# AUTOENCODERS

Lecture 4
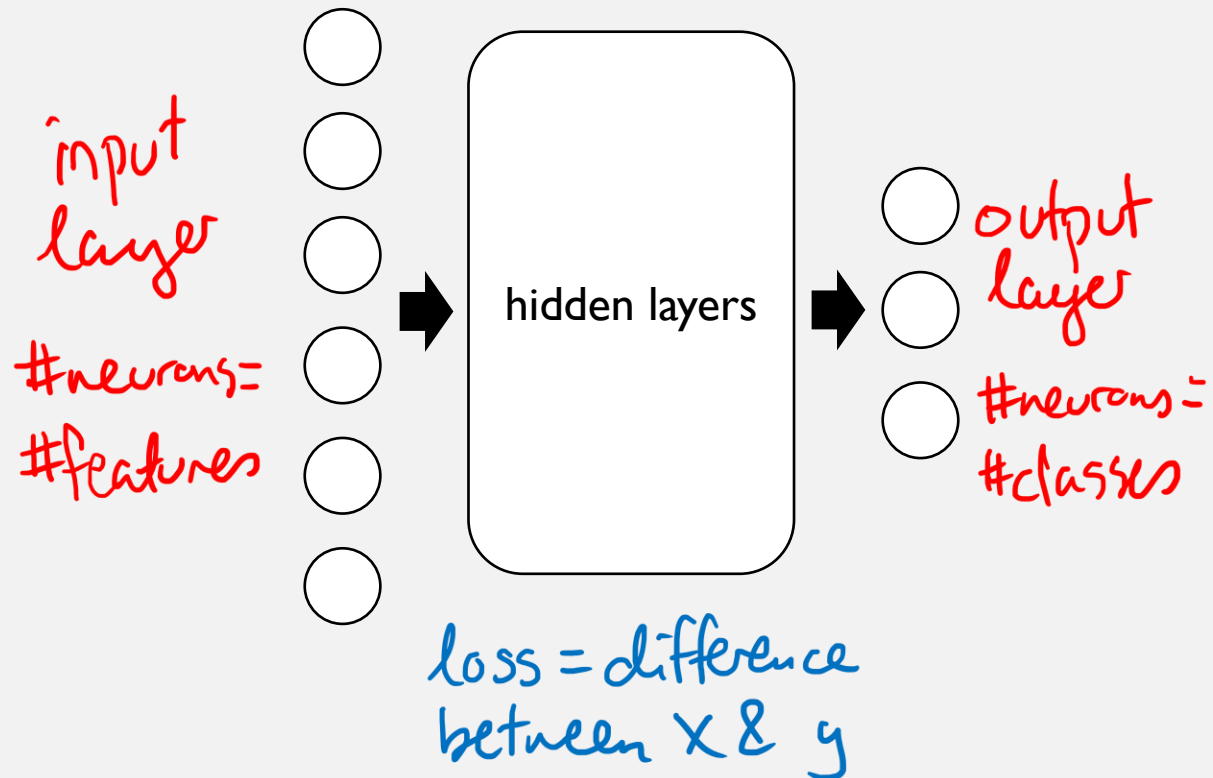
MAL2, Spring 2025

# AUTOENCODERS
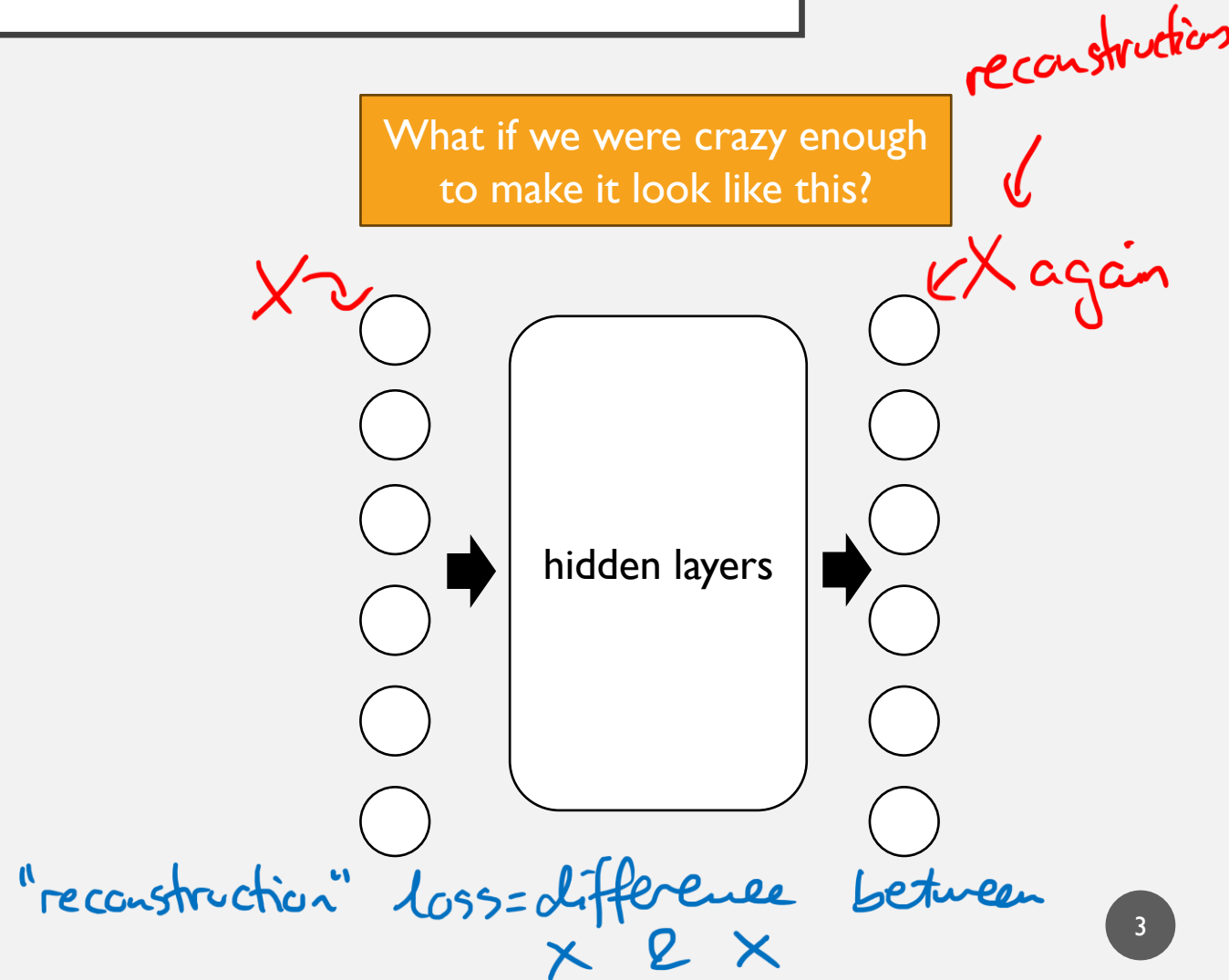
- **What if Y was X?**
- Dimensionality reduction
- Unsupervised pretraining
- Denoising images
- Colorization

# WHAT IF Y WAS X?

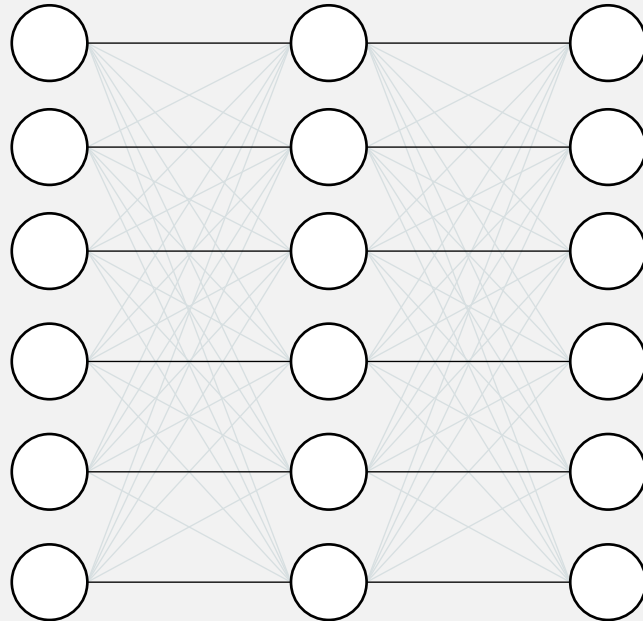Every neural network we have seen so far has looked like this:

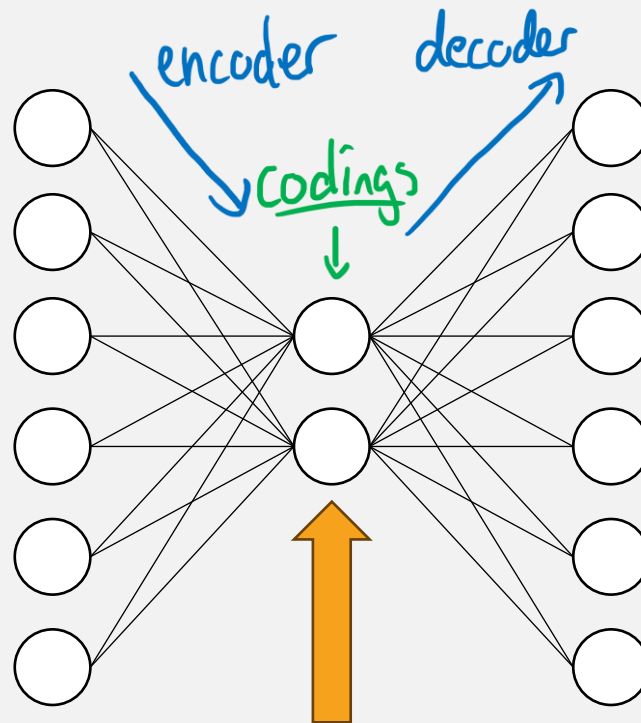What if we were crazy enough to make it look like this?

reconstruction

X~

kX again

input layer

#neurons=
#features

hidden layers

output layer

#neurons=
#classes

loss = difference between X & y

hidden layers

"reconstruction" loss = difference X & X

between
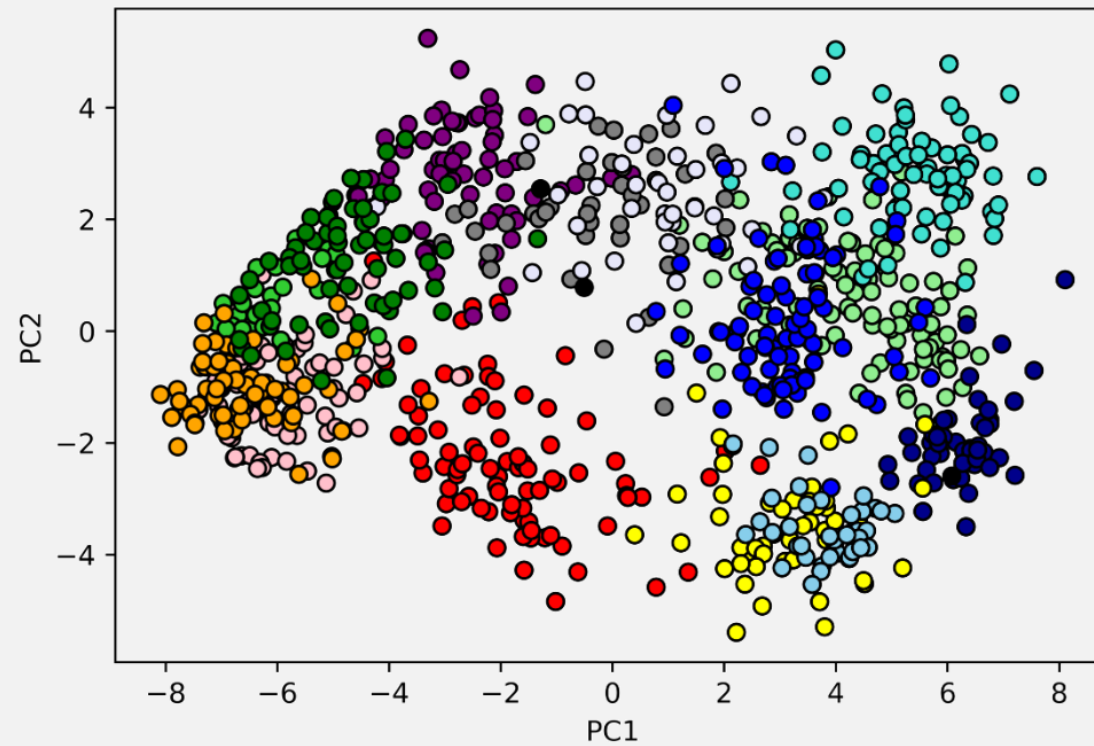
3

# BUT ISN'T THAT EASY?

# AN UNDERCOMPLETE AUTOENCODER



To reconstruct X, we need to express as much information as possible in these two neurons

# AUTOENCODERS

- What if Y was X?
- **Dimensionality reduction**
- Unsupervised pretraining
- Denoising images
- Colorization

# REMEMBER THIS?

# PERFORMING PCA WITH AN AUTOENCODER

*linear like PCA*

```python
encoder = Sequential([Dense(2)])
```
*encoder w/ 2 neurons, no activation*

```python
decoder = Sequential([Dense(49)])
```
*#questions*

```python
autoencoder = Sequential([encoder, decoder])
```
*construct full NN*

```python
optimizer = SGD(learning_rate=0.5)
```
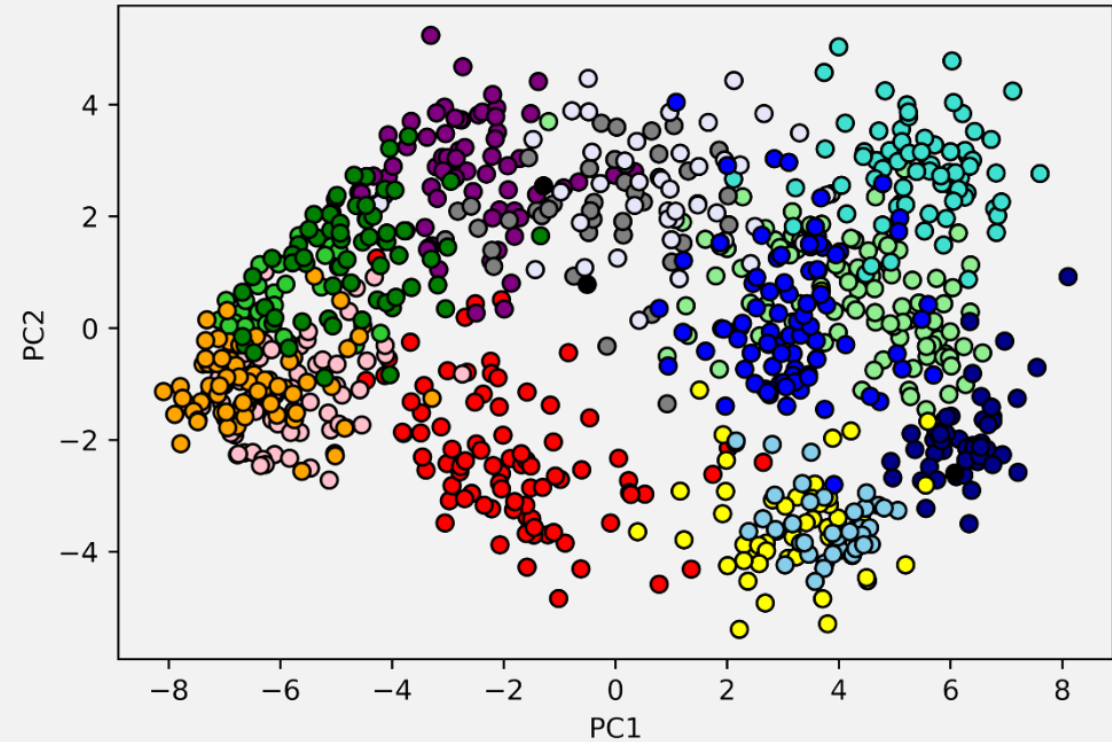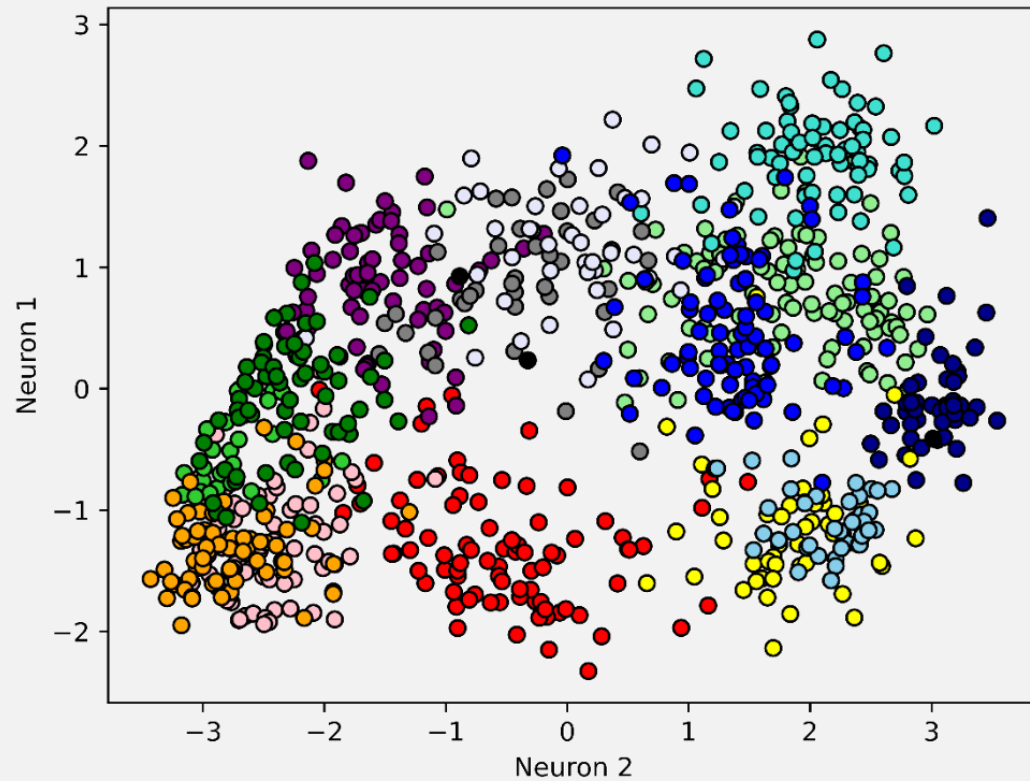*it will be fast, so GD ok*

```python
autoencoder.compile(loss="mse", optimizer=optimizer)
```
*reconstruction loss $(X - \hat{X})^2$*

```python
history = autoencoder.fit(X, X, epochs=500, verbose = False)
```
*long time*

```python
codings = encoder.predict(X)
```
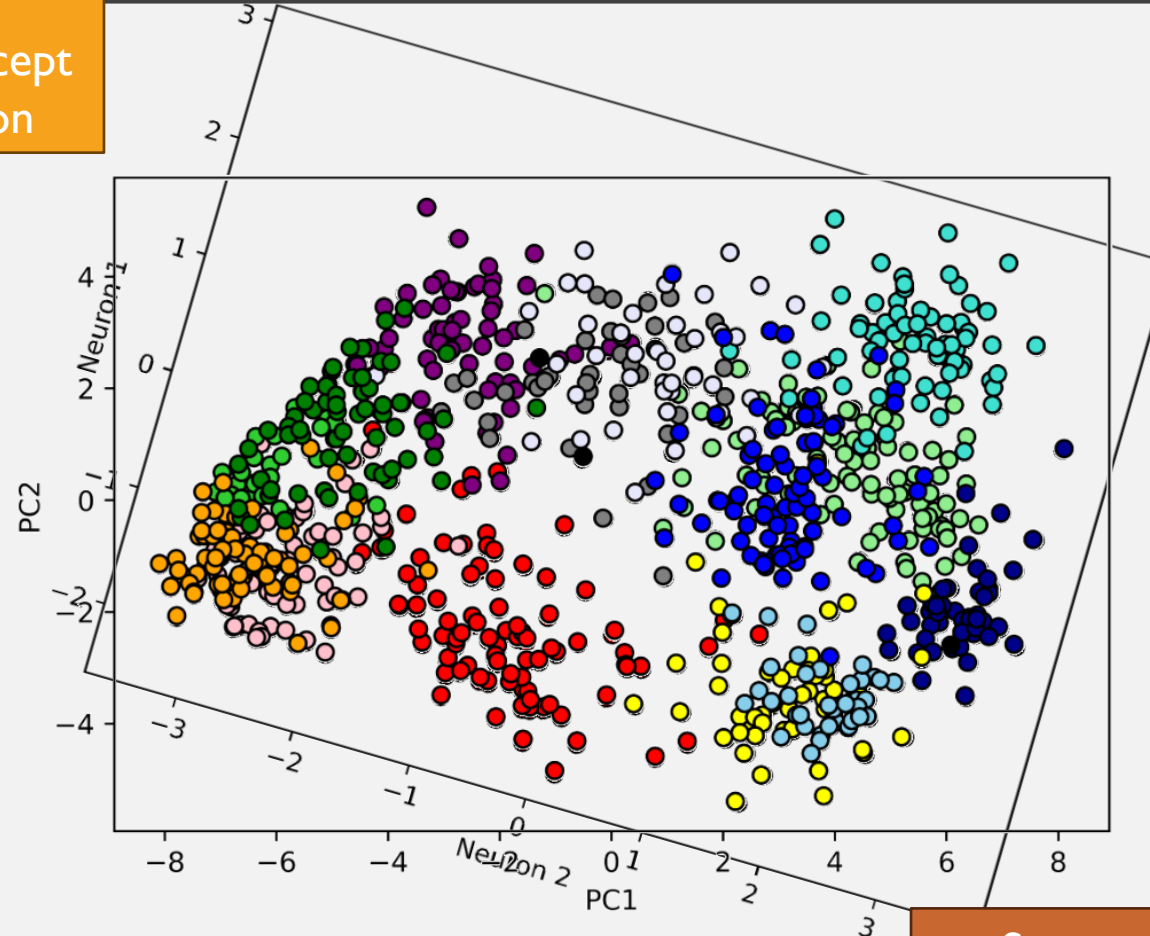
# AUTOENCODERS VS PCA
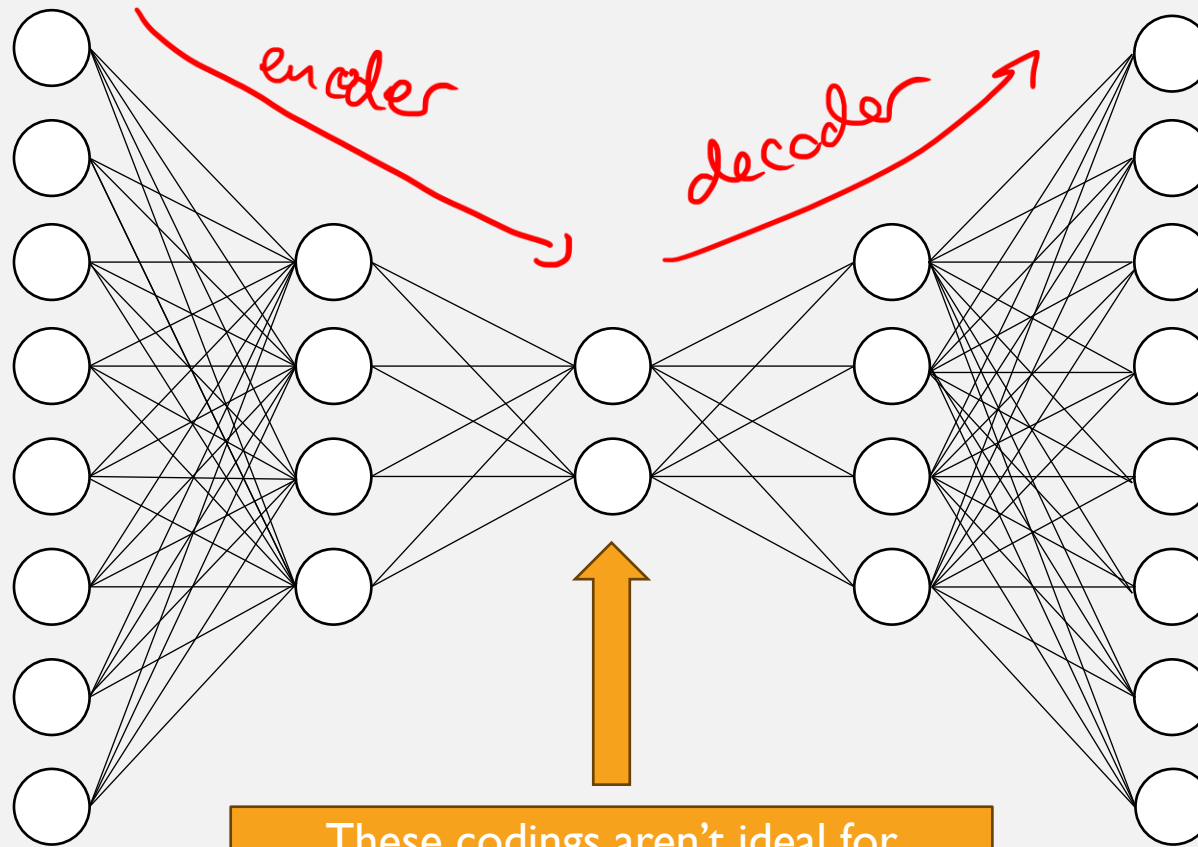
# AUTOENCODERS VS PCA

An autoencoder with linear activation is essentially PCA except for a scale factor and a rotation

linear algebra:
they span
the same
subspace



So now we can do PCA in an overcomplicated way … why care?

# DEEP AUTOENCODERS



encoder

decoder

symmetrical non-linear NN's capable of learning complex codings

These codings aren't ideal for visualization as the transformation to the reconstructions is complex
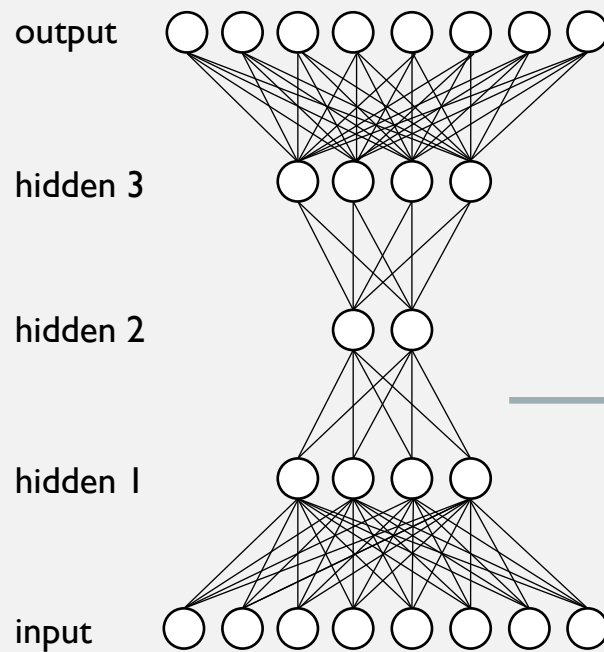
# AUTOENCODERS

- What if Y was X?
- Dimensionality reduction
- **Unsupervised pretraining**
- Denoising images
- Colorization
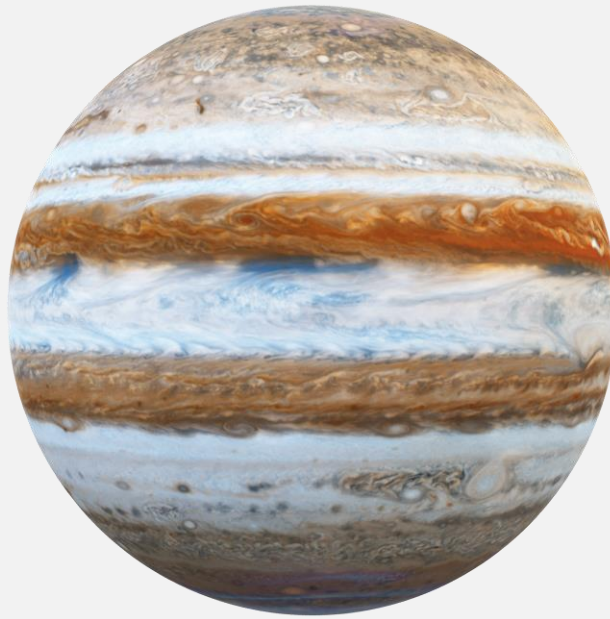
# AUTOENCODERS

- What if Y was X?
- Dimensionality reduction
- Unsupervised pretraining
- Denoising images
- Colorization

# DENOISING

Add noise to the input and try to reconstruct the noise-free image
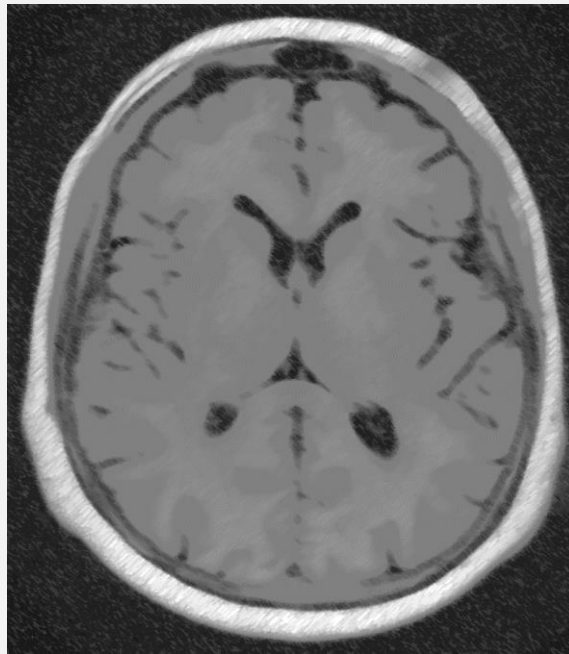
# DENOISING

**Take the autoencoder we just made**
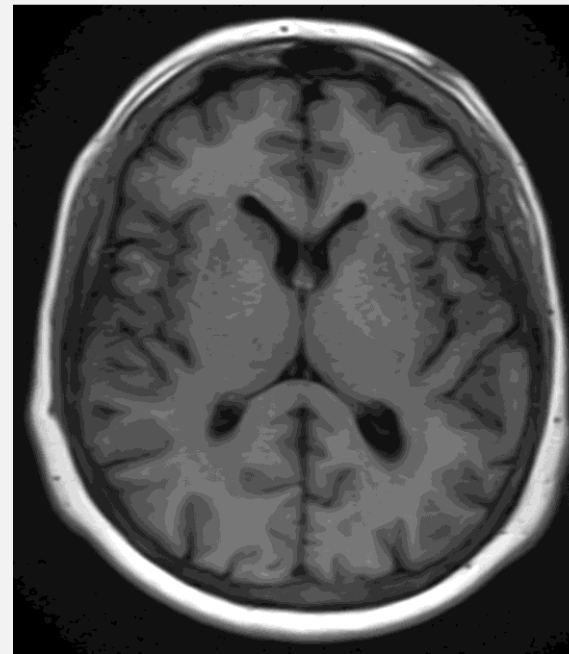


and experiment with the constraint

You have 15 minutes

# AN INTERESTING APPLICATION



low-radiation
dose image of brain

high-radiation
image

train AE to
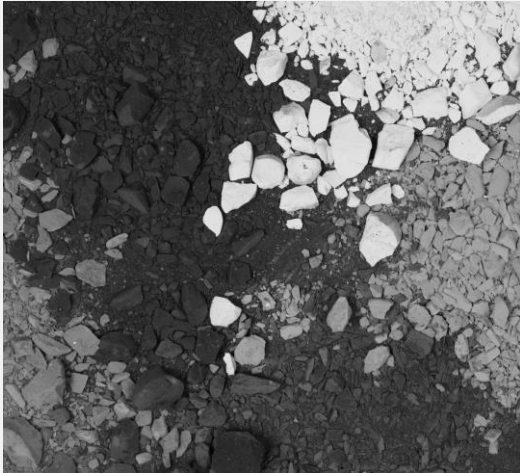turn low-res
into high-res
⟹ high-res
w/ low rad.
dose

# AUTOENCODERS

- What if Y was X?
- Dimensionality reduction
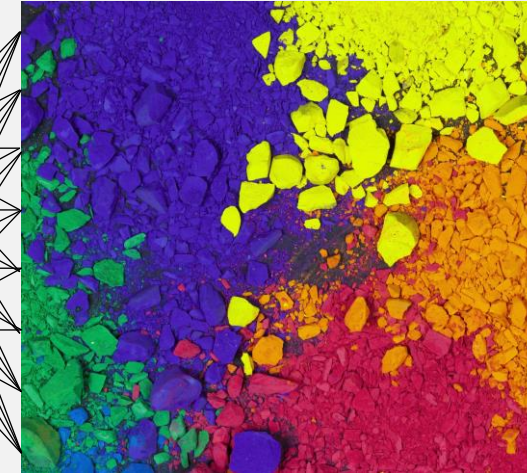- Unsupervised pretraining
- Denoising images
- **Colorization**

# COLORIZATION

2. Make a copy of the dataset and turn each picture black-and-white.

1. Get a dataset of color pictures.



3. Train an autoencoder to turn the black-and-white pictures into the colored ones.

```
conv_encoder = Sequential([
    Conv2D(…),
    MaxPool2D(…),
    …
])
```

```
conv_decoder = Sequential([
    Conv2D(…),
    UpSamling2D(…),      } Conv2DTranspose
    …
])
```

19

# YOUR TICKET OUT THE DOOR

**Scan this QR code**



and tell me about something
you are still unsure about