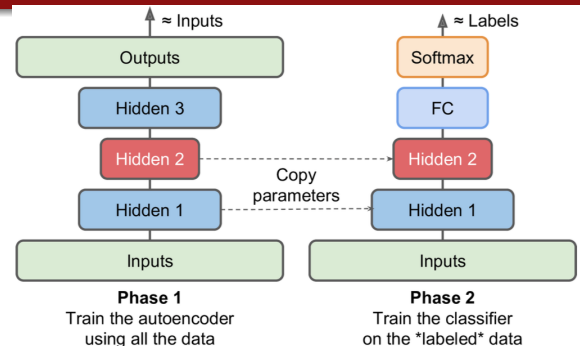# Stacked Autoencoders
## Visualization

Input MNIST Digit

Network Output

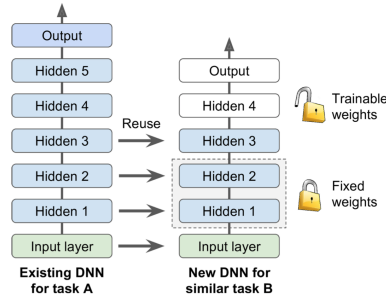Weights (features selected) for five nodes from $H_1$:

---

# Stacked Autoencoders
## Semi-Supervised Learning

≈ Inputs

≈ Labels

Outputs

Softmax

Hidden 3

FC

Hidden 2 — Copy parameters → Hidden 2

Hidden 1

Hidden 1

Inputs

Inputs

**Phase 1**
Train the autoencoder using all the data

**Phase 2**
Train the classifier on the *labeled* data

- Can **pre-train** network with unlabeled data
- ⇒ learn useful features and then train "logic" of dense layer with labeled data

## Slide 1: Transfer Learning from Trained Classifier

### Transfer Learning from Trained Classifier

- Can also transfer from a classifier trained on different task, e.g., transfer a GoogleNet architecture to ultrasound classification
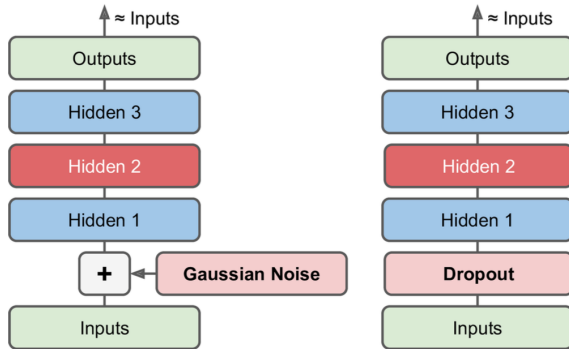


- Often choose existing one from a **model zoo**

## Slide 2: Denoising Autoencoders
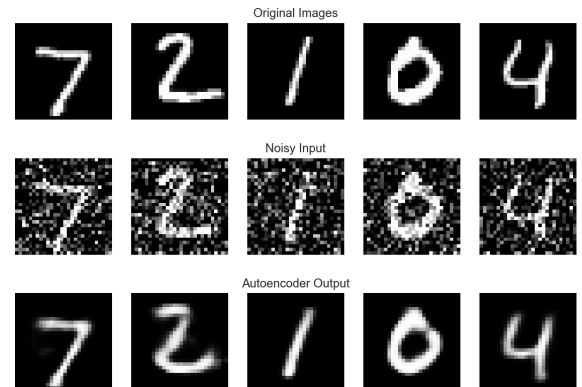
### Denoising Autoencoders
Vincent et al. (2010)

- Can train an autoencoder to learn to **denoise** input by giving input **corrupted** instance $\tilde{x}$ and targeting **uncorrupted** instance $x$
- Example noise models:
    - **Gaussian noise:** $\tilde{x} = x + z$, where $z \sim \mathcal{N}(\mathbf{0}, \sigma^2 I)$
    - **Masking noise:** zero out some fraction $\nu$ of components of $x$
    - **Salt-and-pepper noise:** choose some fraction $\nu$ of components of $x$ and set each to its min or max value (equally likely)

## Slide 3: Denoising Autoencoders

### Denoising Autoencoders

## Slide 4: Denoising Autoencoders — Example

### Denoising Autoencoders
Example

Original Images

Noisy Input

Autoencoder Output

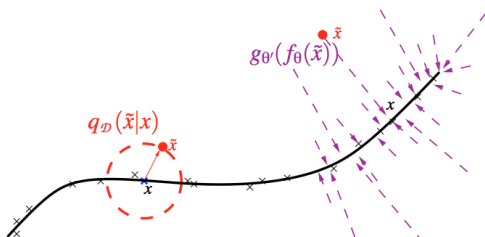## Slide 5: Denoising Autoencoders

### Denoising Autoencoders

- How does it work?
- Even though, e.g., MNIST data are in a 784-dimensional space, they lie on a low-dimensional **manifold** that captures their most important features
- Corruption process moves instance $x$ off of manifold
- Encoder $f_\theta$ and decoder $g_{\theta'}$ are trained to project $\tilde{x}$ back onto manifold

## Slide 6: Sparse Autoencoders

### Sparse Autoencoders

- An overcomplete architecture
- Regularize outputs of hidden layer to enforce **sparsity**:
$$\tilde{\mathcal{J}}(x) = \mathcal{J}(x, g(f(x))) + \alpha \, \Omega(h) \ ,$$
    where $\mathcal{J}$ is loss function, $f$ is encoder, $g$ is decoder, $h = f(x)$, and $\Omega$ penalizes non-sparsity of $h$
- E.g., can use $\Omega(h) = \sum_i |h_i|$ and ReLU activation to force many zero outputs in hidden layer
- Can also measure average activation of $h_i$ across mini-batch and compare it to user-specified **target sparsity** value $p$ (e.g., 0.1) via square error or **Kullback-Leibler divergence**:
$$p \log \frac{p}{q} + (1-p) \log \frac{1-p}{1-q} \ ,$$
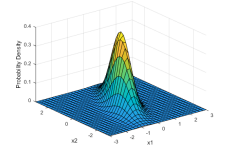    where $q$ is average activation of $h_i$ over mini-batch

## Contractive Autoencoders

- Similar to sparse autoencoder, but use

$$\Omega(\boldsymbol{h}) = \sum_{j=1}^{m} \sum_{i=1}^{n} \left( \frac{\partial h_i}{\partial x_j} \right)^2$$

- I.e., penalize large partial derivatives of encoder outputs wrt input values
- This **contracts** the output space by mapping input points in a neighborhood near $\boldsymbol{x}$ to a smaller output neighborhood near $f(\boldsymbol{x})$
  - $\Rightarrow$ Resists perturbations of input $\boldsymbol{x}$
- If $\boldsymbol{h}$ has sigmoid activation, encoding near binary and a CE pushes embeddings to corners of a hypercube

19 / 34

---
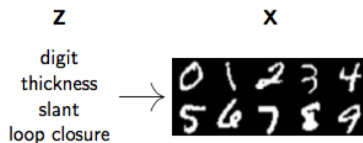
## Variational Autoencoders

- VAE is an autoencoder that is also **generative model**
  - $\Rightarrow$ Can generate new instances according to a probability distribution
  - E.g., hidden Markov models, Bayesian networks
  - Contrast with **discriminative models**, which predict classifications
- Encoder $f$ outputs $[\boldsymbol{\mu}, \boldsymbol{\sigma}]^{\top}$
  - Pair $(\mu_i, \sigma_i)$ parameterizes Gaussian distribution for dimension $i = 1, \ldots, n$
  - Draw $z_i \sim \mathcal{N}(\mu_i, \sigma_i)$
  - Decode this **latent variable** $z$ to get $g(z)$
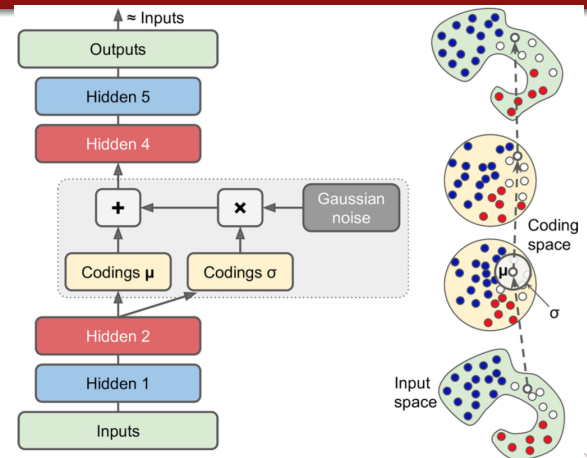


20 / 34

---

## Variational Autoencoders
### Latent Variables

- Independence of $z$ dimensions makes it easy to generate instances wrt complex distributions via decoder $g$
- Latent variables can be thought of as values of attributes describing inputs
  - E.g., for MNIST, latent variables might represent "thickness", "slant", "loop closure"



21 / 34

---

## Variational Autoencoders
### Architecture

22 / 34

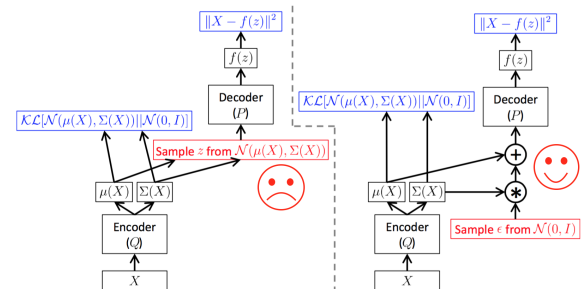---

## Variational Autoencoders
### Optimization

- **Maximum likelihood** (ML) approach for training generative models: find a model ($\theta$) with maximum probability of generating the training set $\mathcal{X}$
- Achieve this by minimizing the sum of:
  - End-to-end AE loss (e.g., square, cross-entropy)
  - Regularizer measuring distance (K-L divergence) from latent distribution $q(z \mid \boldsymbol{x})$ and $\mathcal{N}(\boldsymbol{0}, I)$ (= standard multivariate Gaussian)
- $\mathcal{N}(\boldsymbol{0}, I)$ also considered the **prior distribution** over $z$ (= distribution when no $\boldsymbol{x}$ is known)

```
eps = 1e-10
latent_loss = 0.5 * tf.reduce_sum(
        tf.square(hidden3_sigma) + tf.square(hidden3_mean)
        - 1 - tf.log(eps + tf.square(hidden3_sigma)))
```

23 / 34

---

## Variational Autoencoders
### Reparameterization Trick

- Cannot backprop error signal through random samples
- **Reparameterization trick** emulates $z \sim \mathcal{N}(\mu, \sigma)$ with $\epsilon \sim \mathcal{N}(0, 1)$, $z = \epsilon \sigma + \mu$



24 / 34

## Slide 25

CSCE 496/896 Lecture 5: Autoencoders

Stephen Scott

Introduction

Basic Idea

Stacked AE

Denoising AE

Sparse AE

Contractive AE

Variational AE

GAN

- Draw $z \sim \mathcal{N}(\mathbf{0}, I)$ and display $g(z)$

---

## Slide 26

CSCE 496/896 Lecture 5: Autoencoders

Stephen Scott

Introduction

Basic Idea

Stacked AE

Denoising AE

Sparse AE

Contractive AE

Variational AE

GAN

- Uniformly sample points in $z$ space and decode

---

## Slide 27

CSCE 496/896 Lecture 5: Autoencoders

Stephen Scott

Introduction

Basic Idea

Stacked AE

Denoising AE

Sparse AE

Contractive AE

Variational AE

GAN

- Cluster analysis by digit

---

## Slide 28

CSCE 496/896 Lecture 5: Autoencoders

Stephen Scott

Introduction

Basic Idea

Stacked AE

Denoising AE

Sparse AE

Contractive AE

Variational AE

GAN

- GANs are also generative models, like VAEs
- Models a game between two players
  - **Generator** creates samples intended to come from training distribution
  - **Discriminator** attempts to discern the "real" (original training) samples from the "fake" (generated) ones
- Discriminator trains as a binary classifier, generator trains to fool the discriminator

---

## Slide 29

CSCE 496/896 Lecture 5: Autoencoders

Stephen Scott

Introduction

Basic Idea

Stacked AE

Denoising AE

Sparse AE

Contractive AE

Variational AE

GAN

- Let $D(\boldsymbol{x})$ be discriminator parameterized by $\boldsymbol{\theta}^{(D)}$
  - Goal: Find $\boldsymbol{\theta}^{(D)}$ minimizing $J^{(D)}\left(\boldsymbol{\theta}^{(D)}, \boldsymbol{\theta}^{(G)}\right)$
- Let $G(\boldsymbol{z})$ be generator parameterized by $\boldsymbol{\theta}^{(G)}$
  - Goal: Find $\boldsymbol{\theta}^{(G)}$ minimizing $J^{(G)}\left(\boldsymbol{\theta}^{(D)}, \boldsymbol{\theta}^{(G)}\right)$
- A **Nash equilibrium** of this game is $(\boldsymbol{\theta}^{(D)}, \boldsymbol{\theta}^{(G)})$ such that each $\boldsymbol{\theta}^{(i)}$, $i \in \{D, G\}$ yields a local minimum of its corresponding $J$

---

## Slide 30

CSCE 496/896 Lecture 5: Autoencoders

Stephen Scott

Introduction

Basic Idea

Stacked AE

Denoising AE

Sparse AE

Contractive AE

Variational AE

GAN

- Each training step:
  - Draw a minibatch of $\boldsymbol{x}$ values from dataset
  - Draw a minibatch of $\boldsymbol{z}$ values from prior (e.g., $\mathcal{N}(\mathbf{0}, I)$)
  - Simultaneously update $\boldsymbol{\theta}^{(G)}$ to reduce $J^{(G)}$ and $\boldsymbol{\theta}^{(D)}$ to reduce $J^{(D)}$, via, e.g., Adam
- For $J^{(D)}$, common to use cross-entropy where label is 1 for real and 0 for fake
- Since generator wants to trick discriminator, can use $J^{(G)} = -J^{(D)}$
  - Others exist that are generally better in practice, e.g., based on ML

CSCE
496/896
Lecture 5:
Autoencoders

Stephen Scott

Introduction

Basic Idea

Stacked AE

Denoising AE

Sparse AE

Contractive AE

Variational AE

GAN

- "Deep, convolution GAN"
- Generator uses **transposed convolutions** (e.g.,
  `tf.layers.conv2d_transpose`) without pooling to
  upsample images for input to discriminator



31 / 34

CSCE
496/896
Lecture 5:
Autoencoders

Stephen Scott

Introduction

Basic Idea

Stacked AE

Denoising AE

Sparse AE

Contractive AE

Variational AE

GAN

Trained from LSUN dataset, sampled $z$ space



32 / 34

CSCE
496/896
Lecture 5:
Autoencoders

Stephen Scott

Introduction

Basic Idea

Stacked AE

Denoising AE

Sparse AE

Contractive AE

Variational AE

GAN

Trained from frame grabs of interview, sampled $z$ space



33 / 34

CSCE
496/896
Lecture 5:
Autoencoders

Stephen Scott

Introduction

Basic Idea

Stacked AE

Denoising AE

Sparse AE

Contractive AE

Variational AE

GAN

Performed semantic arithmetic in $z$ space!



(Non-center images have noise added in $z$ space; center is
noise-free)

34 / 34