

# Day 9: Unsupervised learning, dimensionality reduction

**Introduction to Machine Learning Summer School**  
**June 18, 2018 - June 29, 2018, Chicago**

Instructor: Suriya Gunasekar, TTI Chicago

28 June 2018



THE UNIVERSITY OF  
CHICAGO



# Topics so far

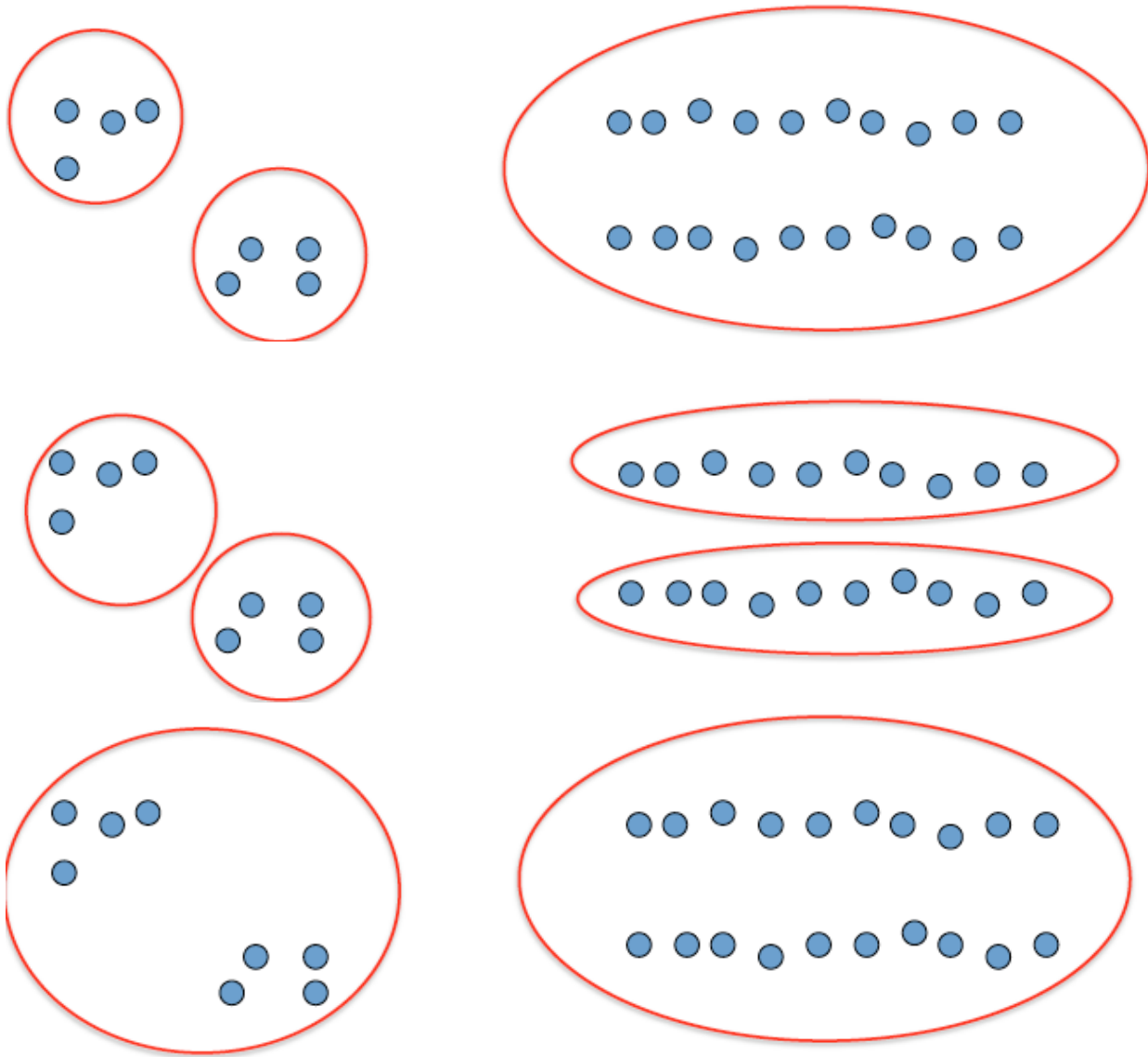
- Linear regression
- Classification
  - Logistic regression
  - Maximum margin classifiers, kernel trick
  - Generative models
  - Neural networks
  - Ensemble methods
- Today and Tomorrow
  - Unsupervised learning – dimensionality reduction, clustering
  - Review

# Unsupervised learning

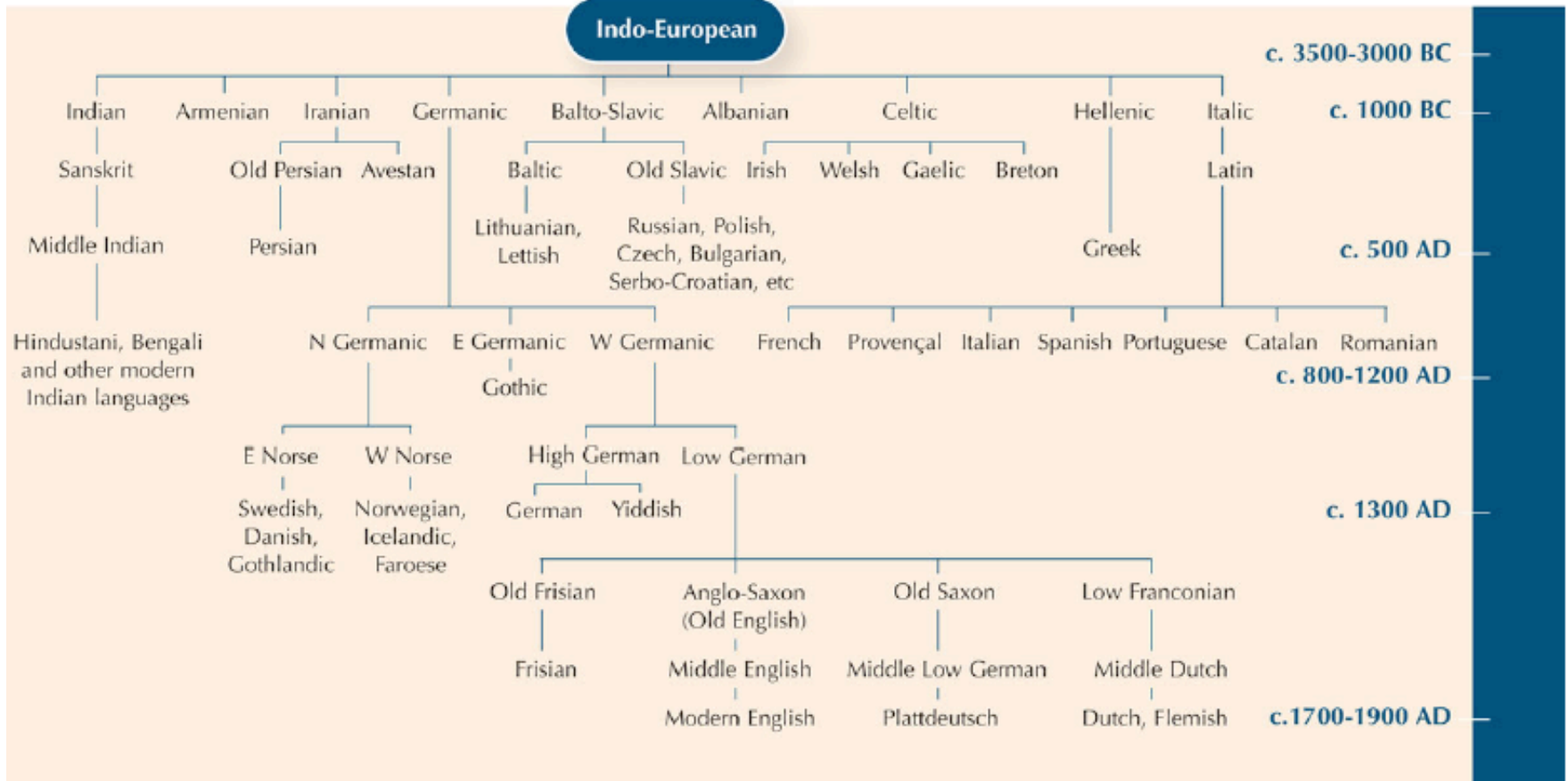
- **Unsupervised learning:**  
Requires data  $x \in \mathcal{X}$ , but no labels
- **Goal?:** Compact representation of the data by detecting patterns
  - e.g. Group emails by topic
- **Useful when we don't know what we are looking for**
  - makes evaluation tricky
- **Applications in visualization, exploratory data analysis, semi-supervised learning**



# Clustering

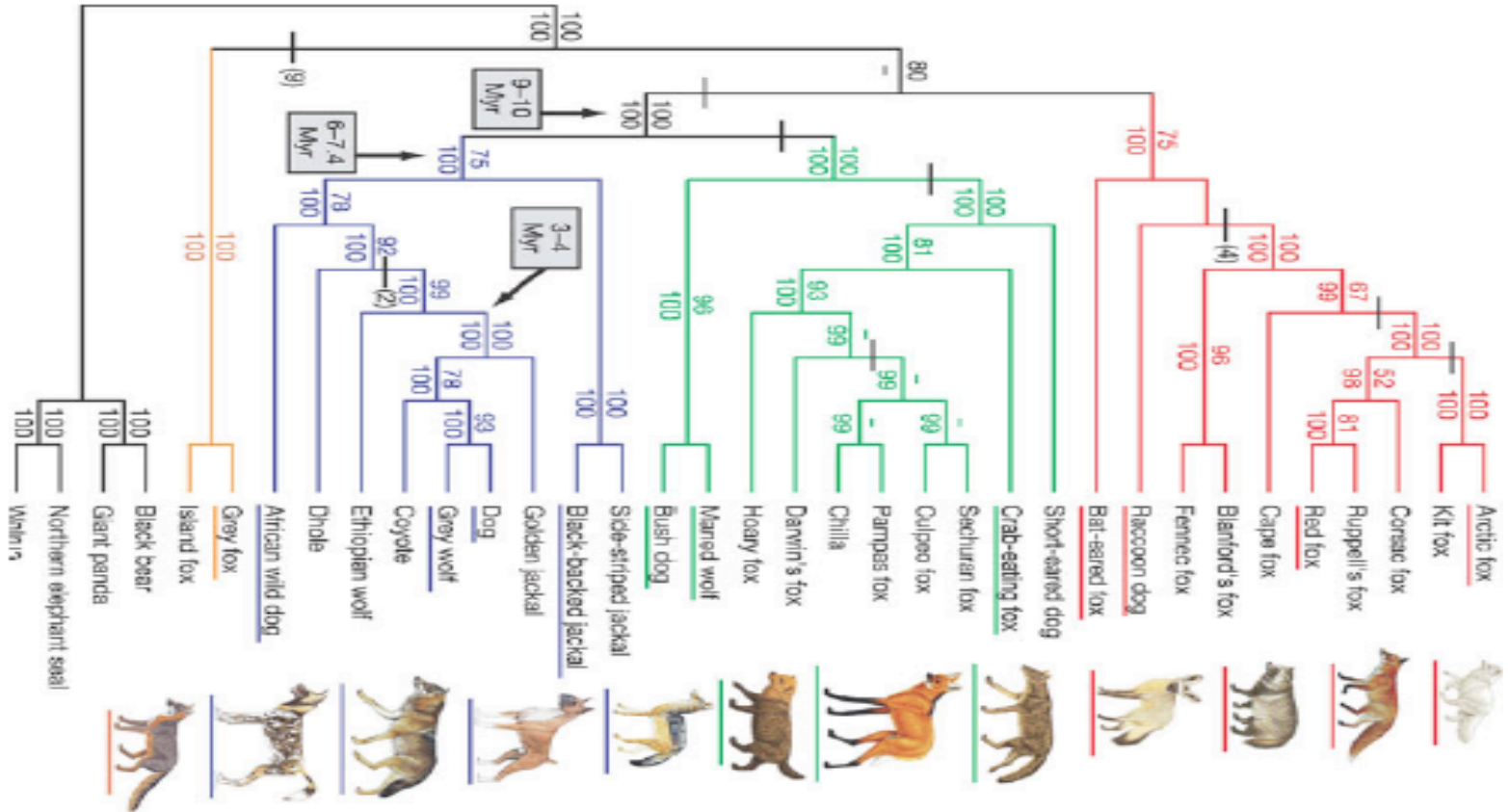


# Clustering languages

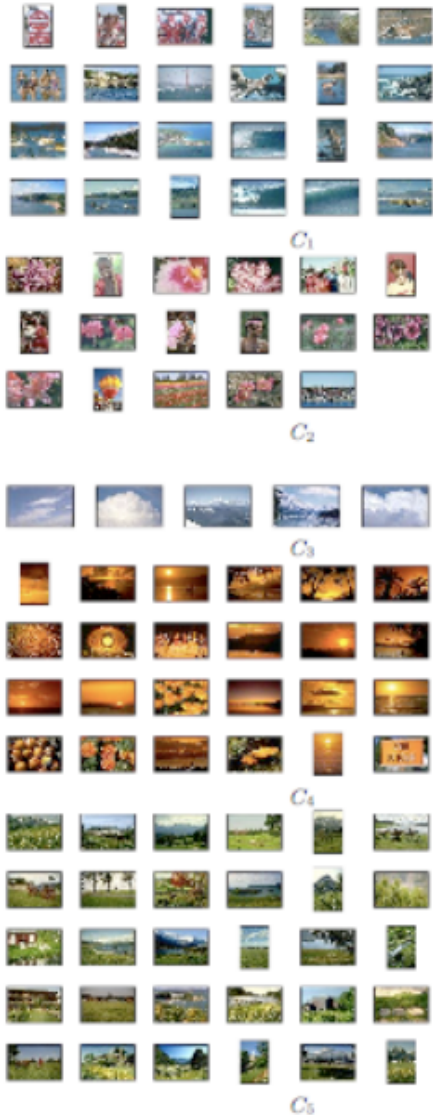


[Image from scienceinschool.org]

# Clustering species (phylogeny)



# Image clustering/segmentation



[Goldberger et al.]

Goal: Break up the image into meaningful or perceptually similar regions



[Slide from James Hayes]

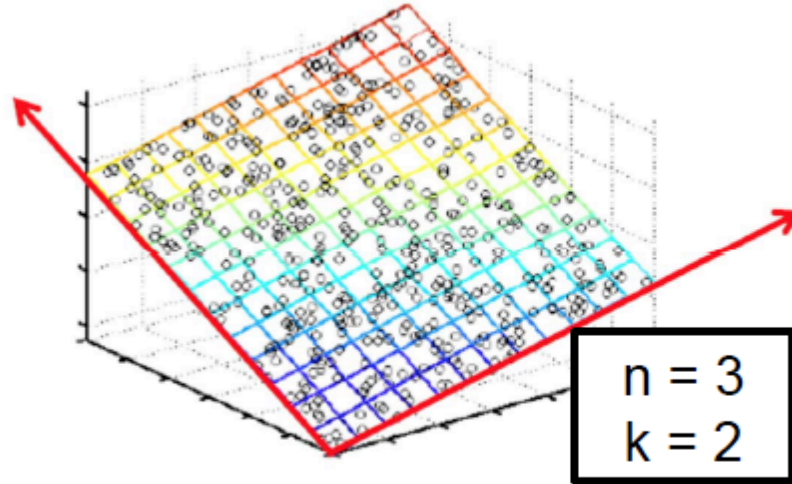
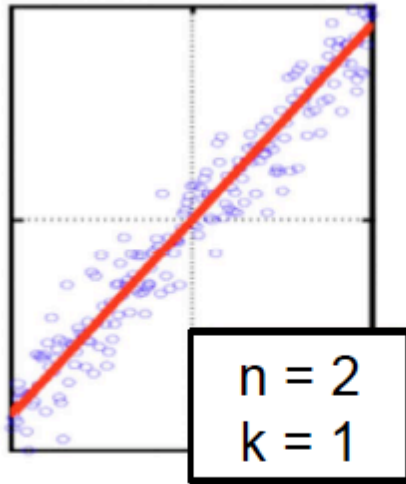
Current trend is  
to use datasets  
with labels for  
such task  
e.g., MS COCO

# Dimensionality reduction

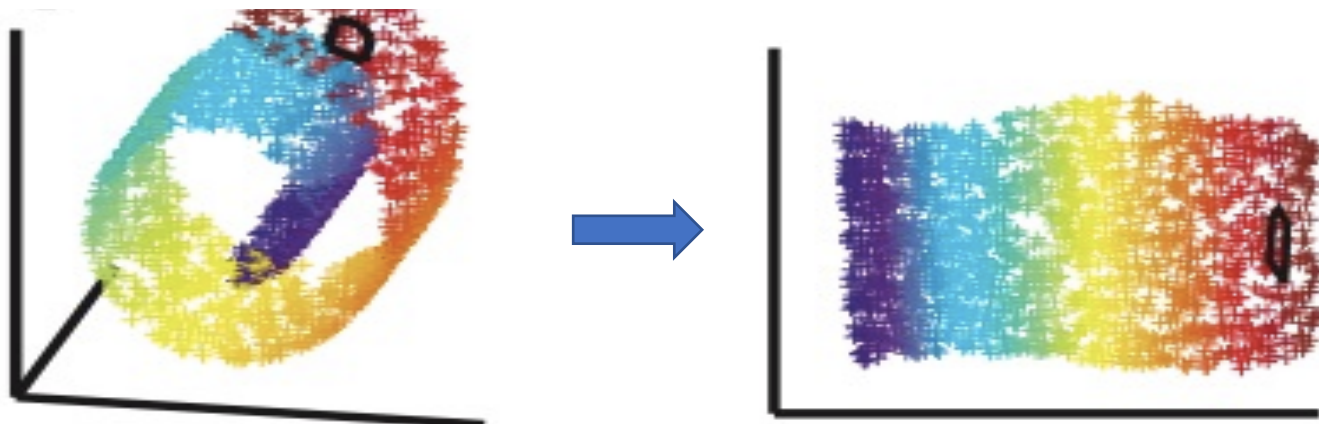
- Input data  $x \in \mathcal{X}$  may have thousands or millions of dimensions!
  - e.g., text data represented as bag of words
  - e.g., video stream of images
  - e.g., fMRI data  $\# \text{voxels} \times \# \text{timesteps}$
- Dimensionality reduction: represent data with fewer dimensions
  - easier learning in subsequent tasks (preprocessing)
  - visualization
  - discover intrinsic patterns in the data



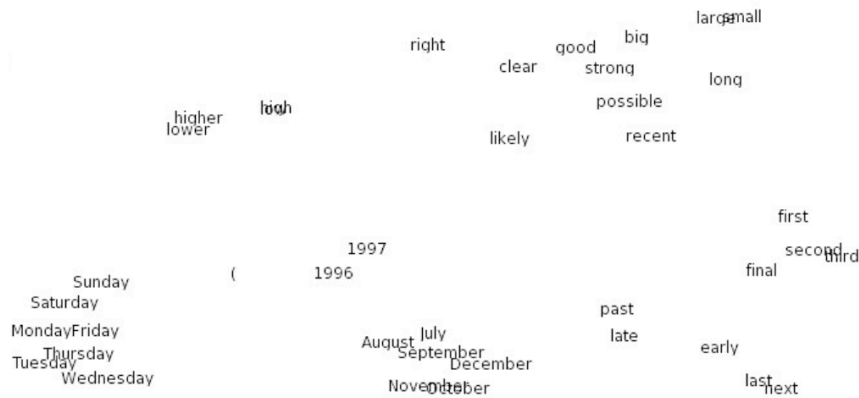
# Manifolds



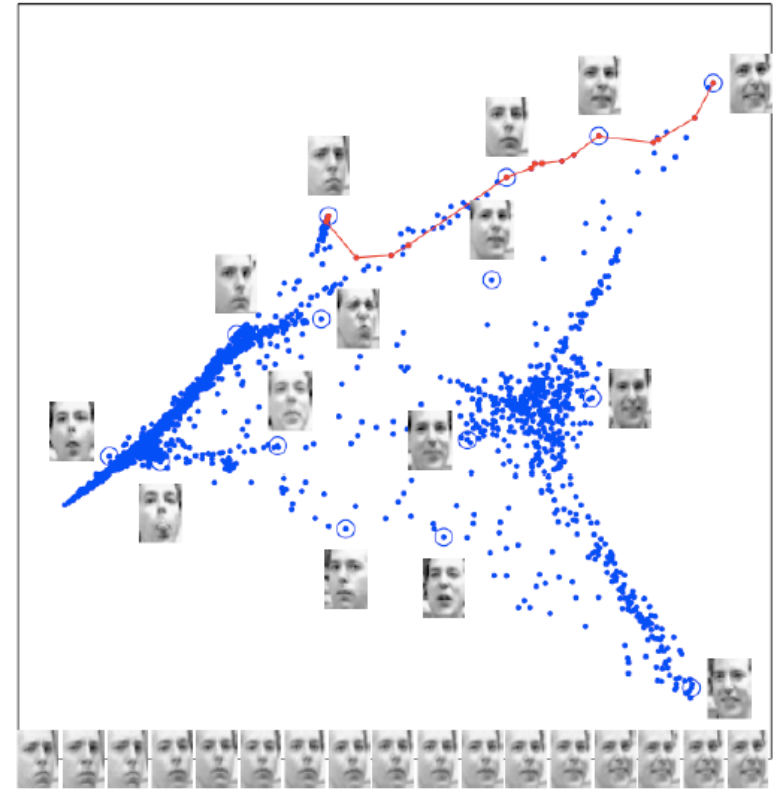
Slide from Yi Zhang



# Embeddings



t-SNE visualization from Turian et al. (2010)



# Low dimensional embedding

- Given high dimensional feature

$$\mathbf{x} = [x_1, x_2, \dots, x_d]$$

find transformations

$$z(\mathbf{x}) = [z_1(\mathbf{x}), z_2(\mathbf{x}), \dots, z_k(\mathbf{x})]$$

so that “almost all useful information” about  $\mathbf{x}$  is retained in  $z(\mathbf{x})$

- In general  $k \ll d$ , and  $z(\mathbf{x})$  is not invertible
- Transformation learned from a dataset of examples of  $x$

$$S = \{\mathbf{x}^{(i)} \in \mathbb{R}^d : i = 1, 2, \dots, N\}$$

- Note: typically no labels  $y$

# Linear dimensionality reduction

- Given high dimensional feature

$$\mathbf{x} = [x_1, x_2, \dots, x_d]$$

find transformations

$$\mathbf{z} = z(\mathbf{x}) = [z_1(\mathbf{x}), z_2(\mathbf{x}), \dots, z_k(\mathbf{x})]$$

- Restrict  $z(\mathbf{x})$  to be a linear function of  $\mathbf{x}$

$$z_1 = \mathbf{w}_1 \cdot \mathbf{x}$$

$$z_2 = \mathbf{w}_2 \cdot \mathbf{x}$$

$\vdots$

$$z_k = \mathbf{w}_k \cdot \mathbf{x}$$

$$\begin{array}{c} z_1 \\ z_2 \\ \vdots \\ z_k \end{array} = \begin{array}{c} \leftarrow \mathbf{w}^{(1)} \rightarrow \\ \leftarrow \mathbf{w}^{(2)} \rightarrow \\ \vdots \\ \leftarrow \mathbf{w}^{(k)} \rightarrow \end{array} \begin{array}{c} x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_d \end{array}$$

$$\mathbf{z} = \mathbf{W}\mathbf{x}$$

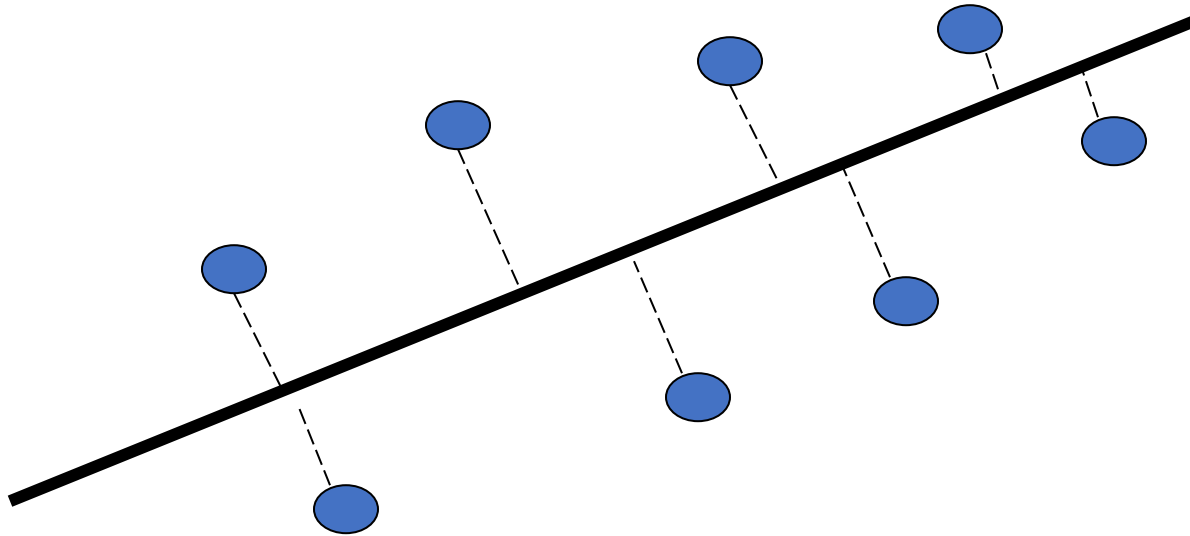
where

$$\mathbf{z} \in \mathbb{R}^k,$$
$$\mathbf{W} \in \mathbb{R}^{k \times d},$$
$$\mathbf{x} \in \mathbb{R}^d$$

only question  
is which  $\mathbf{W}$ ?

# Linear dimensionality 2D example

- Given points  $S = \{\mathbf{x}^{(i)} : i = 1, 2, \dots, N\}$  in 2D, we want a 1D representation
  - project  $\{\mathbf{x}^{(i)}\}$  onto a line  $\mathbf{w} \cdot \mathbf{x} = 0$

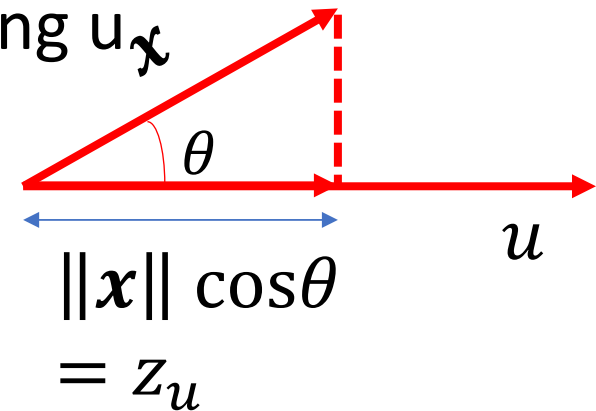


- Find  $\mathbf{w}$  to minimize the sum of squared distances to the line

# Vector projections

- $\mathbf{x} \cdot \mathbf{u} = \|\mathbf{x}\| \|\mathbf{u}\| \cos \theta$
- Assuming  $\|\mathbf{u}\| = 1$ ,
- $\mathbf{x} \cdot \mathbf{u} = \|\mathbf{x}\| \cos \theta = z_u \rightarrow$  value of  $x$  along  $u$

- distance of  $\mathbf{x}$  to projection is  
 $\|z_u \mathbf{u} - \mathbf{x}\| = \|(\mathbf{x} \cdot \mathbf{u})\mathbf{u} - \mathbf{x}\|$



# Principal component analysis

- For a 1D embedding along direction  $\mathbf{u}$ , distance of  $\mathbf{x}$  to the projection along  $\mathbf{u}$  is given by

$$\|z_{\mathbf{u}} \mathbf{u} - \mathbf{x}\| = \|(\mathbf{x} \cdot \mathbf{u})\mathbf{u} - \mathbf{x}\|$$

- More generally for  $k$  dimensional embedding:
  - find orthonormal basis of the  $k$  dimensional subspace  
 $\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_k \in \mathbb{R}^d$ , i.e.,  $\mathbf{u}_i \cdot \mathbf{u}_j = 1$  if  $i = j$ , and 0 otherwise
  - let  $\mathbf{U} \in \mathbb{R}^{k \times d}$  be the matrix with  $\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_k$  along rows
  - distance of projection of  $\mathbf{x}$  to  $\text{span}\{\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_k\}$   
 $\|\mathbf{U}^T \mathbf{U} \mathbf{x} - \mathbf{x}\|$
  - also from orthonormality of  $\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_k$ , check  $\mathbf{U} \mathbf{U}^T = \mathbf{I}$

- PCA objective

$$\min_{\mathbf{U} \in \mathbb{R}^{k \times d}} \sum_{i=1}^N \|\mathbf{U}^T \mathbf{U} \mathbf{x}^{(i)} - \mathbf{x}^{(i)}\|^2 \quad \text{s.t.} \quad \mathbf{U} \mathbf{U}^T = \mathbf{I}$$

# PCA

- PCA objective

$$\min_{U \in \mathbb{R}^{k \times d}} \frac{1}{N} \sum_{i=1}^N \|U^T U x^{(i)} - x^{(i)}\|^2 \quad s.t. \quad UU^T = I$$

- Also, for all  $UU^T = I$

$$\begin{aligned} \|U^T U x - x\|^2 &= \|x\|^2 + x^T U^T U U^T U x - 2x^T U^T U x \\ &= \|x\|^2 - x^T U^T U x = \|x\|^2 - \|Ux\|^2 \end{aligned}$$

- Equivalent PCA objective

$$\max_U \frac{1}{N} \sum_{i=1}^N \|Ux^{(i)}\|^2 = \sum_{j \in [k]} u_j^T \hat{\Sigma}_{xx} u_j \quad s.t. \quad UU^T = I$$

where  $\hat{\Sigma}_{xx} = \frac{1}{N} \sum_{i=1}^N x^{(i)} x^{(i)T}$  (derivation in board)

- This is the same as finding top  $k$  eigenvectors of  $\hat{\Sigma}_{xx}$

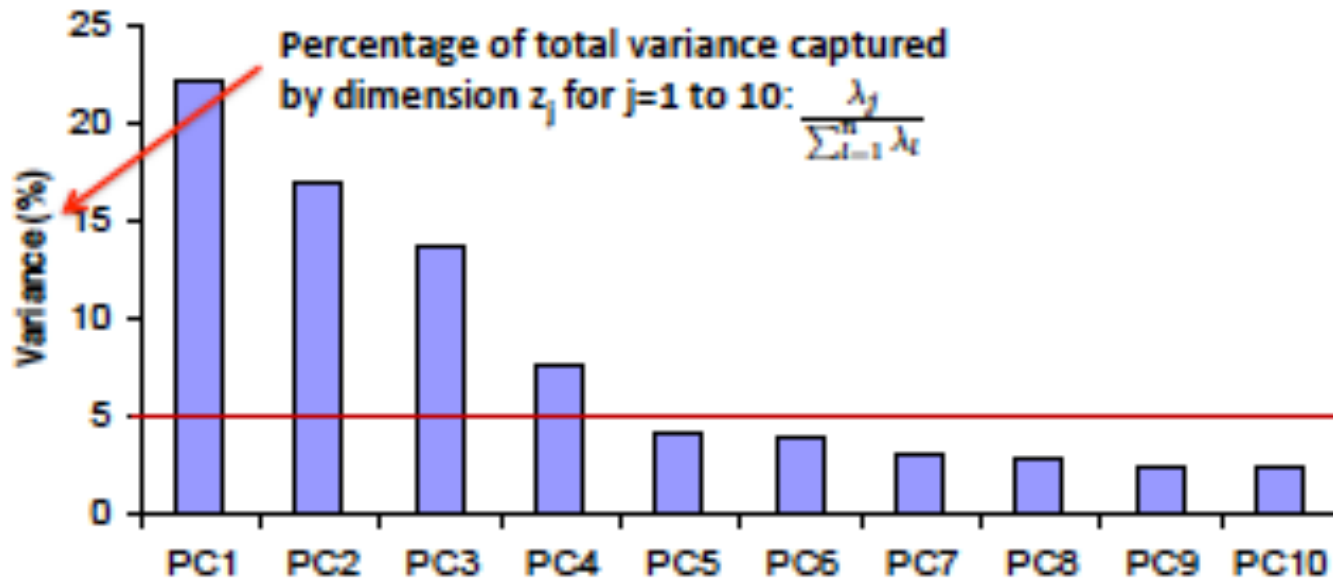


# PCA algorithm

- Given  $S = \{\mathbf{x}^{(i)} \in \mathbb{R}^d : i = 1, 2, \dots, N\}$
- Let  $\mathbf{X} \in \mathbb{R}^{N \times d}$  be data matrix
  - make sure  $\mathbf{X}$  is re-centered so that column mean is 0
- $\hat{\Sigma}_{xx} = \frac{1}{N} \sum_{i=1}^N \mathbf{x}^{(i)} \mathbf{x}^{(i)\top} = \frac{1}{N} \mathbf{X}^\top \mathbf{X} \in \mathbb{R}^{d \times d}$
- $\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_k \in \mathbb{R}^d$  are top  $k$  eigenvectors of  $\hat{\Sigma}_{xx}$

# How to pick $k$ ?

- Data assumed to be low dimensional projection + noise
- Only keep projections onto components with large eigenvalues and ignore the rest

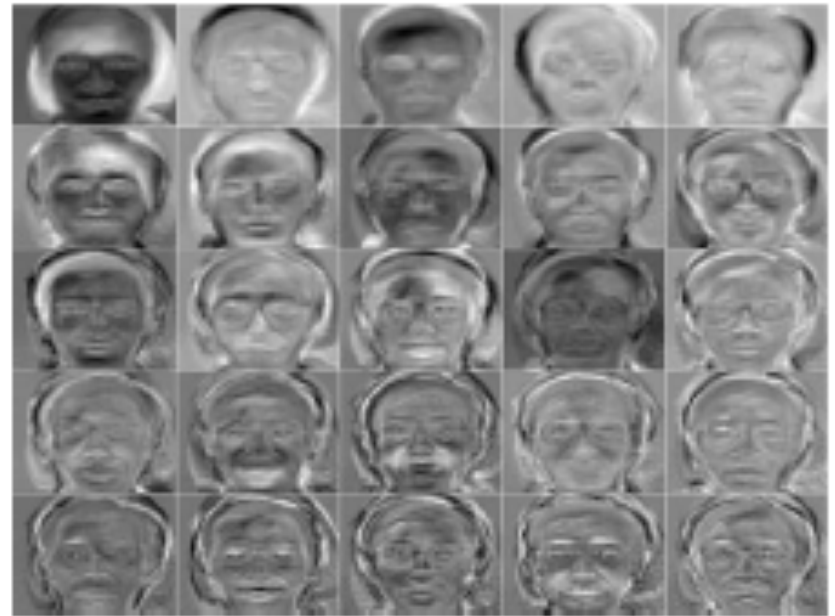


# Eigenfaces

- Input images:



- Principal components:



- Turk and Pentland '91

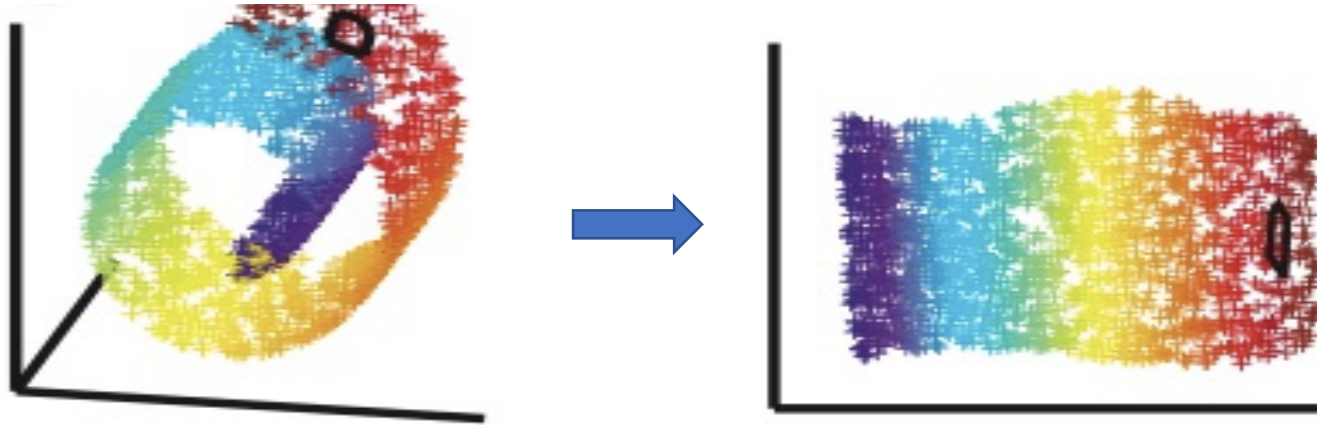
# SVD version

- Given  $S = \{\mathbf{x}^{(i)} \in \mathbb{R}^d : i = 1, 2, \dots, N\}$
- Let  $\mathbf{X} \in \mathbb{R}^{N \times d}$  be data matrix
  - make sure  $\mathbf{X}$  is re-centered so that column mean is 0
- $\mathbf{X} = \bar{\mathbf{V}}\bar{\mathbf{S}}\bar{\mathbf{U}}^\top$  be the **Singular Value Decomposition (SVD)** of  $\mathbf{X}$ , where
  - $\bar{\mathbf{V}} \in \mathbb{R}^{N \times d}$  have orthonormal columns, i.e.,  $\bar{\mathbf{V}}^\top \bar{\mathbf{V}} = \mathbf{I}$ 
    - columns of  $\bar{\mathbf{V}}$  are called left singular vectors
  - $\bar{\mathbf{U}} \in \mathbb{R}^{d \times d}$  also has orthonormal columns, i.e.,  $\bar{\mathbf{U}}^\top \bar{\mathbf{U}} = \mathbf{I}$ 
    - columns of  $\bar{\mathbf{U}}$  are called right singular vectors
  - $\bar{\mathbf{S}} = \text{diagonal}(\sigma_1, \sigma_2, \dots, \sigma_d) \in \mathbb{R}^{d \times d}$ 
    - $\sigma_1, \sigma_2, \dots, \sigma_d$  are called the singular values
- First  $k$  columns of  $\bar{\mathbf{U}}$  are the  $\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_k$  we want.
- Representation of  $\mathbf{x} \in \mathbb{R}^d$  as  $z(\mathbf{x}) \in \mathbb{R}^k$  is given by
$$z(\mathbf{x})_j = \sigma_j \mathbf{u}_j \cdot \mathbf{x} \quad \text{for } j = 1, 2, \dots, k$$

# Other linear dimensionality reduction

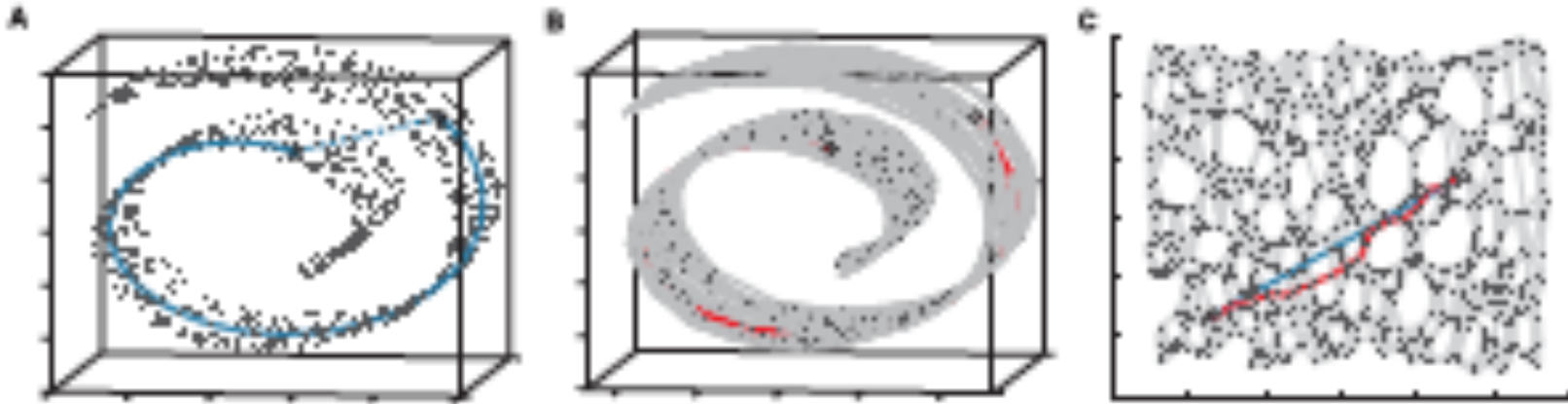
- **PCA:** given data  $x \in \mathbb{R}^d$ , find  $U \in \mathbb{R}^{k \times d}$  to minimize
$$\min_U \|U^T U x - x\|_2^2 \text{ s.t. } U U^T = I$$
- **Canonical correlation analysis:** given two “views” of data  $x \in \mathbb{R}^d$  and  $x' \in \mathbb{R}^{d'}$ , find  $U \in \mathbb{R}^{k \times d}$ ,  $U' \in \mathbb{R}^{k \times d'}$  to minimize
$$\|U x - U' x'\|_2^2 \text{ s.t. } U U^T = U' U'^T = I$$
- **Sparse dictionary learning:** learn a sparse representation of  $x$  as a linear combination of over-complete dictionary
$$x \rightarrow D z \text{ where } D \in \mathbb{R}^{d \times m}, z \in \mathbb{R}^m$$
  - unlike PCA, here  $m \gg d$  so  $z$  is higher dimensional, but learned to be sparse!
- **Independent component analysis**
- **Factor analysis**
- **Linear discriminant analysis**

# Non linear dimensionality reduction



- **Isomap**
- **Autoencoders**
- Kernel PCA
- Local linear embedding
- Check out t-SNE for 2D visualization
- ...

# Isomap



[Tenenbaum, Silva, Langford. Science 2000]

# Isomap – algorithm

- Dataset of  $N$  points  $S = \{\mathbf{x}^{(i)} \in \mathbb{R}^d : i = 1, 2, \dots, N\}$
- Represent the points as a kNN-graph with weights proportional to distance between the points
- The geodesic distance  $d(x, x')$  between points in the manifold is the length of shortest path in the graph
- Use any shortest path algorithm can be used to construct a matrix  $M \in \mathbb{R}^{N \times N}$  of  $d(x^{(i)}, x^{(j)})$  for all  $x^{(i)}, x^{(j)} \in S$
- **MDS**: Find a (low dimensional) embedding  $z(x)$  of  $x$  so that distances are preserved

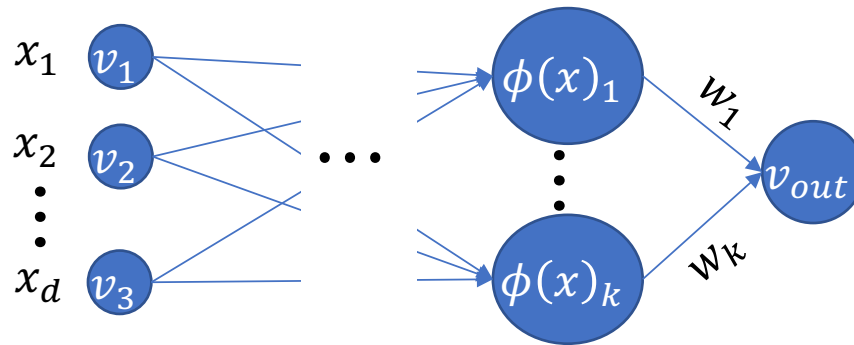
$$\min_z \sum_{i,j \in [N]} (\|z(x^{(i)}) - z(x^{(j)})\| - M_{ij})^2$$

- sometimes  $\min_z \sum_{i,j \in [N]} \frac{(\|z(x^{(i)}) - z(x^{(j)})\| - M_{ij})^2}{M_{ij}^2}$



# Autoencoders

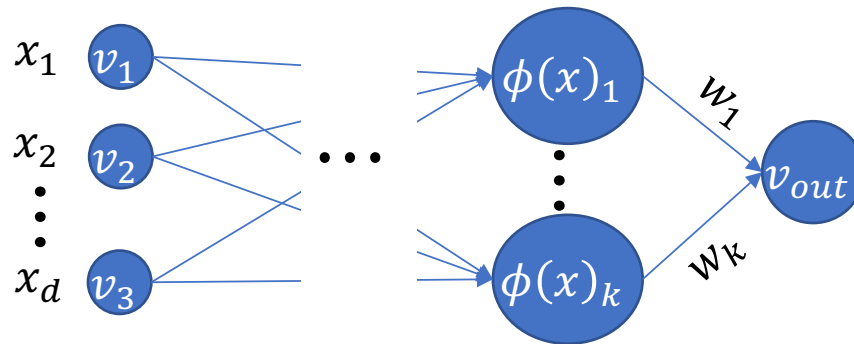
- Recall neural networks as feature learning



- was learned for some supervised learning task
- weights learned by minimizing  $\ell(v_{out}, y)$
- but we don't have  $y$  anymore!

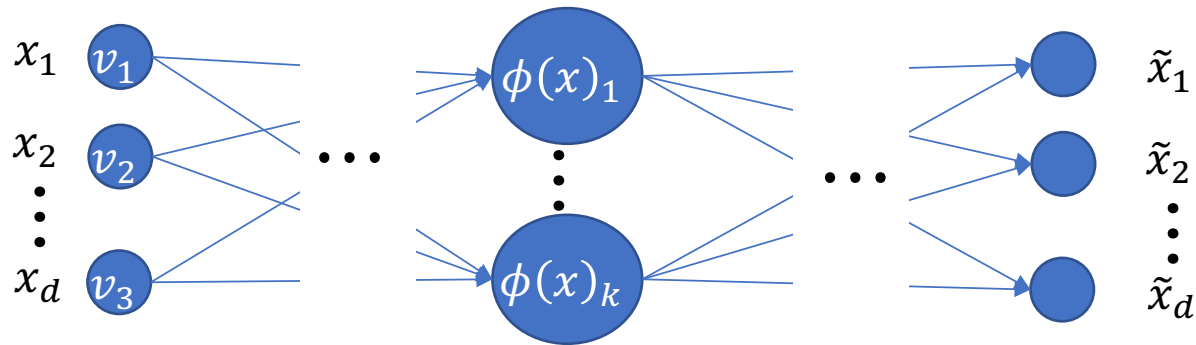
# Autoencoders

- Recall neural networks as feature learning



- was learned for some supervised learning task
- weights learned by minimizing  $\ell(v_{out}, y)$
- but we don't have  $y$  anymore!
- instead use another “decoder” network to reconstruct  $x$

# Autoencoders



- $\phi(x) = f_{W_1}(x)$
- $\tilde{x} = f_{W_2}(\phi(x))$
- some loss  $\ell(\tilde{x}, x)$

$$\hat{W}_1, \hat{W}_2 = \min_{W_1, W_2} \sum_{i=1}^N \ell \left( f_{W_2} \left( f_{W_1}(\mathbf{x}^{(i)}) \right), \mathbf{x}^{(i)} \right)$$

- learn using SGD with backpropagation