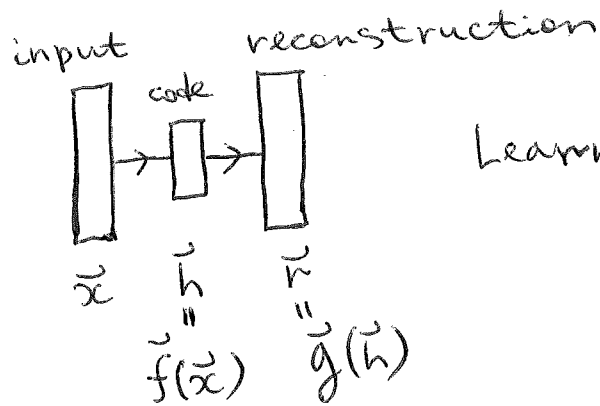


Lecture 21 Autoencoders



minimize loss function $L(\vec{x}, g(f(\vec{x})))$

Learn $g(f(\vec{x})) = \vec{x}$

copies input to output

We're interested in the code \vec{h} :
dimensionality reduction, feature learning

We need to avoid overfitting (i.e., just copying input to output):

use undercomplete autoencoder (code dim'n $<$ input dim'n)

OR overcomplete autoencoder ($D_c > D_i$)

+ regularization for sparsity of the code, robustness to noise, robustness to translations/rotations, etc.

Sparse autoencoders: minimize $L(\vec{x}, g(f(\vec{x}))) + \lambda R(\vec{h})$

could be used e.g. for feature learning which then provide input to classification. For example, one could use ReLU activation units

in the code layer \oplus $\mathcal{L}(\vec{h}) = |\vec{h}|$
 (i.e., to encode \vec{h}) L_1 penalty.

This should produce \emptyset 's in the \vec{h} code.

Denoising autoencoders: minimize

$$L(\vec{x}, \vec{g}(\vec{f}(\vec{x}'))), \text{ where}$$

$$\vec{x}' = \vec{x} + \underbrace{\vec{E}}_{\text{noise}} \leftarrow \text{perturbed input}$$

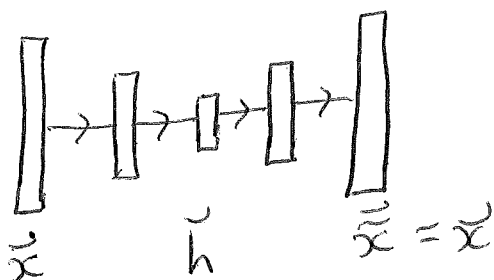
This forces autoencoders to undo the noise corruption.

Contractive autoencoders: minimize

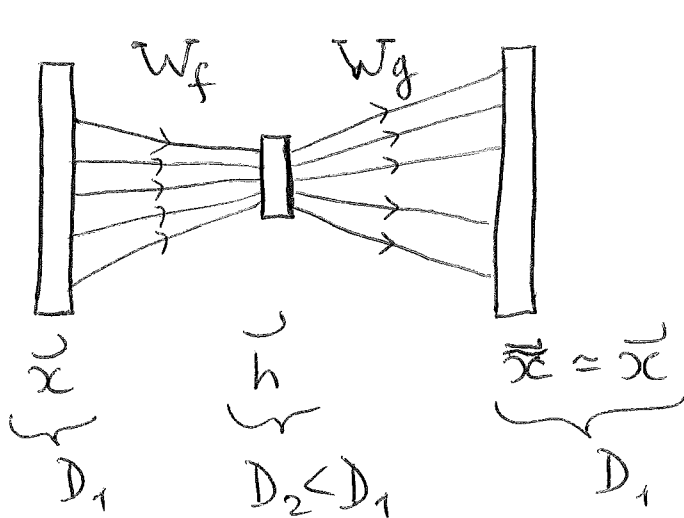
$$L(\vec{x}, \vec{g}(\vec{f}(\vec{x}))) + \lambda \sum_i \underbrace{\|\nabla_{\vec{x}} h_i\|^2}_{\mathcal{L}(\vec{h}, \vec{x})}$$

$h_i = f_i(\vec{x})$

Autoencoders \rightarrow shallow (as discussed thus far)
 \rightarrow deep (multiple layers to represent \vec{f} & \vec{g}):



Ex.: linear autoencoder



$$\vec{h} = W_f \vec{x}$$

$$\vec{\hat{x}} = W_g \vec{h}$$

$$\min_{W_f, W_g} \left\{ \frac{1}{2} \sum_n | \underbrace{W_g W_f \vec{x}_n}_{\text{datapoint}} - \vec{x}_n |^2 \right\}$$

$L = \text{loss function}$

Equiv. to PCA (principal component analysis)

non-linear autoencoder

$$\min_{W_f, W_g} \left\{ \frac{1}{2} \sum_n | \vec{g}(\vec{f}(\vec{x}_n; W_f); W_g) - \vec{x}_n |^2 \right\}$$

↑ NN weights + biases

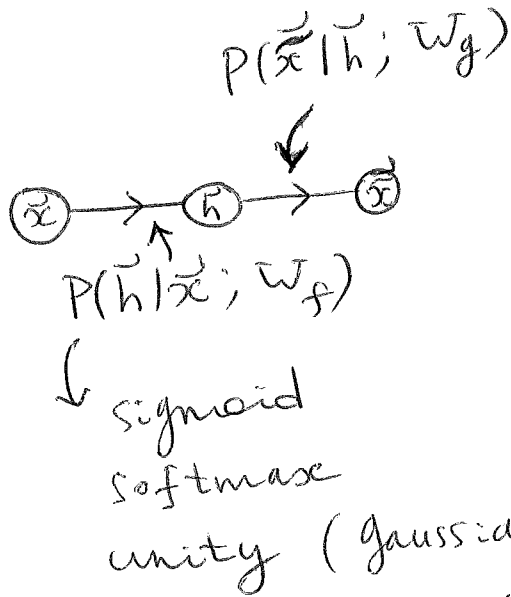
[Projection onto lower-D non-linear manifold.]

Probabilistic autoencoders

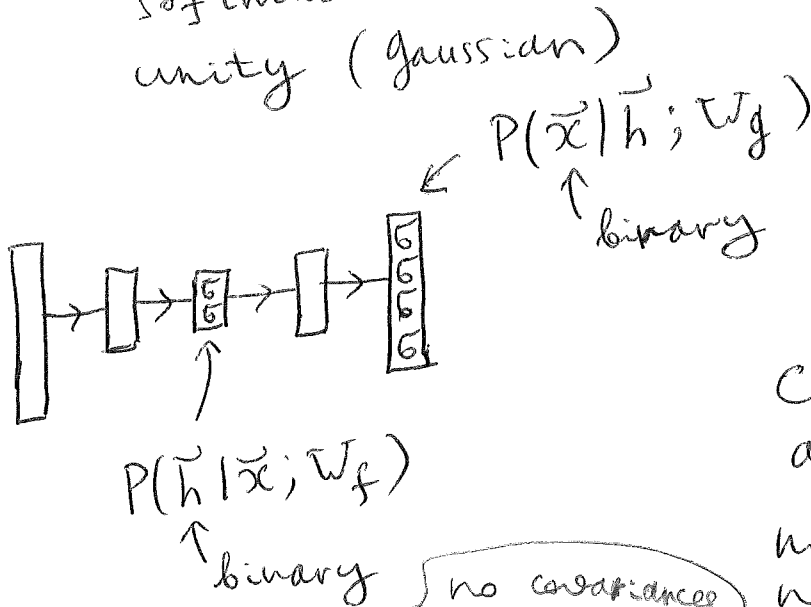
f: $\Pr(\vec{h} | \vec{x}; W_f)$

g: $\Pr(\vec{x} | \vec{h}; W_g)$

use sigmoid, softmax or linear units at the code & output layers to produce gaussian μ, σ
 ↗ discard "-" & negative



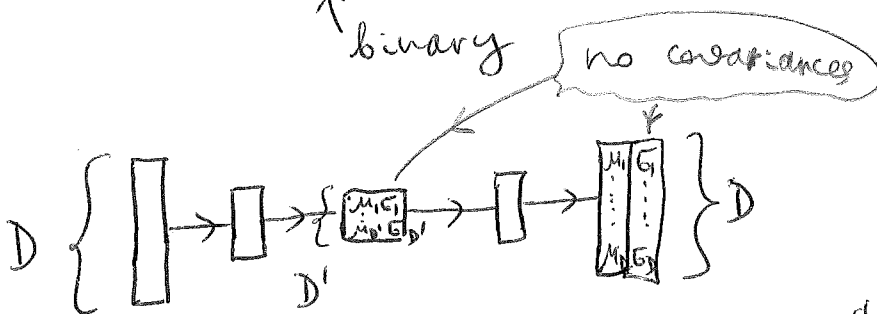
$\vec{x}, \vec{x}, \vec{h}$ are random vars



Can use these as generative models to produce novel outputs:

1. Sample \vec{h} from $P(\vec{h})$

2. Sample \vec{x} from decoder: $P(\vec{x} | \vec{h}; W_g)$



[Denoising autoencoders:
(detailed discussion)]

1. Sample \vec{x} from training data
2. Sample noisy/corrupted \vec{x}' : ^(version of)

$$\vec{x}' = \vec{x} + \vec{\epsilon} \Rightarrow P(\vec{x}' | \vec{x})$$

3. Use (\vec{x}, \vec{x}') to construct $P_{\text{decoder}}(\vec{x} | \vec{h} = f(\vec{x}')) \equiv P_{\text{reconstruct}}(\vec{x} | \vec{x}')$

↳ minimize

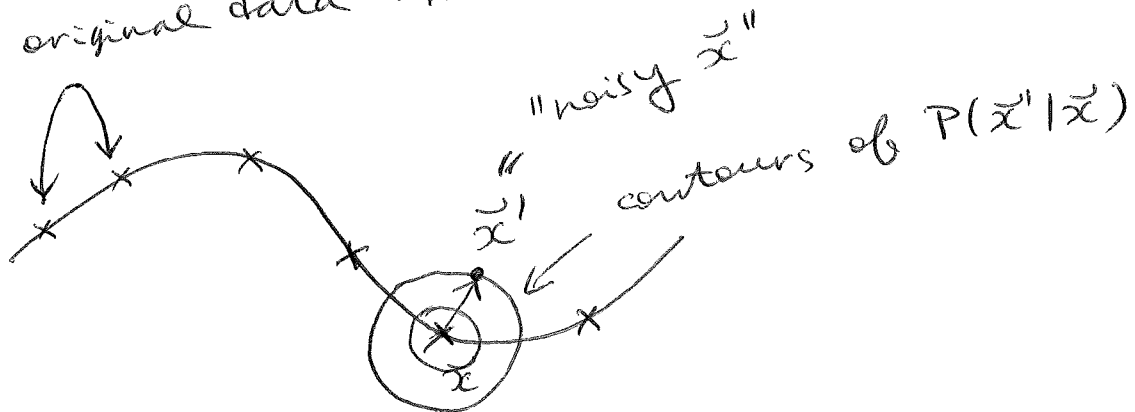
$$-\log P_{\text{decoder}}(\vec{x} | \vec{h})$$

The error function is given by

$$E = - \int d\vec{x} \int d\vec{x}' p(\vec{x}) P(\vec{x}' | \vec{x}) \log P_{\text{decoder}}(\vec{x} | \vec{h} = f(\vec{x}'))$$

"summed" over many datapoints \vec{x} & , for each \vec{x} , many corrupted/noisy \vec{x}'

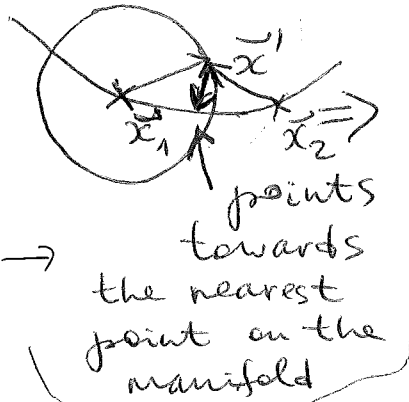
original data on low-D manifold



Now, consider quadratic error:

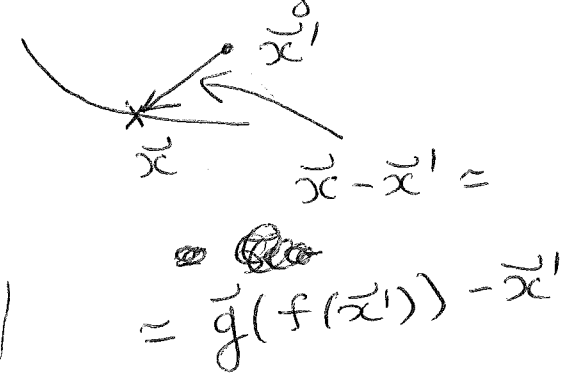
$$E = \int d\vec{x} \int d\vec{x}' p(\vec{x}) P(\vec{x}'|\vec{x}) \underbrace{|\vec{g}(f(\vec{x}')) - \vec{x}|^2}$$

learns to be as close to \vec{x} as possible, by averaging over many \vec{x}' realizations



averaged over many \vec{x} (same \vec{x}'), requires local linearity

points towards the nearest point on the manifold



local linearity of the "true" manifold

represents a gradient flow towards the "true" \vec{x} -manifold