# Chen Data Science & AI for Neuroscience Summer School

CHEN TIANQIAO & CHRISSY INSTITUTE

Caltech

# Autoencoders

Sabera Talukder

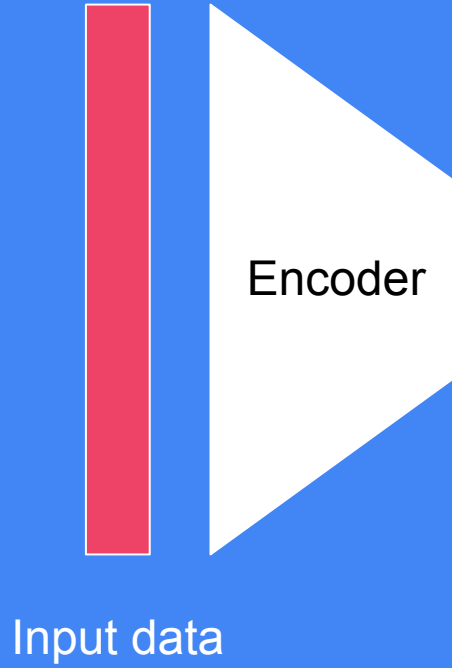Representations in the context of architectures!

# Principal Components Analysis
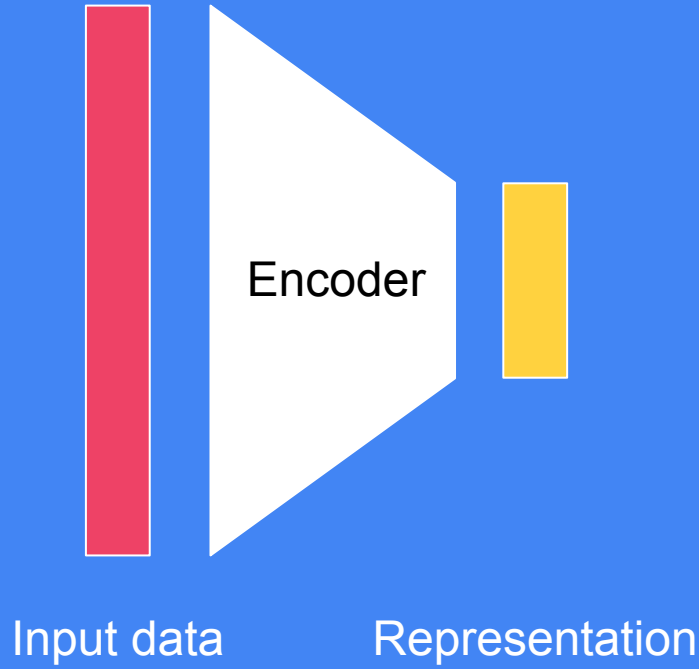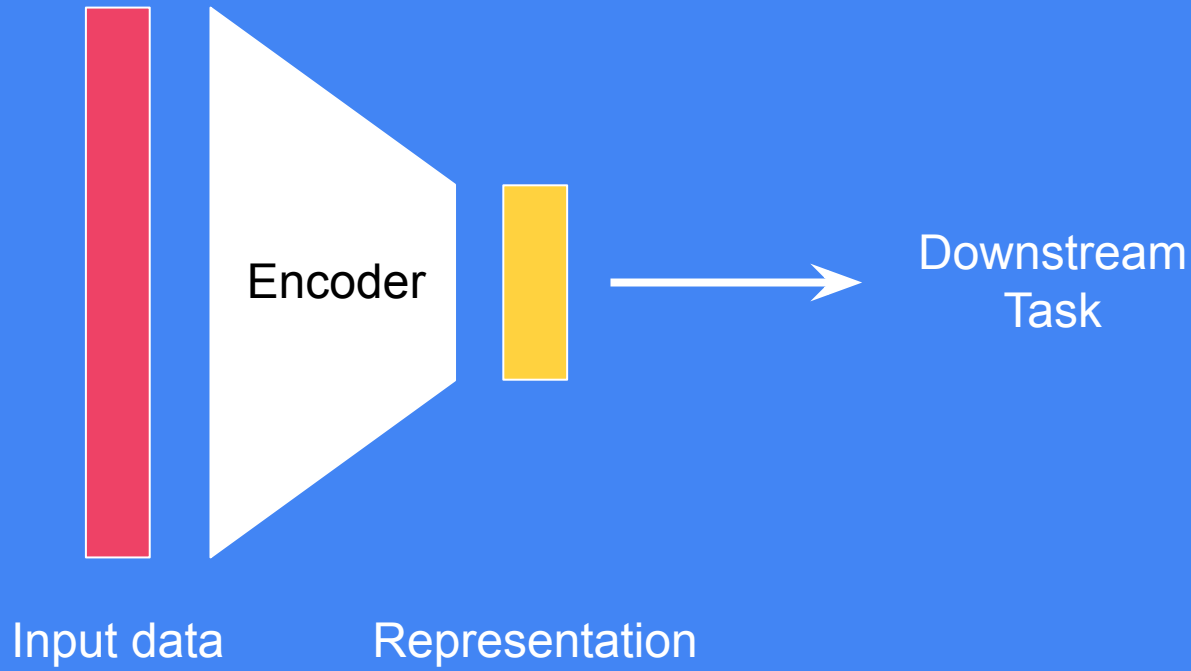
# Principal Components Analysis



Input data

# Principal Components Analysis



Encoder

Input data

# Principal Components Analysis



Encoder

Input data

Representation

# Principal Components Analysis



Input data

Encoder

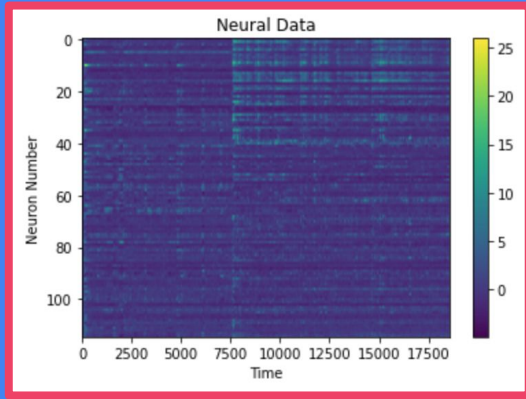Representation

Downstream Task

# Principal Components Analysis
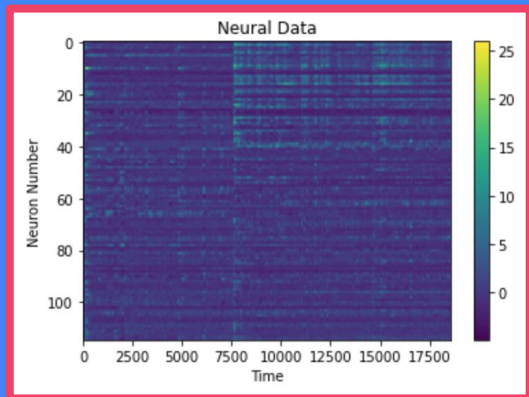


Input data

# Principal Components Analysis
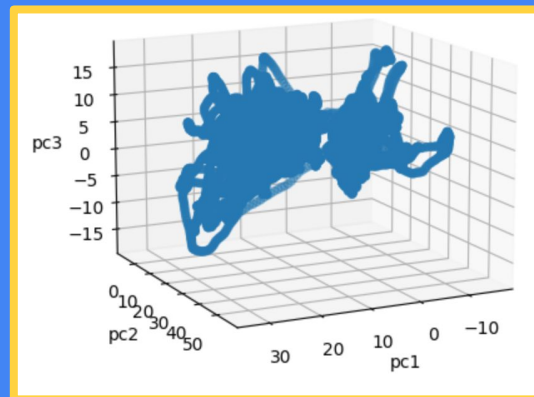


Input data

Encoder
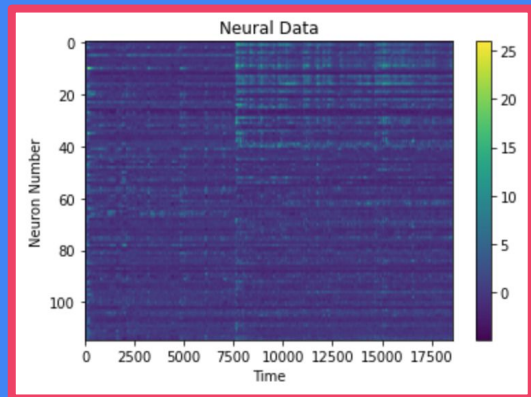is PCA

# Principal Components Analysis



Input data

Encoder is PCA
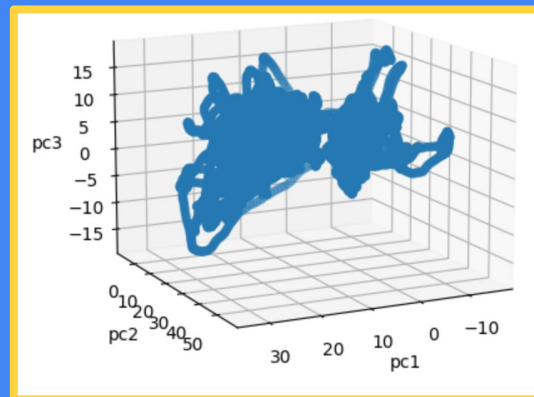
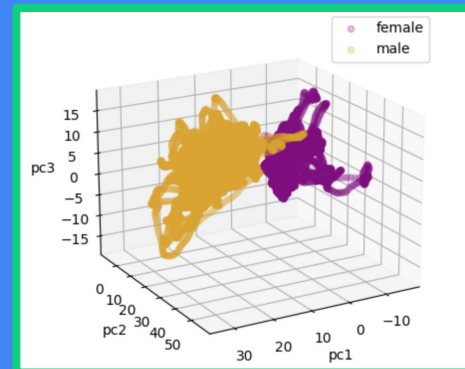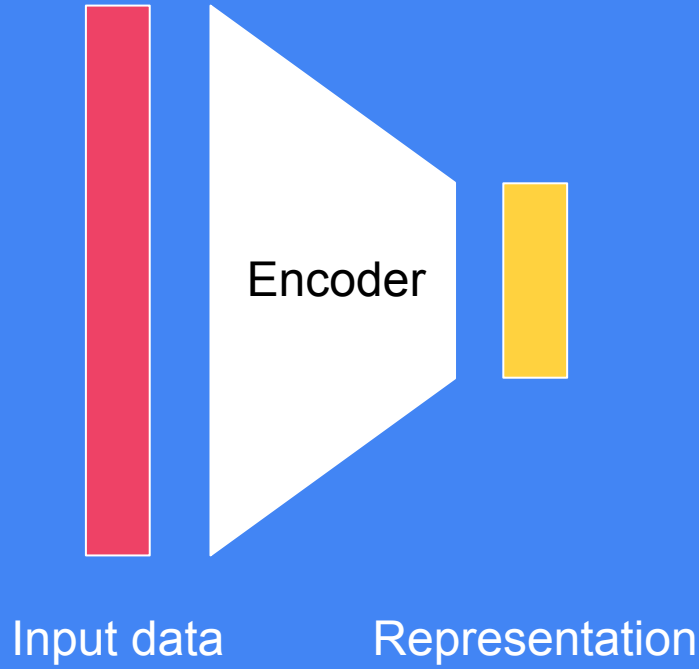Representation

# Principal Components Analysis



Input data

Encoder is PCA

Representation
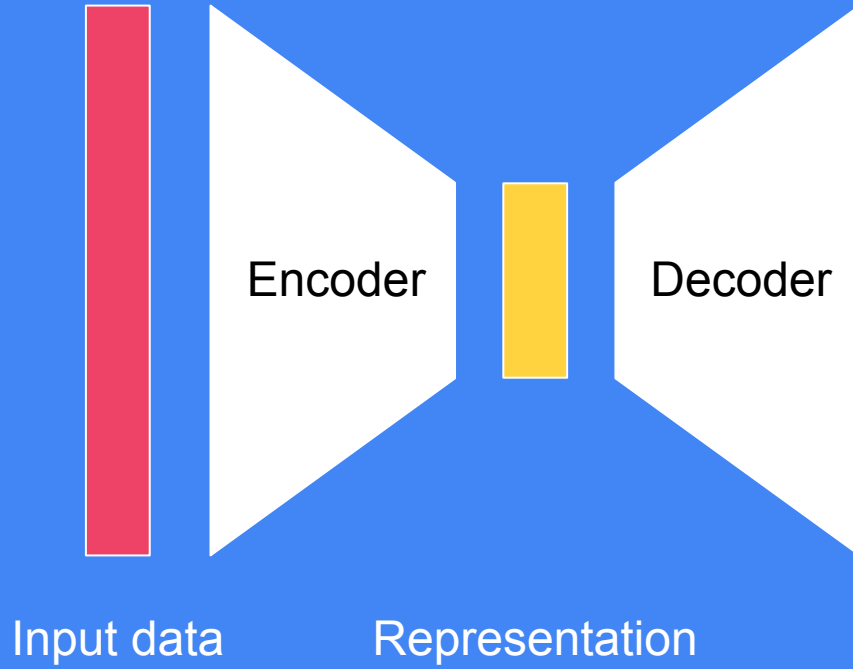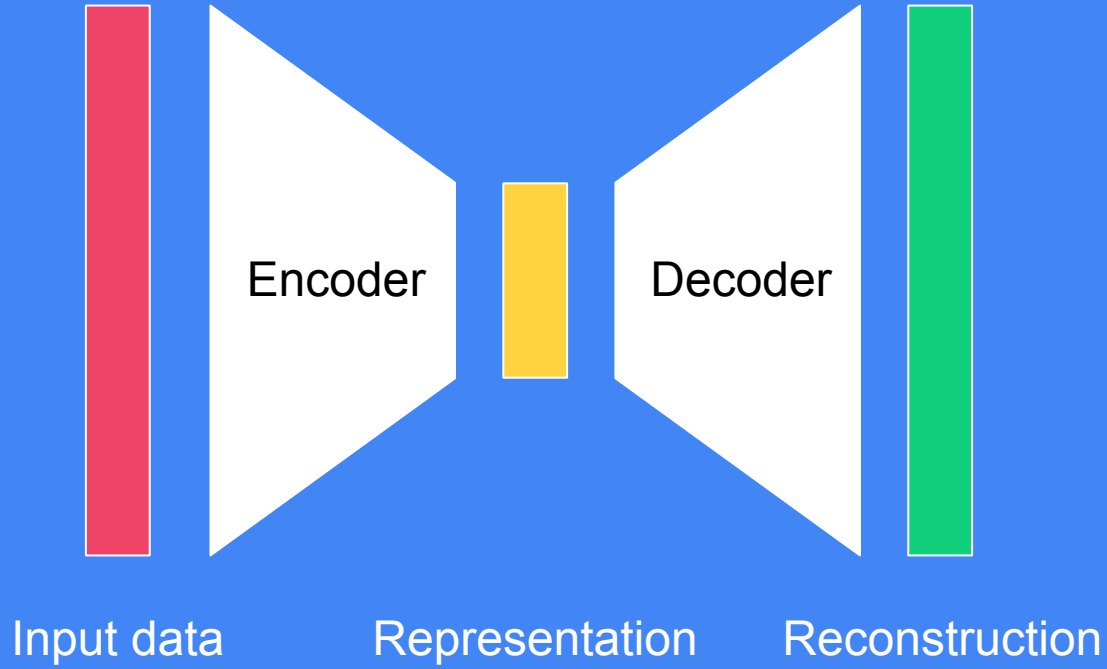
Downstream Task

# Autoencoder
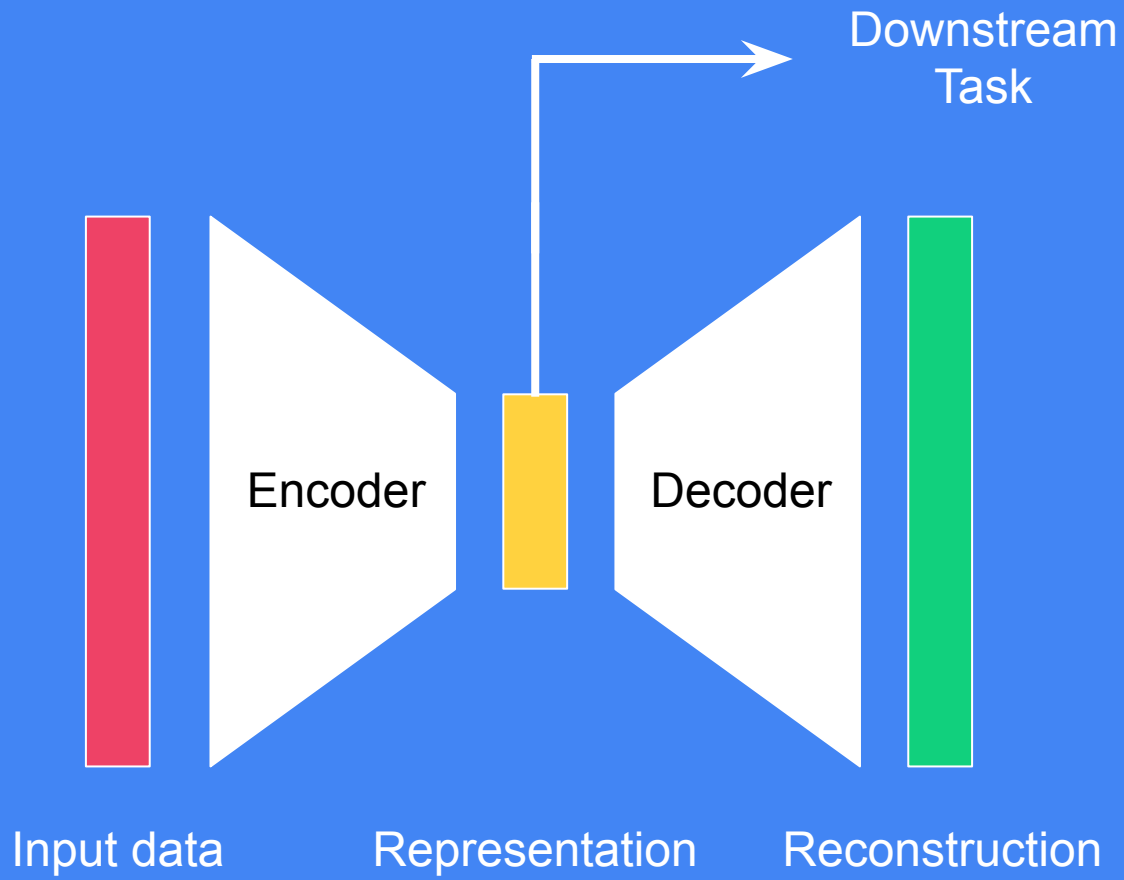


Input data          Representation

# Autoencoder



Input data       Representation

# Autoencoder



Input data      Representation      Reconstruction

Encoder      Decoder

# Autoencoder



Input data       Representation       Reconstruction

Encoder       Decoder

Downstream Task

# Autoencoder



Downstream Task

Encoder is a Neural Network!

Decoder is a Neural Network!

Input data

Representation

Reconstruction

# MNIST Example: Dataset
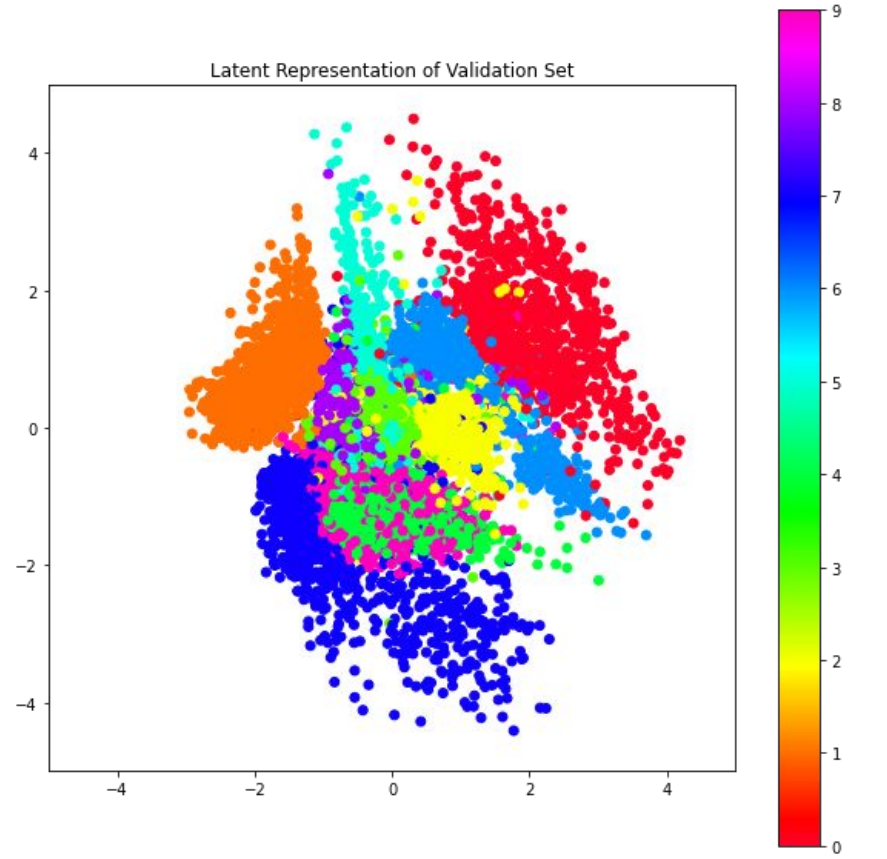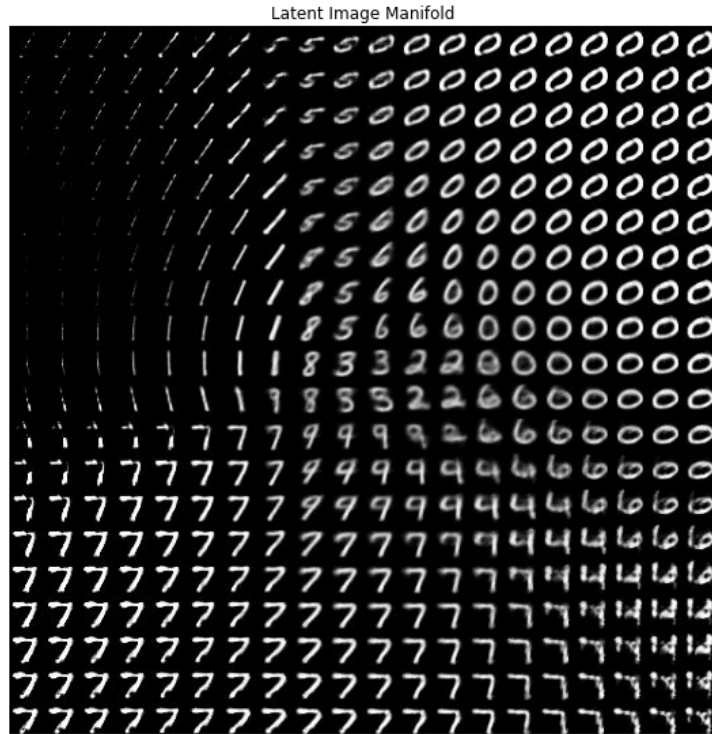
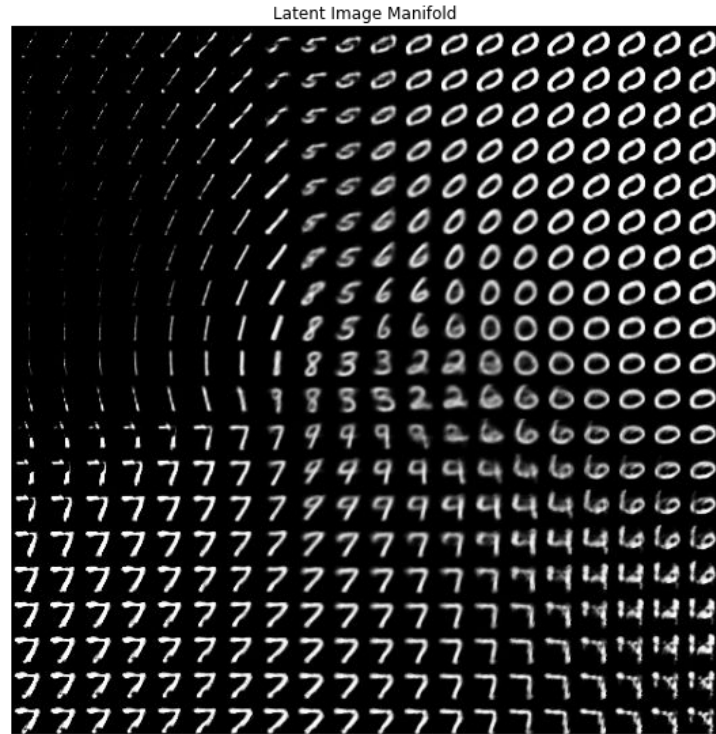# MNIST Example: Dataset

# MNIST Example: Dataset



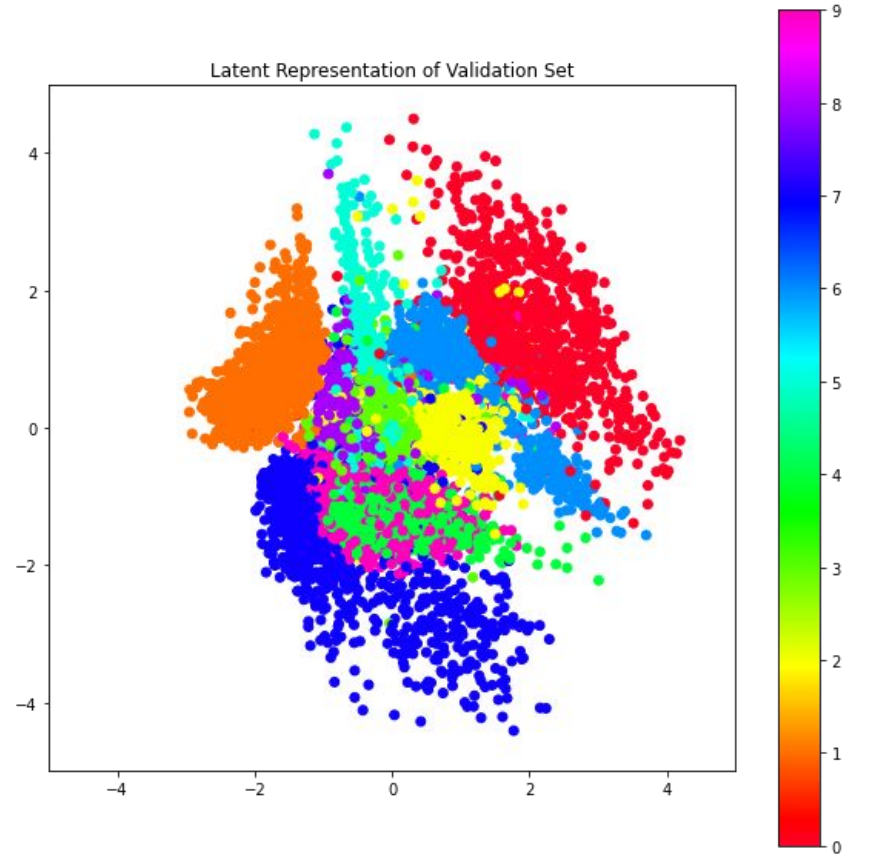Goal of using an autoencoder on MNIST is to generate handwritten digits!
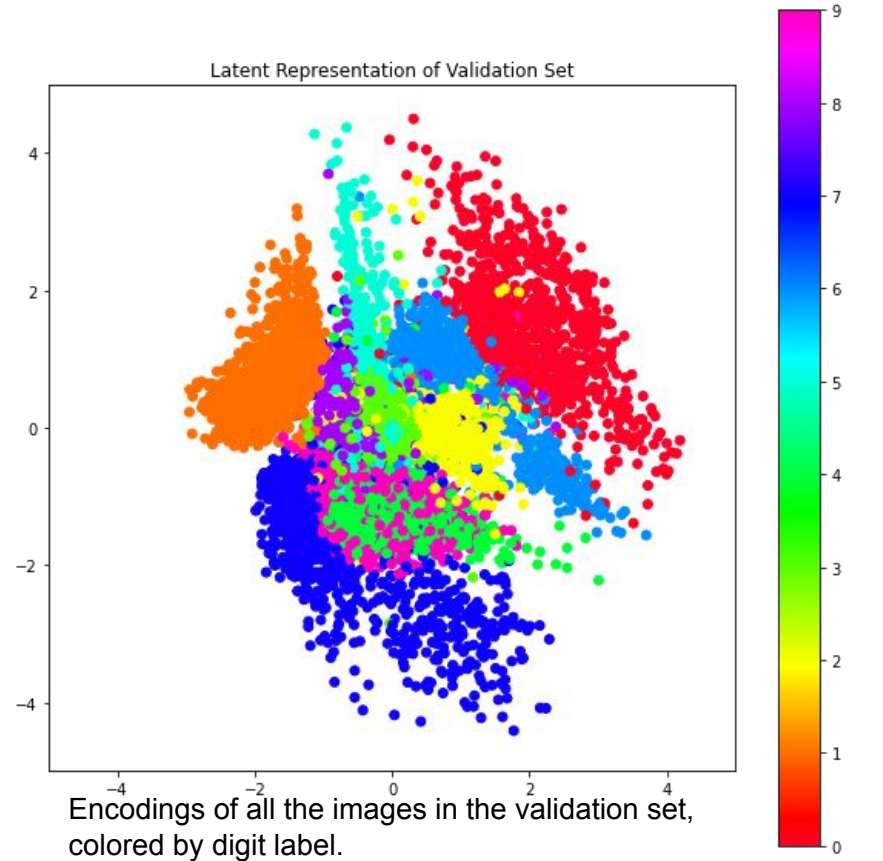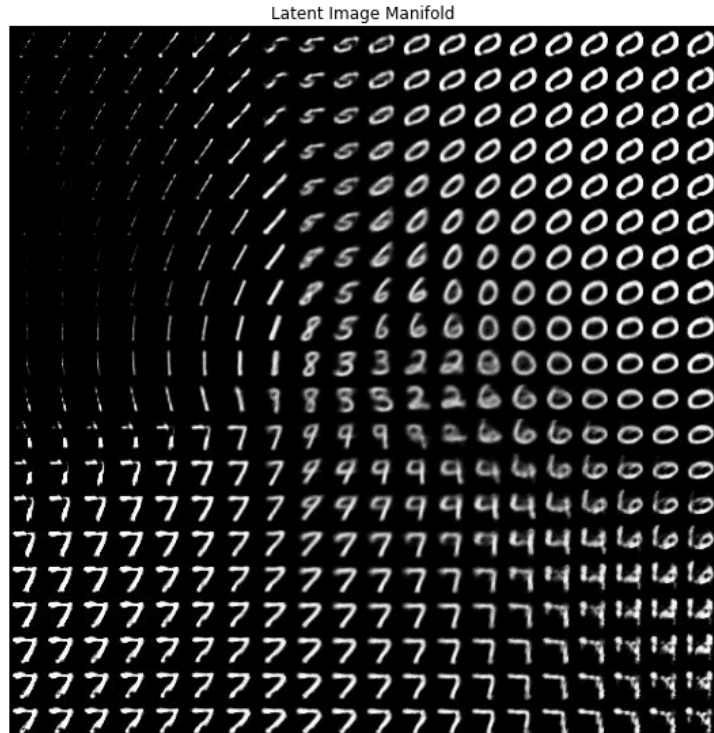
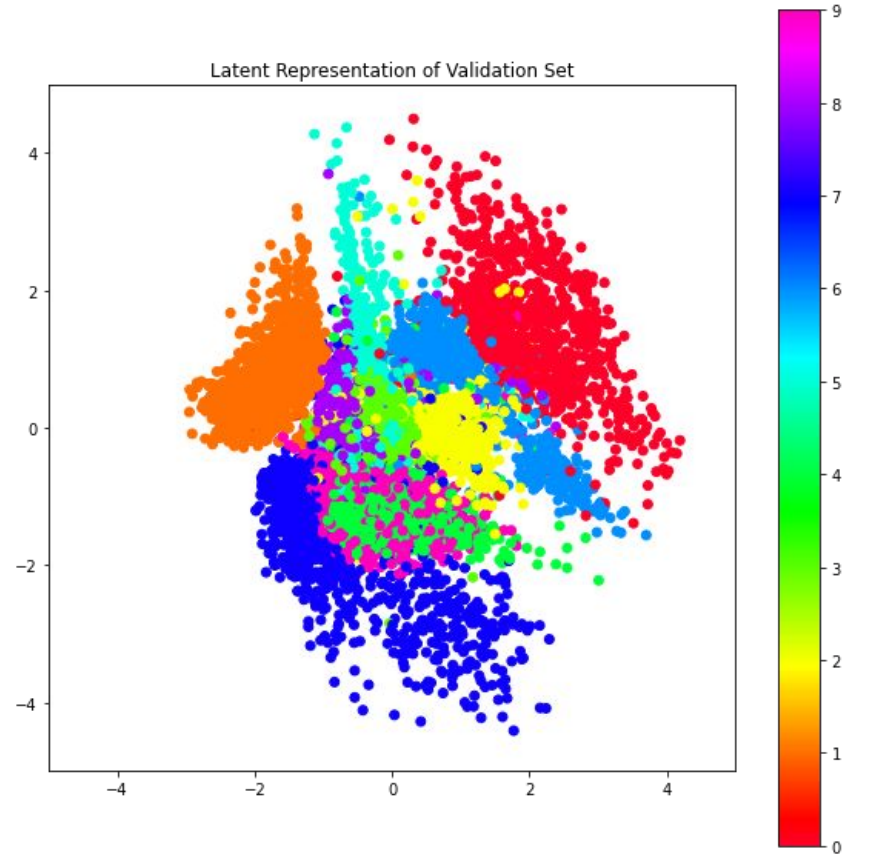# MNIST Example: 2D Latent Space

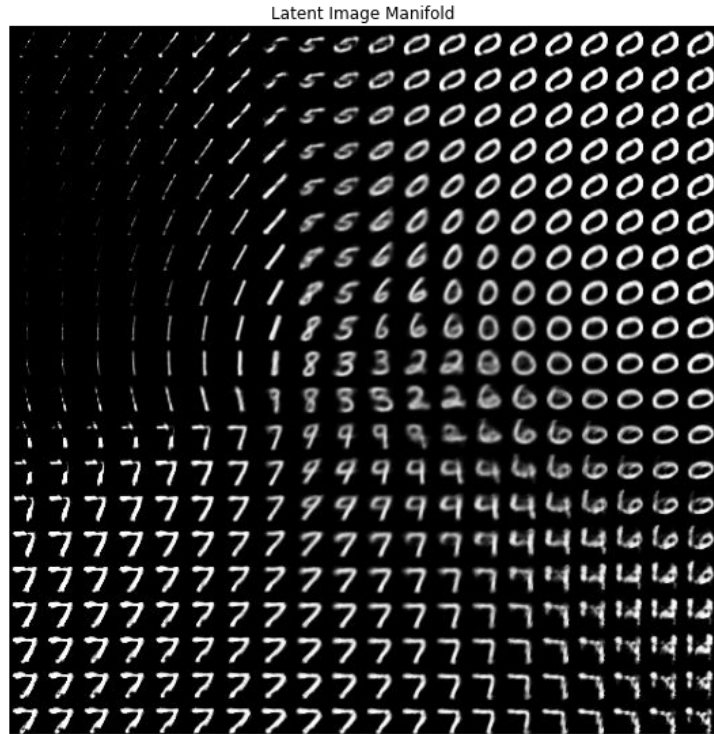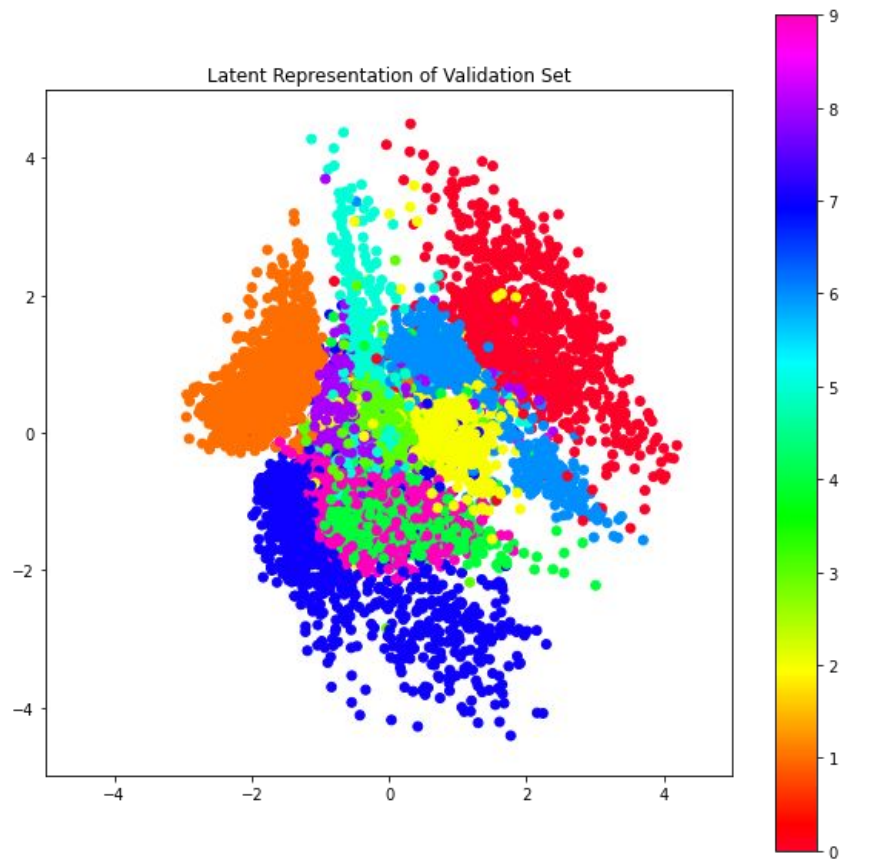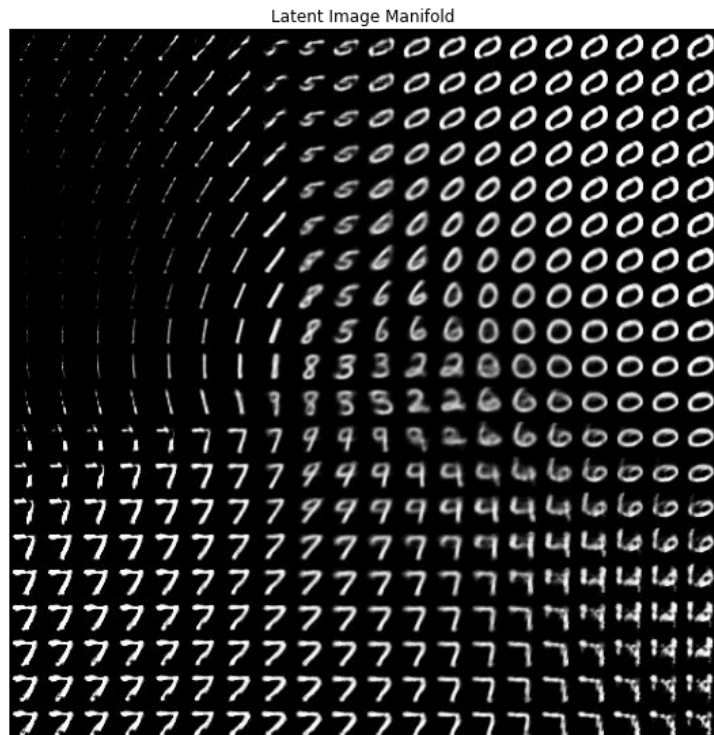# MNIST Example: 2D Latent Space



Sample images generated by picking points in the 2D latent space on a grid and then decoding them.

# MNIST Example: 2D Latent Space



Latent Image Manifold

Latent Representation of Validation Set

Encodings of all the images in the validation set, colored by digit label.

# MNIST Example: 2D Latent Space

**Takeaway 2:** This model has no prior concept of what each digit is. However, it learns to essentially embed digit types in localized regions of the latent space. Why? Because this representation is a useful one that it learns that helps with the task of reconstruction, which is what we want from good features.

# Autoencoder

# Autoencoder



Input data

Encoder
is a
Neural
Network!

# Autoencoder



Input data

Encoder is a Neural Network!

Representation

# Autoencoder



Input data

Encoder is a Neural Network!

Decoder is a Neural Network!

Reconstruction

Representation

# Autoencoder



Input data

Representation

Reconstruction
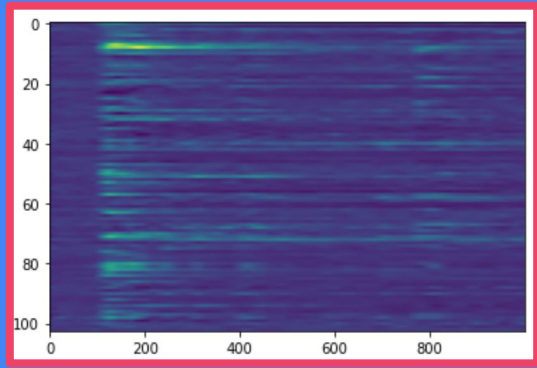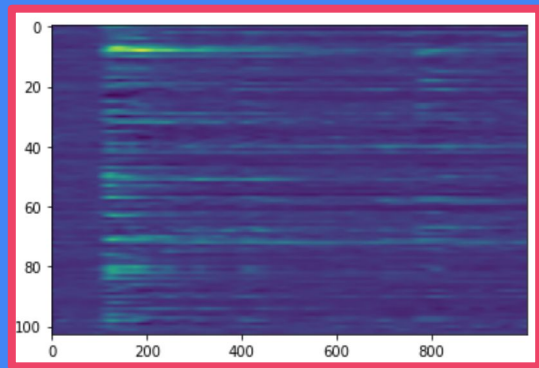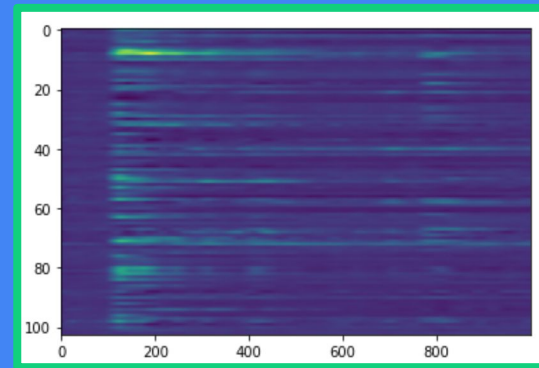
# Linear Layer

# Linear Layer

A "linear layer" is just a matrix that multiplies the input and produces some output.

# Linear Layer

A "linear layer" is just a matrix that multiplies the input and produces some output.

Fact: the composition of two linear layers is going to be linear. Why is this important? If we only use linear operations in our model, then the expressiveness of our model is going to be restricted significantly.
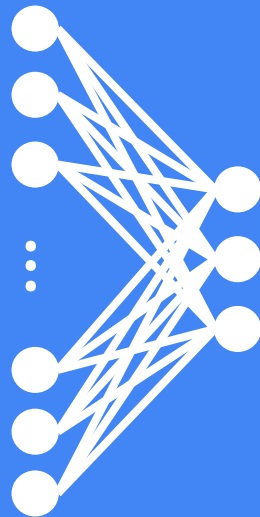
# Linear Layer

A "linear layer" is just a matrix that multiplies the input and produces some output.

Fact: the composition of two linear layers is going to be linear. Why is this important? If we only use linear operations in our model, then the expressiveness of our model is going to be restricted significantly.

Nonlinearities go between the linear layers (functions that go between the layers of a neural network are called activation functions) so that when we compose all the layers together, we get a nonlinear function. This allows us to build universal function approximators!
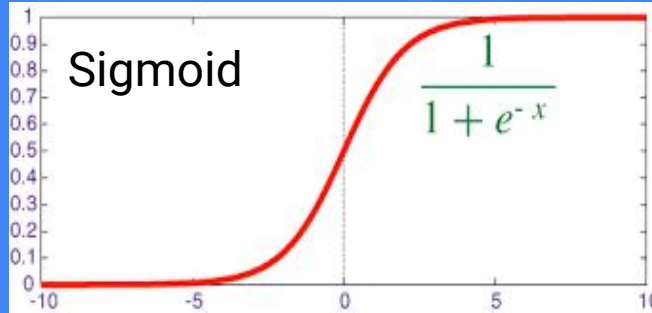
# Nonlinearities

We introduce nonlinearities in our model in the form of activation functions.
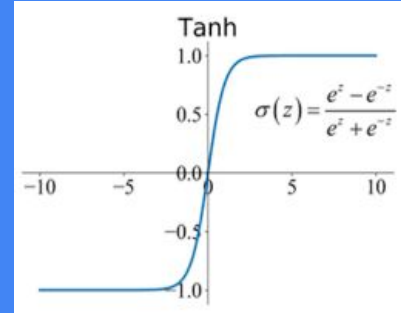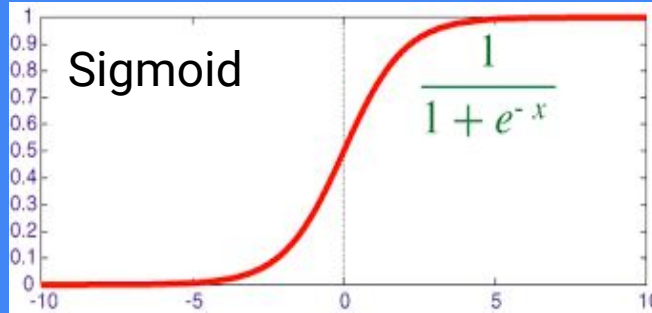
# Nonlinearities

We introduce nonlinearities in our model in the form of activation functions.



Sigmoid

$$\frac{1}{1 + e^{-x}}$$

# Nonlinearities

We introduce nonlinearities in our model in the form of activation functions.

# Nonlinearities

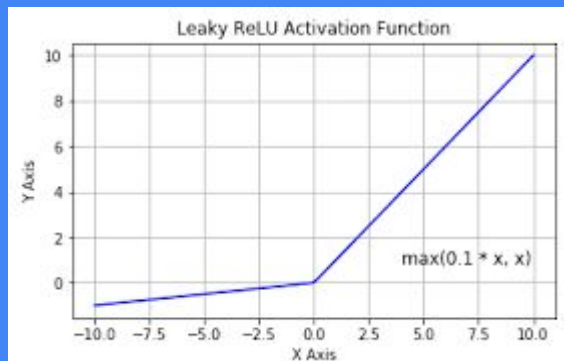We introduce nonlinearities in our model in the form of activation functions.



Sigmoid

$$\frac{1}{1 + e^{-x}}$$



Tanh

$$\sigma(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$



Leaky ReLU Activation Function

max(0.1 * x, x)

# Nonlinearities

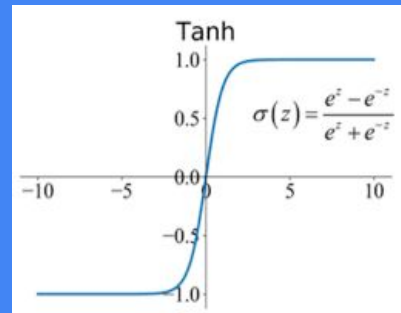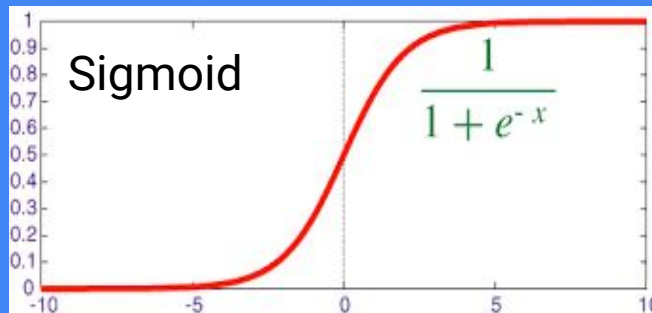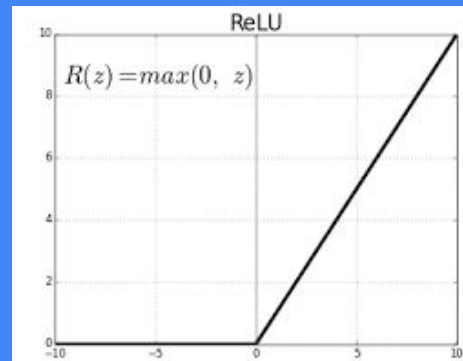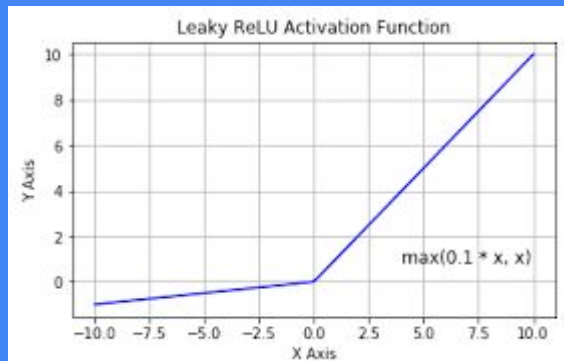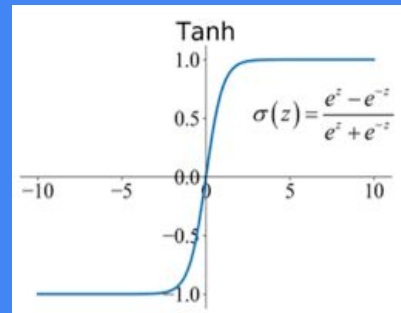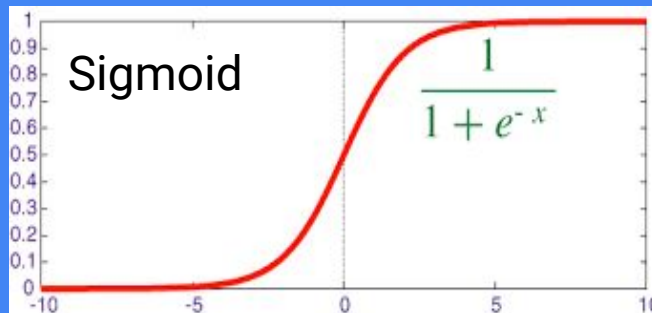We introduce nonlinearities in our model in the form of activation functions.



Sigmoid $\dfrac{1}{1 + e^{-x}}$

Tanh $\sigma(z) = \dfrac{e^z - e^{-z}}{e^z + e^{-z}}$

Leaky ReLU Activation Function $\max(0.1 * x, x)$
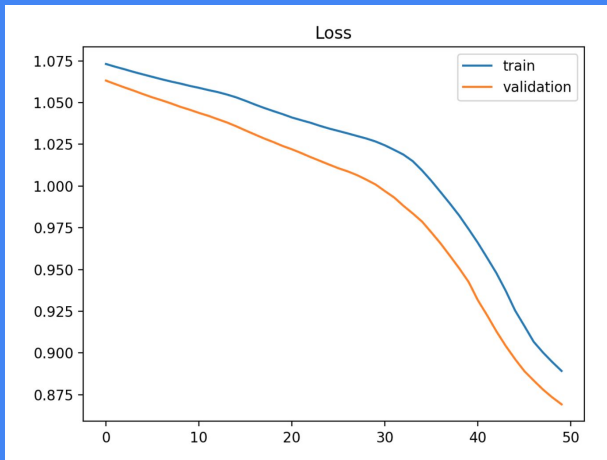
ReLU $R(z) = max(0, \; z)$

Jump ahead to training!

# Jump ahead to training!

Specifically underfitting and overfitting, and how to read training and validation curves!
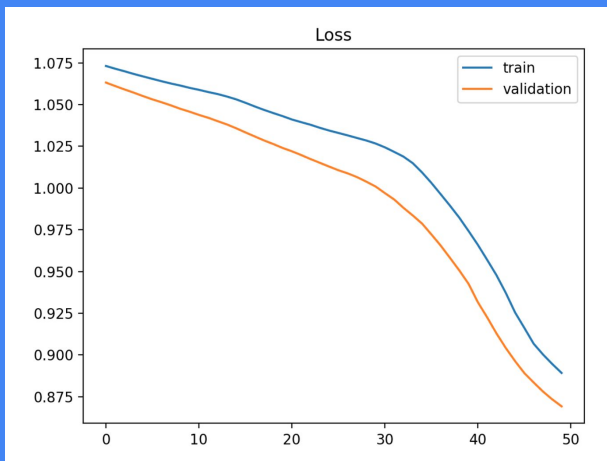
What is the training set versus the validation set?

# Underfitting = when the model doesn't have a low loss on the training set
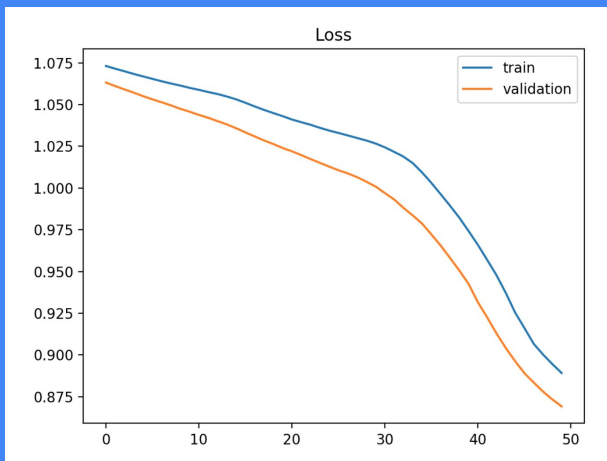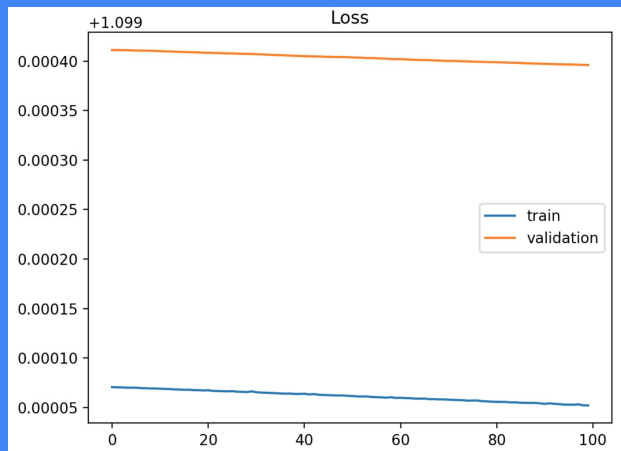
# Underfitting = when the model doesn't have a low loss on the training set

# Overfitting = when the model starts to memorize the training set

Underfitting = when the model doesn't have a low loss on the training set

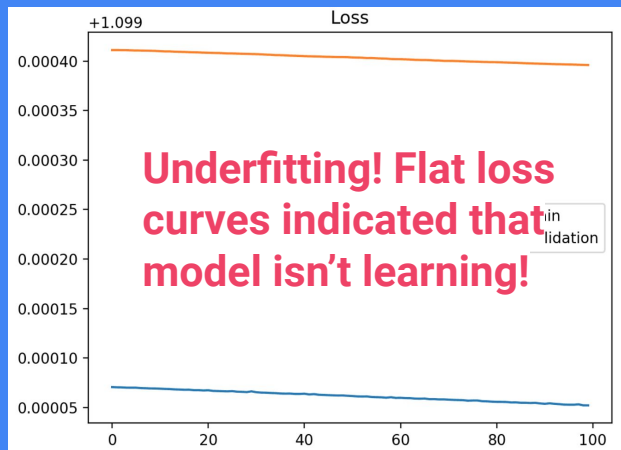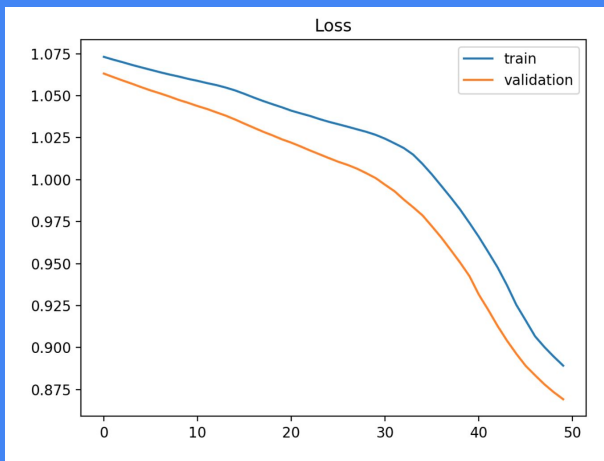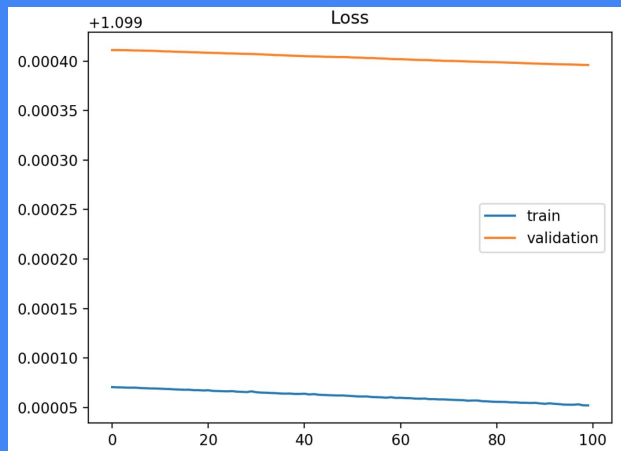Overfitting = when the model starts to memorize the training set

# Underfitting = when the model doesn't have a low loss on the training set

# Overfitting = when the model starts to memorize the training set



Loss

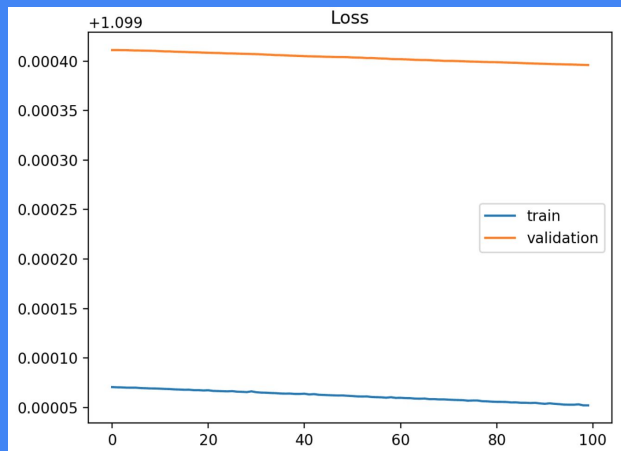**Underfitting! Flat loss curves indicated that model isn't learning!**

# Underfitting = when the model doesn't have a low loss on the training set

# Overfitting = when the model starts to memorize the training set

Jason Brownlee

# Underfitting = when the model doesn't have a low loss on the training set

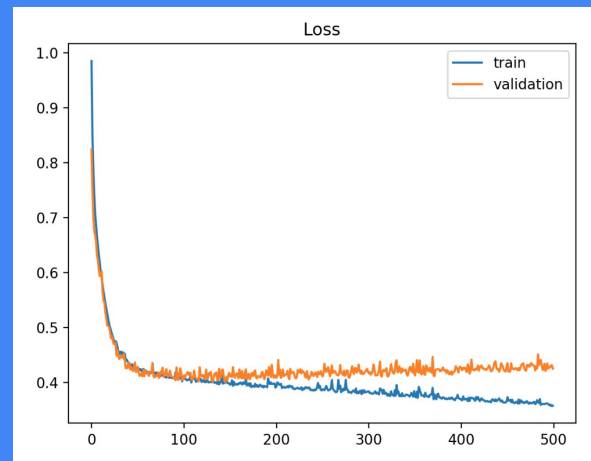# Overfitting = when the model starts to memorize the training set
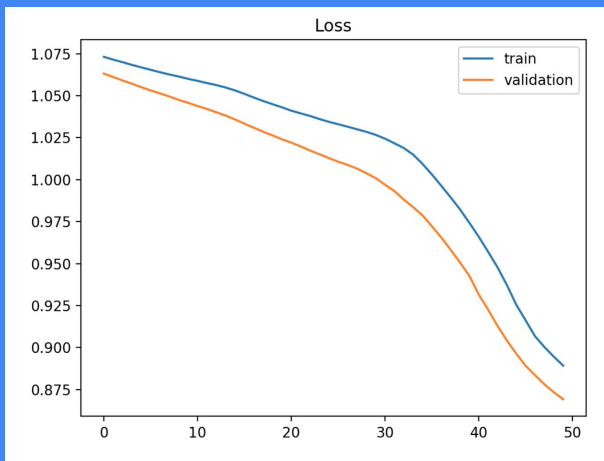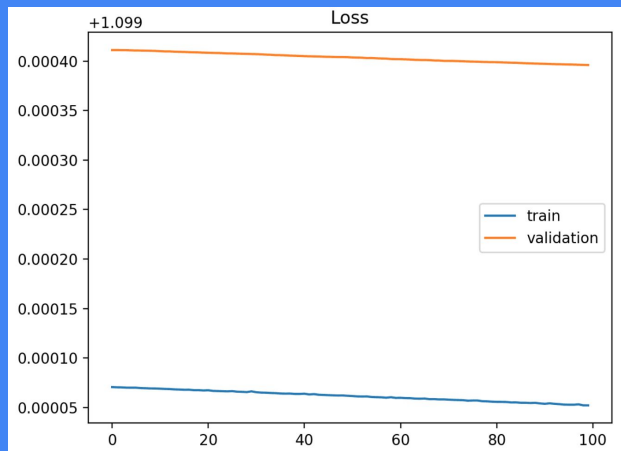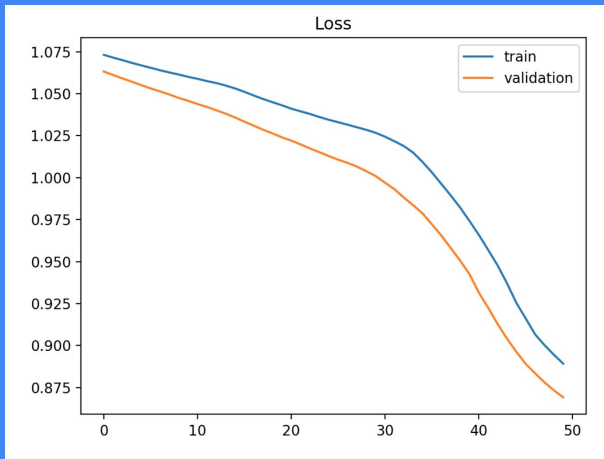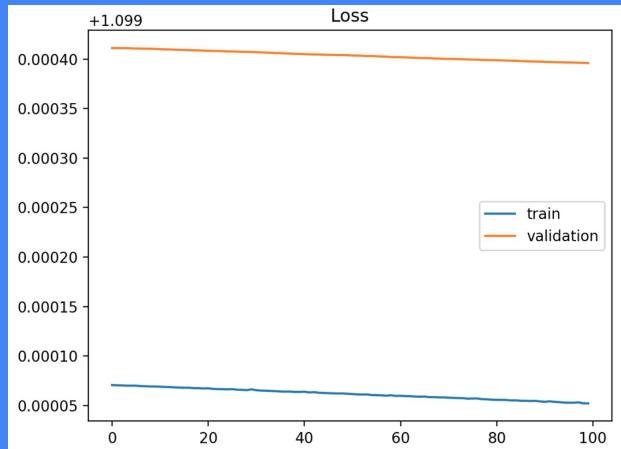
# Underfitting = when the model doesn't have a low loss on the training set

# Overfitting = when the model starts to memorize the training set



**Overfitting! The validation loss is going up while the training loss is going down!**

Jason Brownlee

What do you think a well trained model will look like with respect to the training and validation losses?

# Just Right!