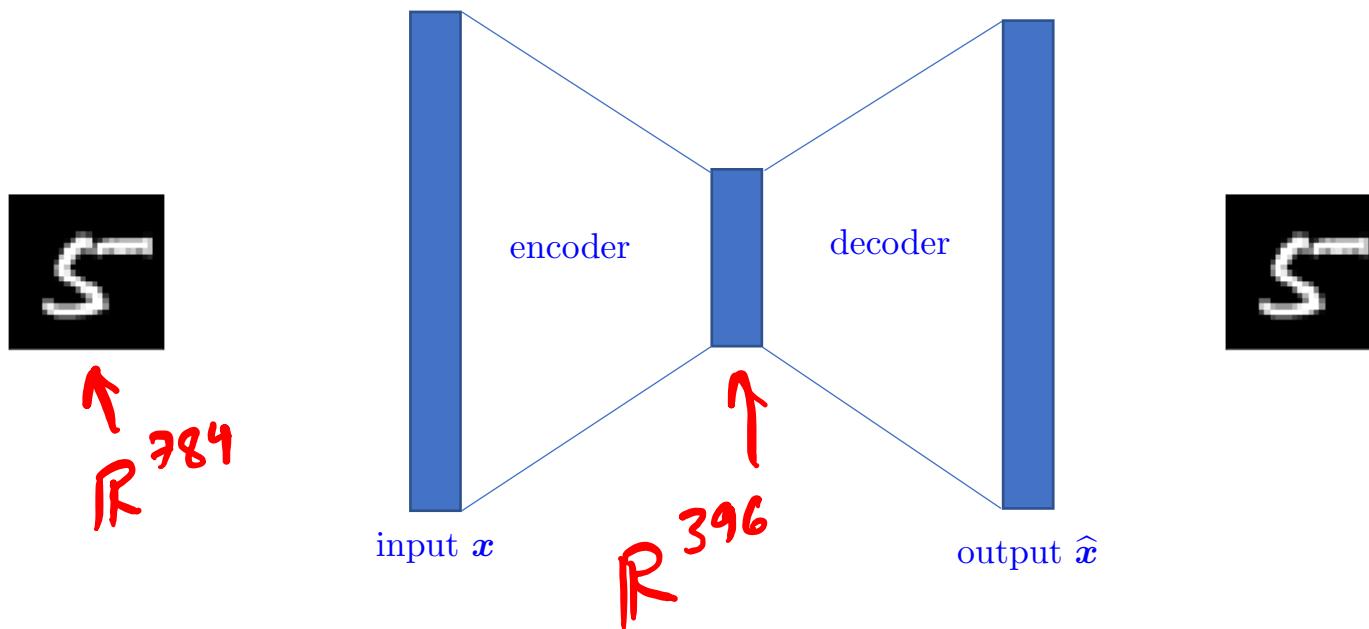


Neural Networks, Autoencoders

- neural network basics
 - autoencoders

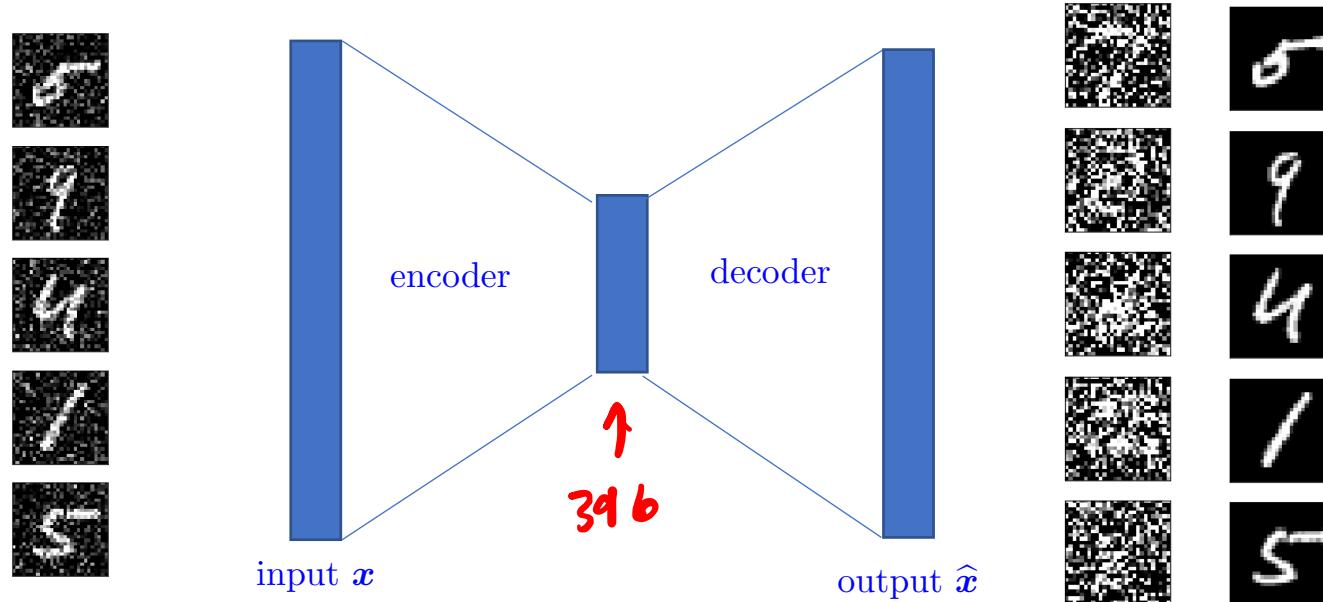
Autoencoders



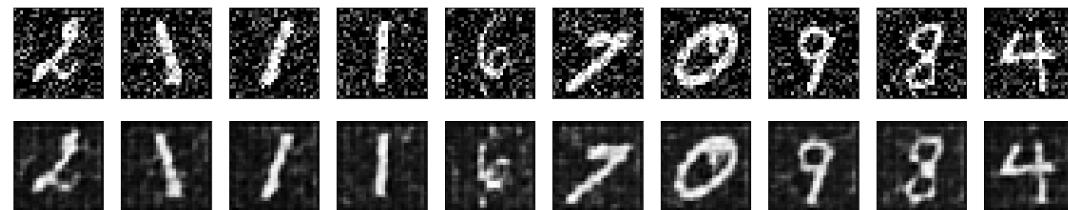
- main applications: denoising, data visualization, manifold learning
- autoencoders are (in general) **not** better than engineered compression algorithms

De-noising Autoencoder

$$\mathcal{D} = \{(\text{noisy image}, \text{clean image})\}$$



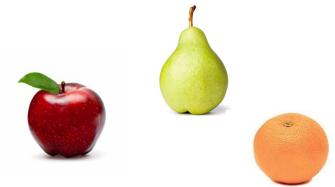
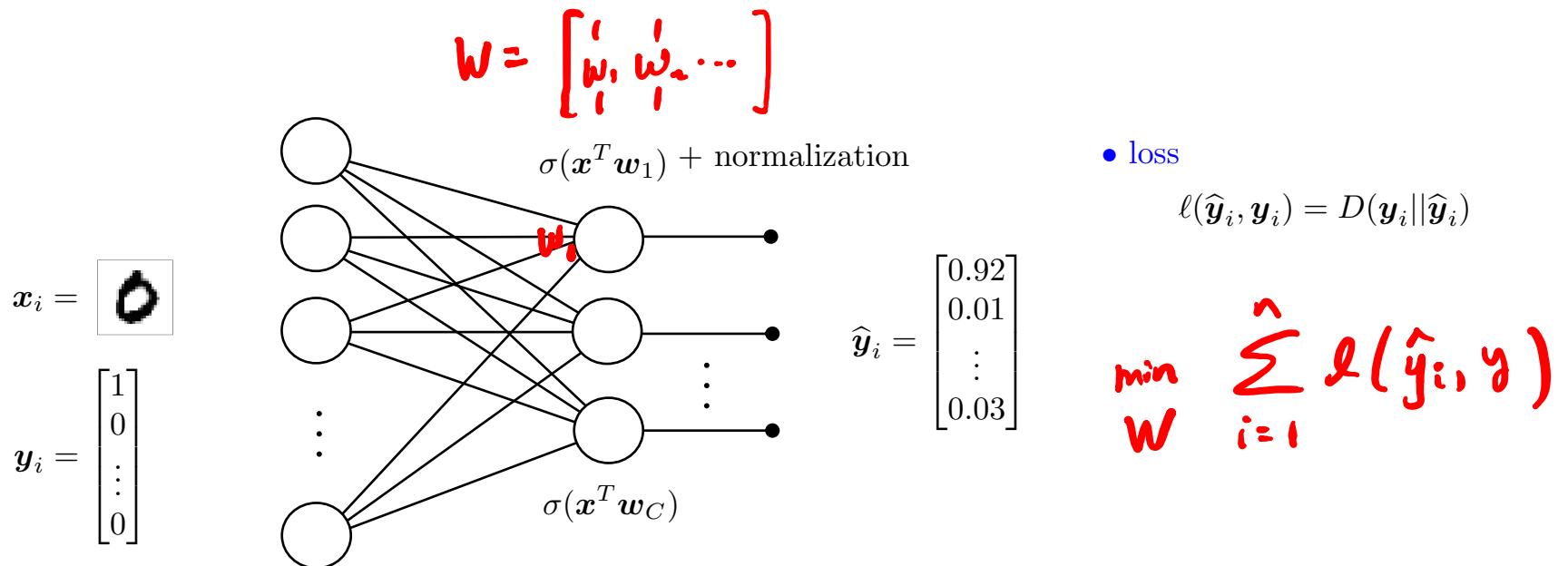
Test



Multinomial Logistic Regression, Single Layer

- logistic regression is a single *single layer of artificial neural network*

multiple neurons, single layer: $\hat{\mathbf{y}} = \sigma(\mathbf{W}^T \mathbf{x})$



- loss

$$\ell(\hat{\mathbf{y}}_i, \mathbf{y}_i) = D(\mathbf{y}_i || \hat{\mathbf{y}}_i)$$

$$\min_{\mathbf{W}} \sum_{i=1}^n \ell(\hat{\mathbf{y}}_i, \mathbf{y}_i)$$

- training: gradient descent/Newton's method
- objective is convex

Artificial Neural Networks and Logistic Regression

- binary logistic regression:

- single *artificial neuron*

$$\text{single layer: } \hat{\mathbf{y}} = \sigma(\mathbf{W}^T \mathbf{x})$$

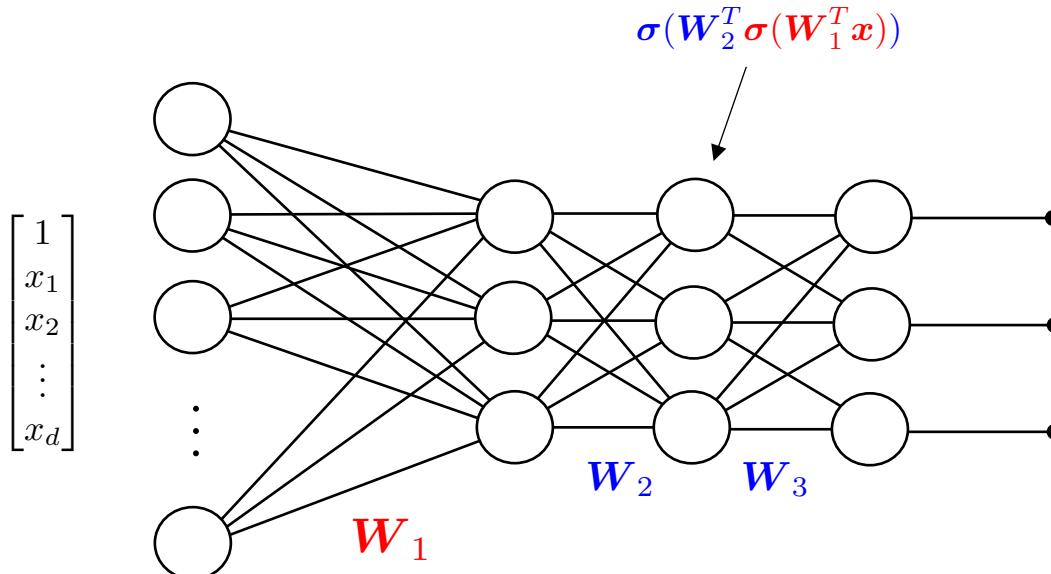
- multiclass logistic regression:

- single layer *artificial neural network*

$$\text{two layers: } \hat{\mathbf{y}} = \sigma(\mathbf{W}_2^T \sigma(\mathbf{W}_1^T \mathbf{x}))$$

$$\ell\text{-layers: } \hat{\mathbf{y}} = \sigma(\mathbf{W}_\ell^T \dots \sigma(\mathbf{W}_2^T \sigma(\mathbf{W}_1^T \mathbf{x})) \dots)$$

- *feedforward* neural network:



- softmax:

$$[\sigma(\mathbf{z})]_i = \frac{e^z_i}{\sum_j e^{z_j}}$$

$$\sigma(\mathbf{z}) = \frac{e^{\mathbf{z}}}{\sum_j e^{z_j}}$$

- logistic (unnormalized):

$$\sigma(z) = \frac{1}{1 + e^{-z}}$$

Training

- output

$$\ell\text{-layers: } \hat{\mathbf{y}} = \sigma(\mathbf{W}_\ell^T \dots \sigma(\mathbf{W}_2^T \sigma(\mathbf{W}_1^T \mathbf{x})) \dots)$$

- objective

$$\min_{\mathbf{W}_1, \mathbf{W}_2, \dots} \sum_i \ell(\hat{\mathbf{y}}_i, \mathbf{y}_i)$$

$$\min_{\mathbf{W}_1, \mathbf{W}_2, \dots} \sum_i D(\mathbf{y}_i || \sigma(\mathbf{W}_\ell^T \dots \sigma(\mathbf{W}_2^T \sigma(\mathbf{W}_1^T \mathbf{x}_i)) \dots))$$

- no longer convex
- we can still compute the gradient and head downhill
- stochastic gradient descent:

$$\sum_{i=1}^{100,000} \ell(y_i, \hat{g}_i) \rightarrow$$

$$1) \sum_{i=1}^{100} \ell(y_i, \hat{g}_i) \\ 2) \sum_{i=101}^{200} \ell(y_i, \hat{g}_i) \dots$$

Training

$$\sigma(z) = \frac{1}{1 + e^{-z}} \quad \frac{d\sigma(z)}{dz} = \sigma(z)(1 - \sigma(z))$$

- output

m -layers: $\hat{\mathbf{y}} = \sigma(\mathbf{W}_m^T \dots \sigma(\mathbf{W}_2^T \sigma(\mathbf{W}_1^T \mathbf{x})) \dots)$

- objective

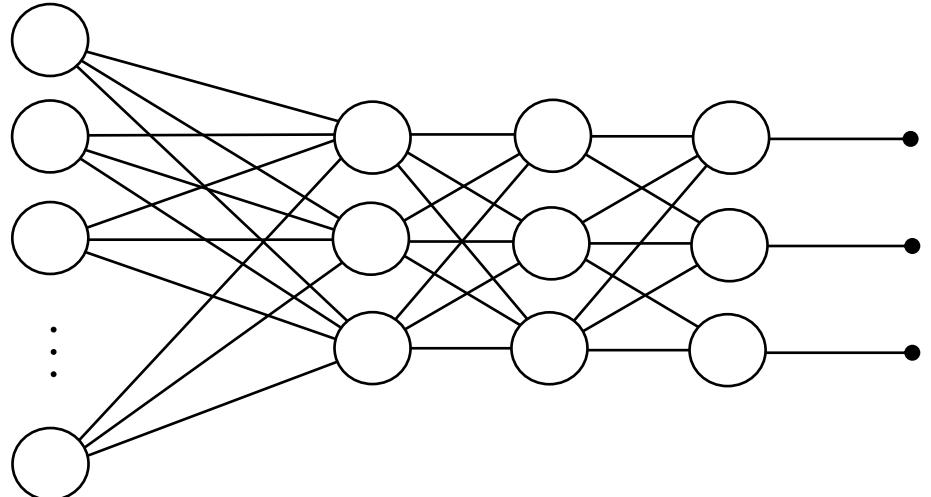
$$\min_{\mathbf{W}_1, \mathbf{W}_2, \dots} \sum_i \ell(\hat{\mathbf{y}}_i, \mathbf{y}_i)$$

- *backpropagation* - smart way to compute gradient

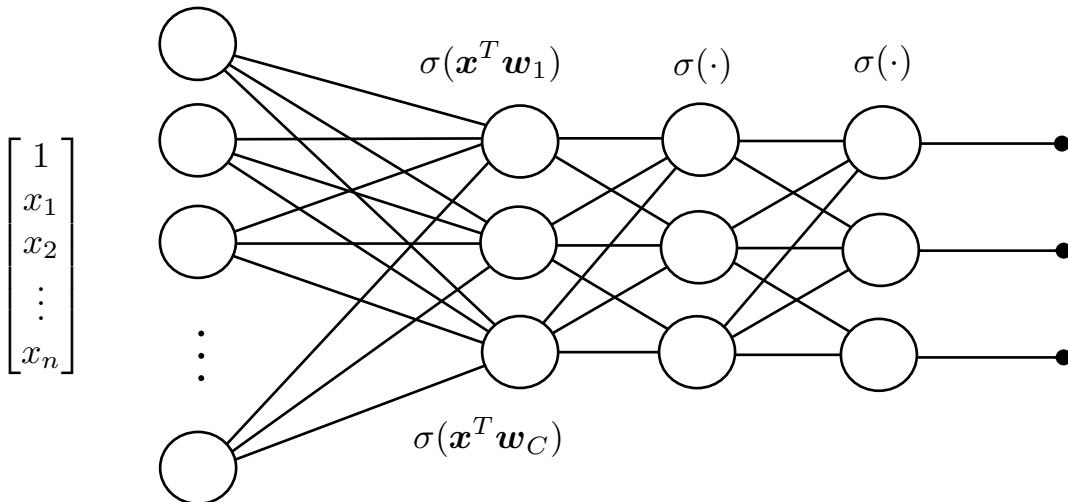
- 1) run \mathbf{x}_i through network
- 2) compute gradient with respect to weights in last layer

\vdots

$$\begin{bmatrix} 1 \\ x_1 \\ x_2 \\ \vdots \\ x_d \end{bmatrix}$$



Activation Functions, Loss functions



- activation functions
 - logistic
 - ReLU
 - linear
 - ...

- loss functions
 - cross entropy
 - KL divergence
 - squared error
 - ...

▷ Available activations
relu function
sigmoid function
softmax function
softplus function
softsign function
tanh function
selu function
elu function
exponential function

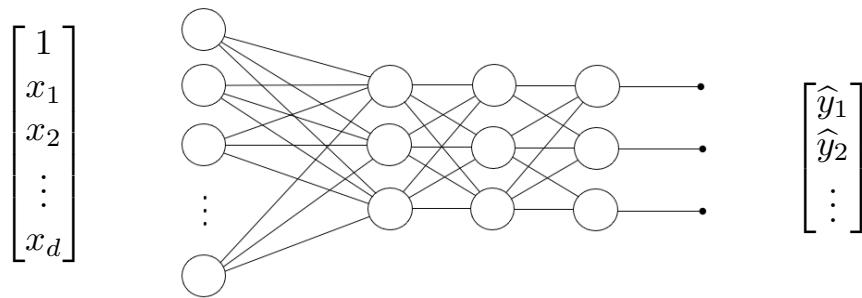
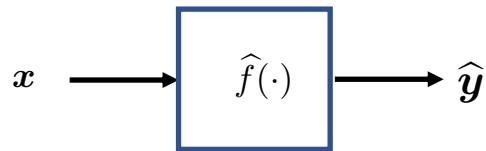
Probabilistic losses

- BinaryCrossentropy class
- CategoricalCrossentropy class
- SparseCategoricalCrossentropy class
- Poisson class
- binary_crossentropy function
- categorical_crossentropy function
- sparse_categorical_crossentropy function
- poisson function
- KLDivergence class
- kl_divergence function

Regression losses

- MeanSquaredError class
- MeanAbsoluteError class
- MeanAbsolutePercentageError class
- MeanSquaredLogarithmicError class
- CosineSimilarity class
- mean_squared_error function
- mean_absolute_error function
- mean_absolute_percentage_error function
- mean_squared_logarithmic_error function
- cosine_similarity function
- Huber class
- huber function
- LogCosh class
- log_cosh function

Function Approximation



- universal approximation theorem

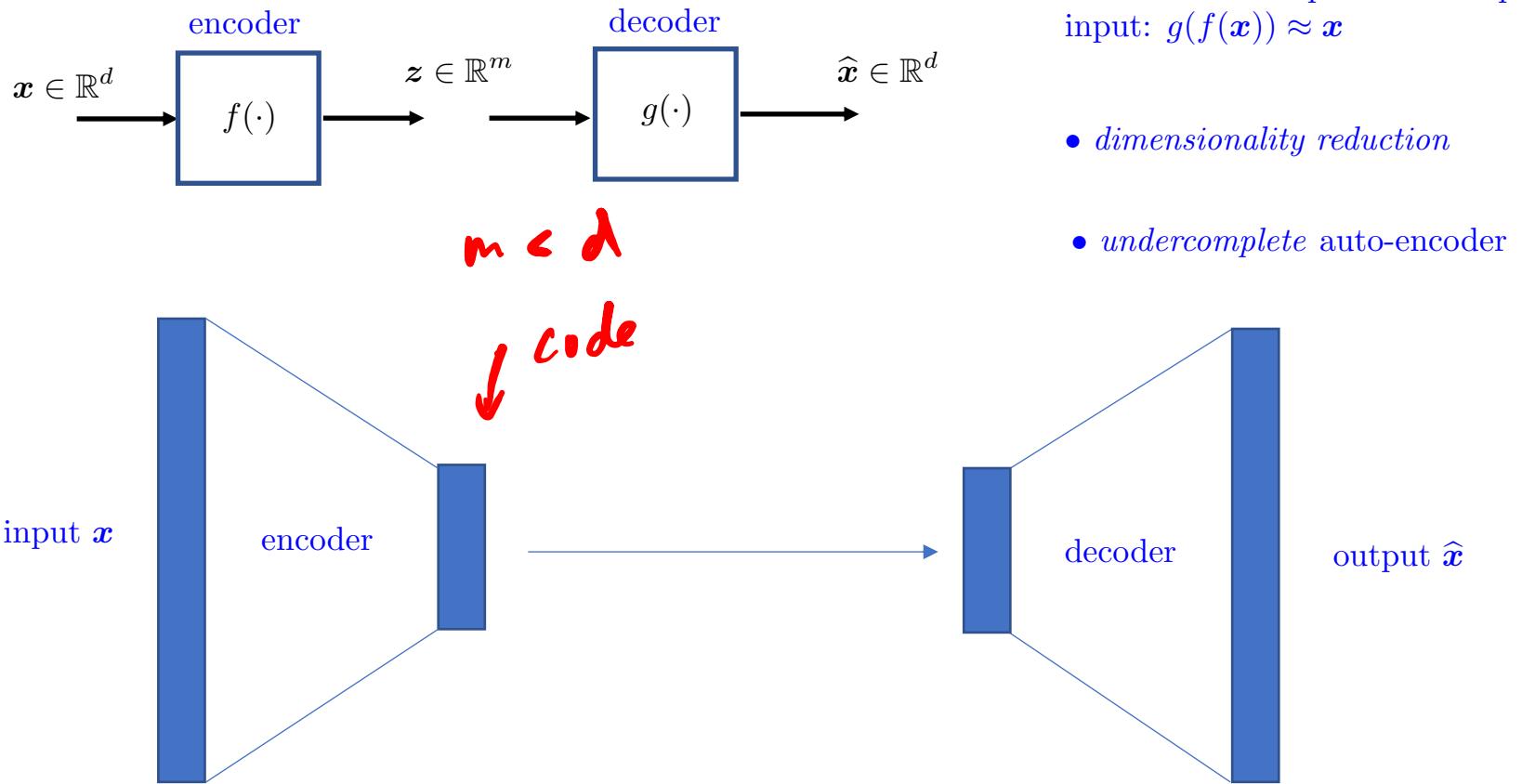
Consider a single layer neural network: $\hat{f}(\mathbf{x}) = \mathbf{W}_2^T \sigma(\mathbf{W}_1^T \mathbf{x})$.

For every continuous function $f : \mathbb{R}^d \rightarrow \mathbb{R}^D$ there exists a set of weights $\mathbf{W}_1, \mathbf{W}_2$ (*for some arbitrarily wide layer*) such that

$$\sup_{\mathbf{x}} ||f(\mathbf{x}) - \hat{f}(\mathbf{x})|| \leq \epsilon$$

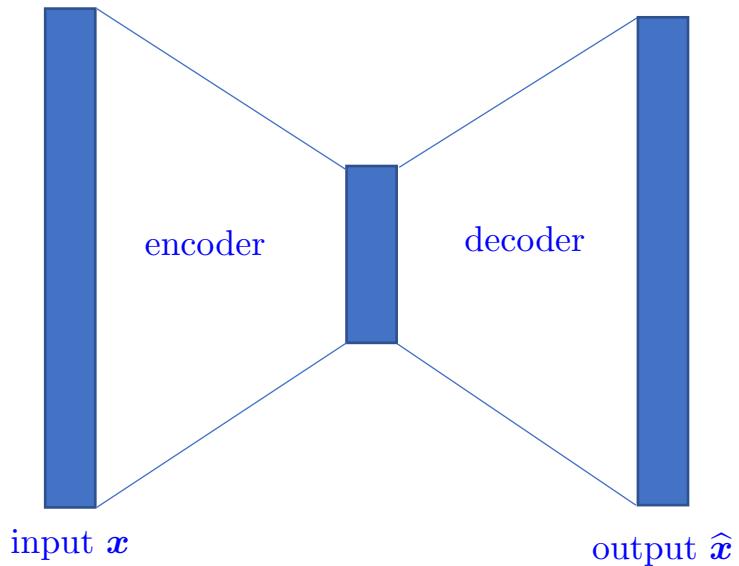
for any $\epsilon > 0$.

Autoencoders



Autoencoders

- autoencoders are lossy
- autoencoders are specific to the types of examples they are trained on
- trained from data as opposed to engineered by a human
- autoencoders are (in general) **not** better than engineered compression algorithms
- main applications: denoising, data visualization, manifold learning



Autoencoders

- training undercomplete autoencoders:

$$\min \sum_i \ell(x_i, \hat{x}_i)$$

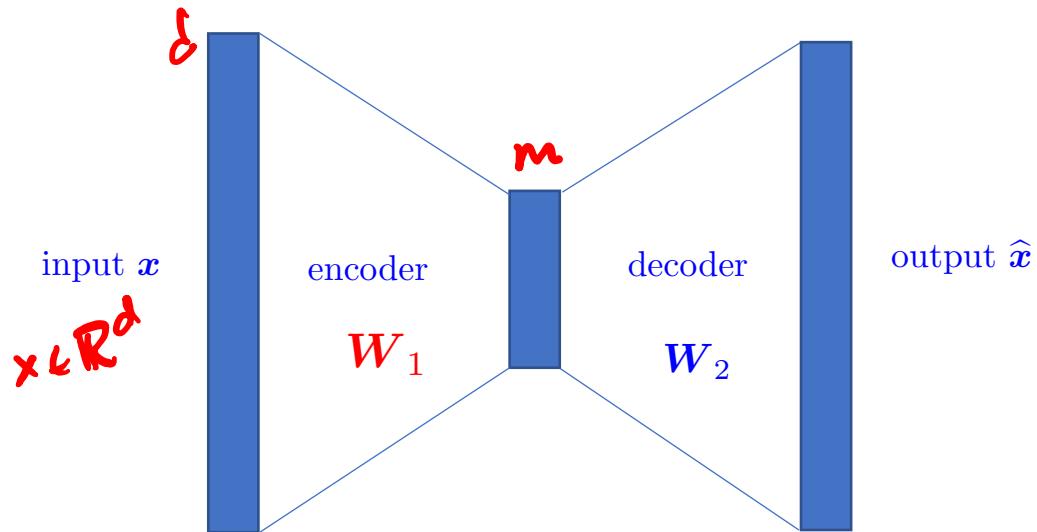
- use SGD and backpropagation

- example: linear activation functions + mean square error loss

$$\hat{x} = \sigma(W_2^T \sigma(W_1^T x))$$

$$\hat{x} = \delta \begin{bmatrix} m \\ W_2 \end{bmatrix} \begin{bmatrix} \delta \\ W_1^T \end{bmatrix} \begin{bmatrix} \delta \\ x \end{bmatrix}$$

$$\tau(z) = z$$



$$\begin{aligned} \min_{W_1, W_2} \sum_i \ell(x_i, \hat{x}_i) &= \min_{W_1, W_2} \sum_i \|x_i - W_2^T W_1^T x_i\|_2^2 \\ &= \min_{W_1, W_2} \|X - \underbrace{W_2^T W_1^T X}_A\|_F^2 \\ &= \min_{A: \text{rank}(A) \leq m} \|X - A\|_F^2 \end{aligned}$$

$$A = \sum_{i=1}^m \tau_i \underline{u_i} \underline{v_i}^\top \quad \text{PLA}$$

Autoencoders

- undercomplete autoencoders

- regularized autoencoders

- sparse autoencoders

$$\min \sum_i \ell(\mathbf{x}_i, \hat{\mathbf{x}}_i) + \lambda \|\mathbf{z}\|_1$$

- de-noising autoencoders

input $\mathbf{x} + \text{noise}$

