



Please scan this code to register your presence on Ilias  
for contact-tracing purposes. Your Ilias identity will be used.  
– Only do this if you are *physically present!* –

# DATA LITERACY

## LECTURE 08

### DIMENSIONALITY REDUCTION

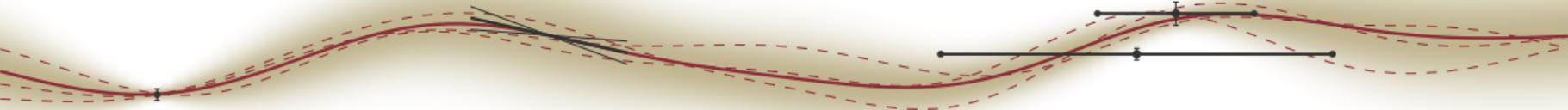
Philipp Hennig

13 December 2020

EBERHARD KARLS  
**UNIVERSITÄT**  
TÜBINGEN



FACULTY OF SCIENCE  
DEPARTMENT OF COMPUTER SCIENCE  
CHAIR FOR THE METHODS OF MACHINE LEARNING





# Challenges arising from Practical Data

## Lecture Outlook

**Data Collection** How do we get data? \_\_\_\_\_ Lectures 2 & 3

- ▶ avoiding *bias* in the collection of data
  - ▶ making sure data is as *informative* as possible
- Experimental Design  
Information

**Estimating from Data** What does the data (not) tell us? \_\_\_\_\_ Lectures 4–9

- ▶ Deriving *estimates* of related quantities from observations Statistical Concepts
- ▶ Can we *trust* the estimate? Variance of Estimators
- ▶ Do data/estimates tell us anything noteworthy? Tests of Significance
- ▶ Did we cheat on ourselves? Properties and Problems of Testing
- ▶ Making Predictions from Data Regression
- ▶ Finding Structure in Data Dimensionality Reduction and Clustering

**Fairness** and other societal aspects \_\_\_\_\_ Lectures 10 & 11

- ▶ Does our analysis put certain people at an disadvantage?
- ▶ Can/should we do anything about this?

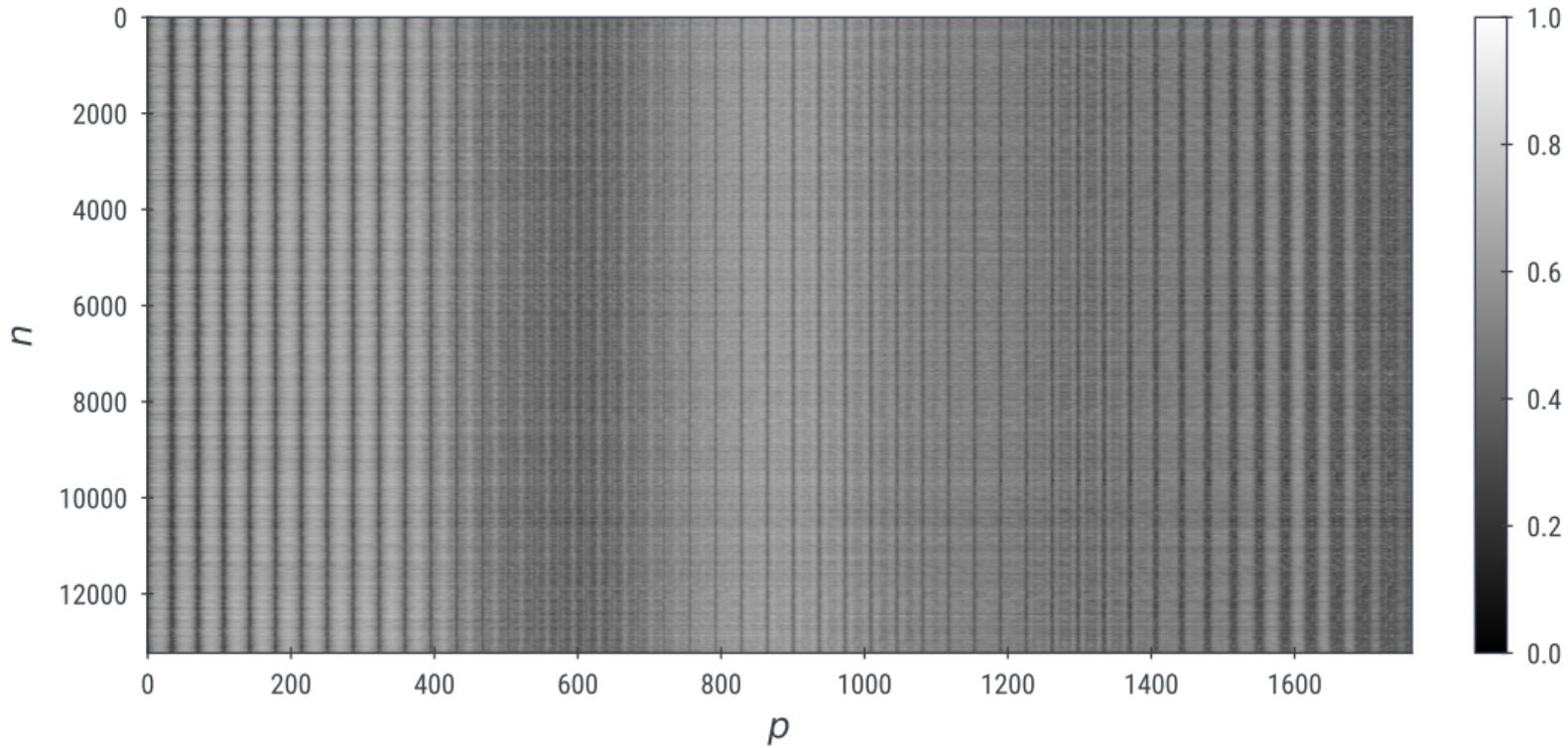
**Style and Techniques** \_\_\_\_\_ Lectures 12–14

- ▶ Collecting, documenting and storing data big and small
- ▶ Designing good code in the presence of data
- ▶ Visualizing and presenting data, analysis and results



# Data is Messy

and usually high-dimensional





Consider a dataset  $X \in \mathbb{R}^{N \times D}$ . **Generative Dimensionality Reduction** aims to find an **encoding**  $\phi : \mathbb{R}^D \rightarrow \mathbb{R}^M$  and a **decoding**  $\psi : \mathbb{R}^M \rightarrow \mathbb{R}^D$  with  $M \ll D$  such that the encoded representation

$$Z := \phi(X) \in \mathbb{R}^{N \times M}$$

is a *good approximation* of  $X$  in the sense that some **reconstruction loss** of  $\tilde{X} = \psi(Z)$ ,

$$\mathcal{L}(X, \psi(Z)) = \mathcal{L}(X, \psi \circ \phi(X))$$

is minimized or small. This may be done, e.g., to

- ▶ save memory
- ▶ construct a low-dimensional visualization
- ▶ “find structure”

Dimensionality reduction does *not* normally use known class labels  $Y \in \mathbb{R}^{N \times C}$ . This is what supervised methods (remember logistic regression!) is for.

# Principal Component Analysis (PCA)

Linear dimensionality reduction

Beltrami, 1873, Jordan, 1874, Pearson, 1901, Schmidt, 1907, Hotelling, 1933, Lanczos, 1950



Consider Data  $X \in \mathbb{R}^{N \times D} = [x_1; \dots; x_N]$ . We are looking for "the best" basis  $U = \{u_i\}_{i=1,\dots,D}$ ,  $u_i^T u_j = \delta_{ij}$  to construct an  $M$ -dimensional approximation to

$$x_n = \sum_{i=1}^D (x_n^T u_i) u_i \quad \text{i.e. } X = (XU)U^T \quad \text{as} \quad \tilde{x}_n := \sum_{i=1}^M a_{ni} u_i + \sum_{i=M+1}^D b_i u_i$$

To define "best", decide to minimize the *square empirical risk*

$$J = \frac{1}{N} \sum_{n=1}^N \|x_n - \tilde{x}_n\|^2 = \frac{1}{N} \sum_{n=1}^N \sum_{j=1}^D \left[ x_n - \sum_{i=1}^M a_{ni} u_i - \sum_{i=M+1}^D b_i u_i \right]_j^2$$

This yields (homework)  $a_{ni} = x_n^T u_i$ , and  $b_i = \bar{x}^T u_i$ . (So  $Z = \phi(X) = XU_{:M} + \frac{1}{N}\mathbf{1}\mathbf{1}^T XU_{M+1:}$ ) and

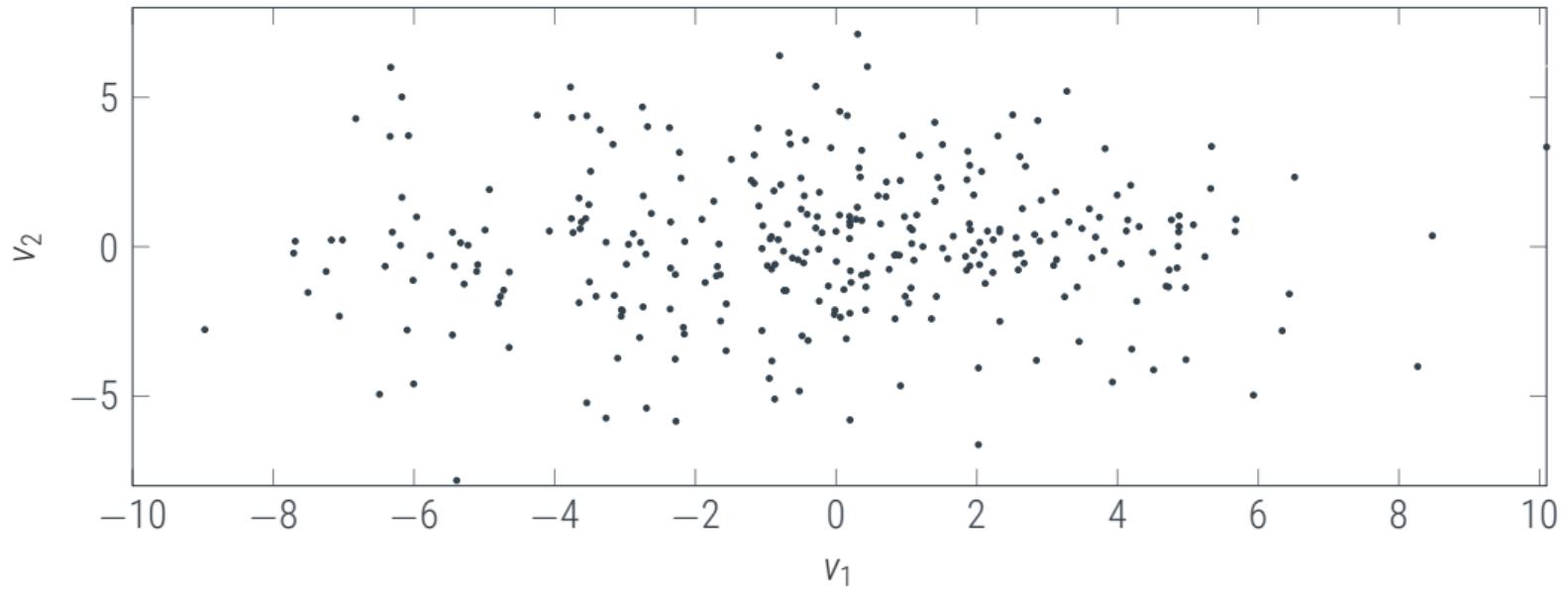
$$J = \frac{1}{N} \sum_{i=M+1}^D \sum_{n=1}^N u_i^T (x_n - \bar{x})(x_n - \bar{x})^T u_i = \sum_{i=M+1}^D u_i^T S u_i \text{ with } S := \frac{1}{N} \sum_{n=1}^N (x_n - \bar{x})(x_n - \bar{x})^T$$

$S$  is the *sample covariance matrix*. Hence, choose  $U$  as the eigenvectors of  $S$ . If we first center the data  $\hat{X} = X - X.\text{mean}()$ , so  $b = 0$ , the  $U$  are the right **singular vectors** of  $\hat{X} = Q\Sigma U^T$ .



# Low-Dimensional Representation

PCA on the mystery dataset



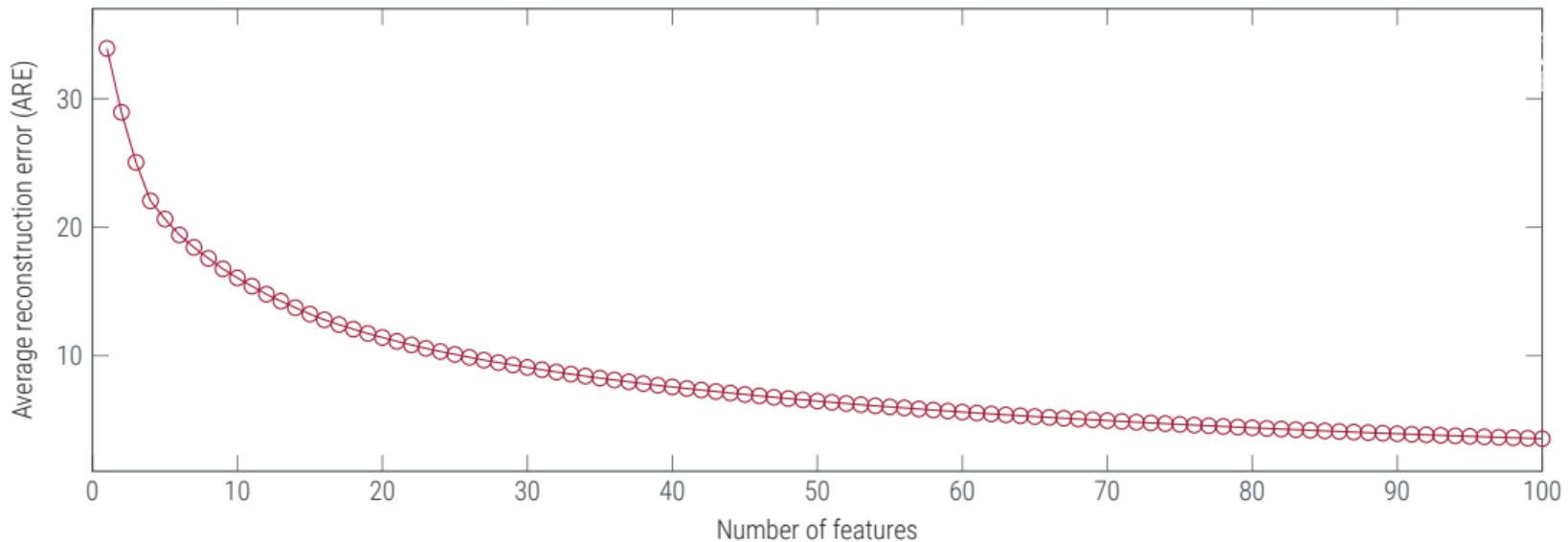
PCA does *not* necessarily reveal “structure” in the data.



# But it *does* compress the data

i.e. what it is designed to do

$$\text{ARE} = \frac{1}{N} \sum_{n=1}^N \|x_n - \tilde{x}_n\|^2 \quad (\text{data has } D = 1764)$$





# Some Datasets are easy on our eyes

Labeled Faces in the Wild (LFW)

Huang et al., 2009. Pre-processed by Samira Samadi (more in next lecture)



$N = 5924$  faces at  $49 \times 36 = 1764$  pixels, gray scale (pre-processed from full dataset)



# Compressing Faces

PCA on *Labeled Faces in the Wild*



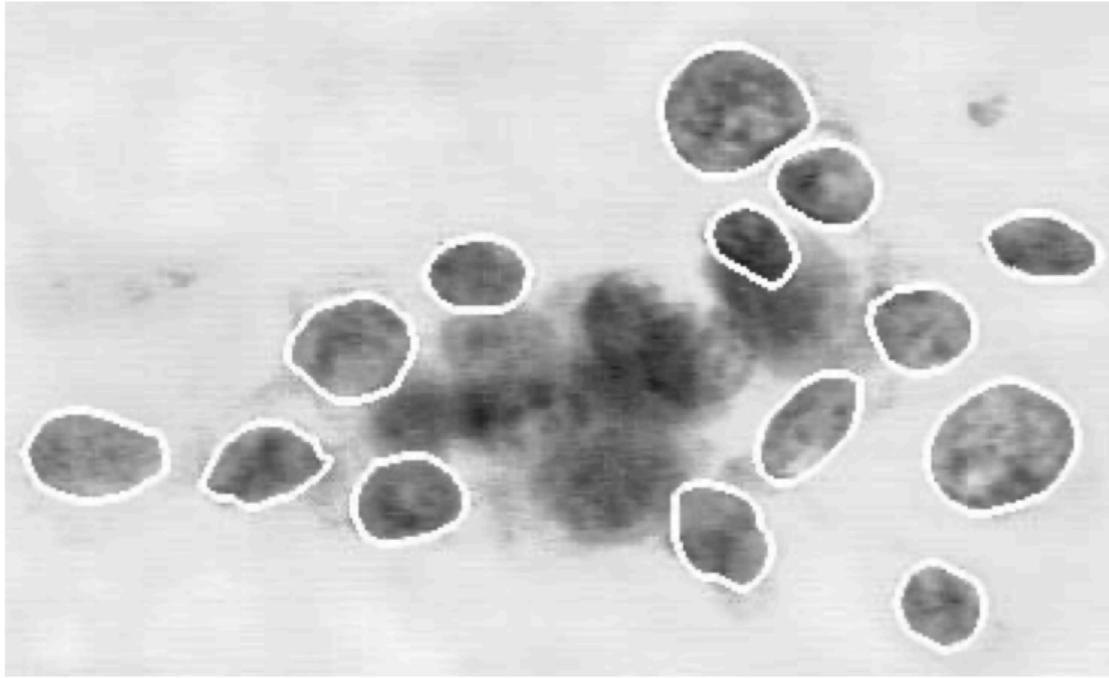
Principal Component Analysis (PCA) constructs the low-dimensional *linear projection* of the data that minimizes the  $\ell_2$  reconstruction error.



# Another Dataset

The Wisconsin Breastcancer Dataset

[https://archive.ics.uci.edu/ml/datasets/Breast+Cancer+Wisconsin+\(Diagnostic\)](https://archive.ics.uci.edu/ml/datasets/Breast+Cancer+Wisconsin+(Diagnostic))



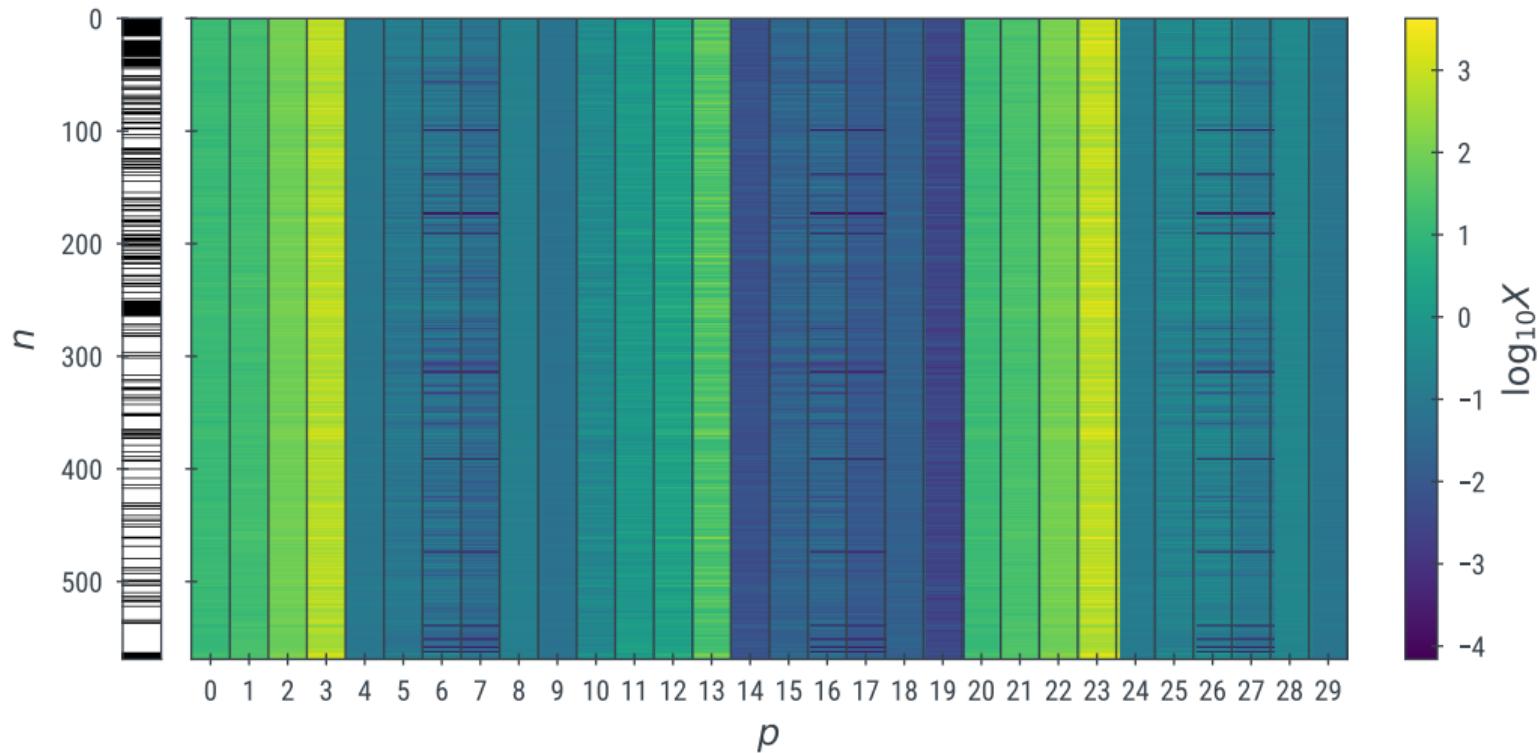
Street, Wolberg, Mangasarian. *Nuclear Feature Extraction For Breast Tumor Diagnosis.*

Proc. SPIE Int. Symp. on Electronic Imaging, vol. 1905, pp. 861–870. 1992



# The Data

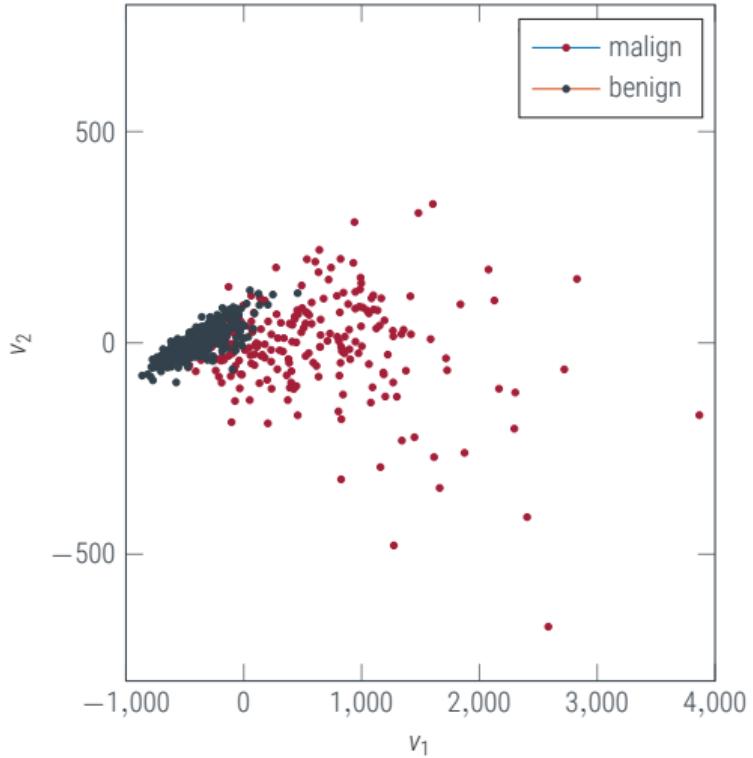
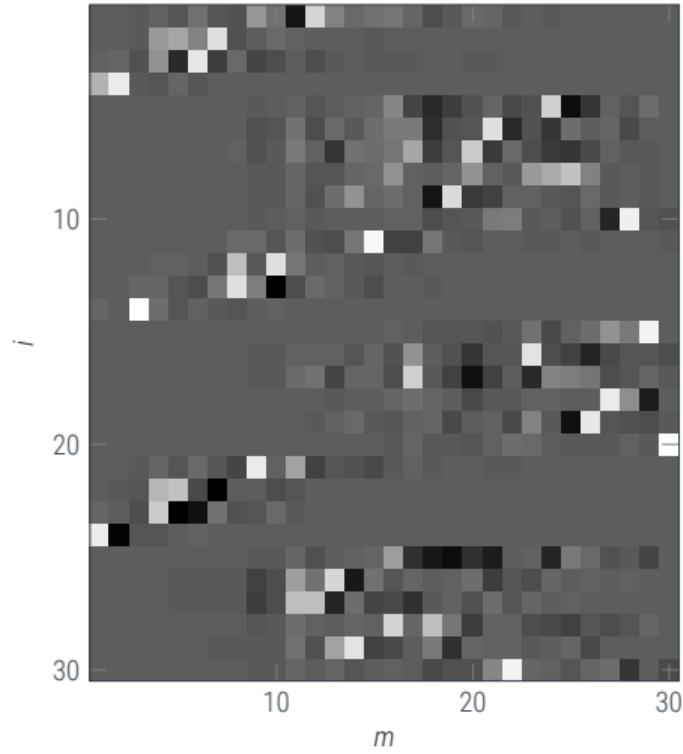
not much to see here





# Visualized

First two principal components





# The Meta Data

for the Wisconsin Breast Cancer Dataset

## Attribute Information:

1) ID number 2) Diagnosis (M = malignant, B = benign)

3-32) Ten real-valued features are computed for each cell nucleus:

1) radius (mean of distances from center to points on the perimeter)

2) texture (standard deviation of gray-scale values)

3) perimeter

4) area

5) smoothness (local variation in radius lengths)

6) compactness ( $\text{perimeter}^2 / \text{area} - 1.0$ )

7) concavity (severity of concave portions of the contour)

8) concave points (number of concave portions of the contour)

9) symmetry

10) fractal dimension ("coastline approximation" - 1)

The mean, standard error and "worst" or largest (mean of the three largest values) of these features were computed for each image, resulting in 30 features.

For instance, field 3 is Mean Radius, field 13 is Radius SE, field 23 is Worst Radius.



# All Data has Units!

even if they are not recorded in the metadata!

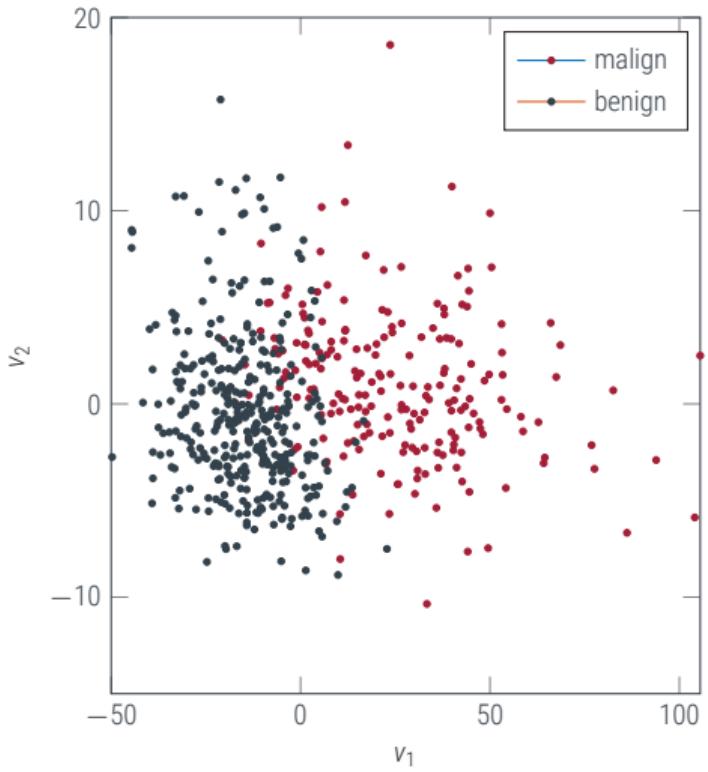
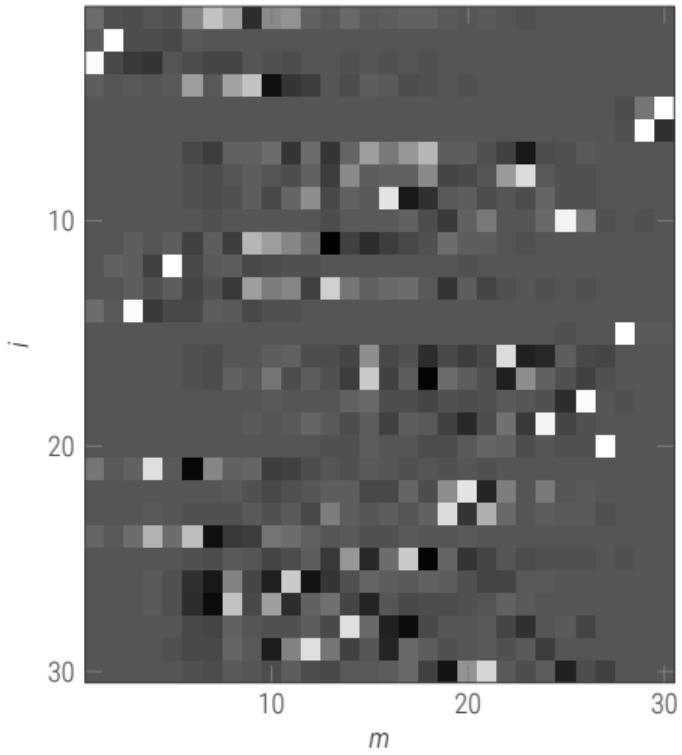
1. radius (mean of distances from center to points on the perimeter) [cm]
2. texture (standard deviation of gray-scale values)  $\in [0, 255]$
3. perimeter [cm?]
4. area [cm<sup>2</sup>]
5. smoothness (local variation in radius lengths) [cm]
6. compactness (perimeter<sup>2</sup> / area - 1.0)  $\in [1, \infty]$
7. concavity (severity of concave portions of the contour) [???
8. concave points (number of concave portions of the contour) [???
9. symmetry [???
10. fractal dimension ("coastline approximation" - 1) [???

Note that  $X \leftarrow X \cdot \text{diag}(\mathbf{s})$  changes the singular vectors of  $X$ !



# PCA, in other units

rescaled  $d_2/255, d_5/10, d_8/200$





## Principal Component Analysis

- ▶ amounts to *projecting* the (centered) data  $X$  onto its first  $M$  right **singular vectors**

$$\hat{X} = X - \mathbf{1}\bar{x}^T, \quad \hat{X} = Q\Sigma U^T = \text{svd}(\hat{X}), \quad \tilde{X} = (\hat{X}U_{1:M})U_{1:M}^T$$

- ▶ PCA is **meaningless** if the individual data dimensions do not all have the same units.

Before reducing the dimensionality of your data, think about *why* you want to do so?

- ▶ Are you hoping to see structure?
- ▶ Do you want to save memory?
- ▶ Are there classes in the data, and you want to know how they differ?

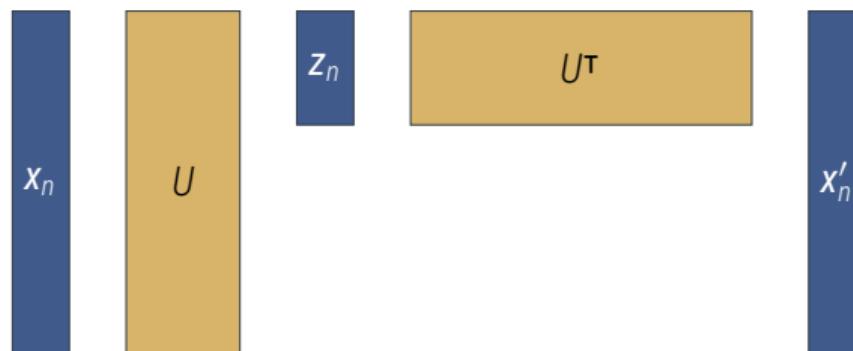
# Nonlinear Dimensionality Reduction

beyond PCA

- ▶ PCA minimizes the **square ( $\ell_2$ ) reconstruction error** of a **linear projection**  $z_n = x_n^\top U$

$$J(U) = \frac{1}{N} \sum_{n=1}^N \|x_n - \tilde{x}_n\|^2 = \frac{1}{N} \sum_{n=1}^N \|x_n - x_n U U^\top\|^2 = \frac{1}{N} \sum_{n=1}^N \ell_2(x_n, \psi(\phi(x_n)))$$

with *encoder*  $\phi(x) = U^\top (x - \bar{x})^\top$  and *decoder*  $\psi(z) = Uz + \bar{x} = (x')^\top$



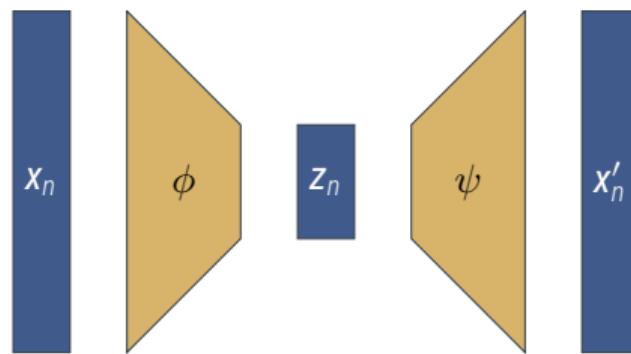
# Nonlinear Dimensionality Reduction

beyond PCA

- ▶ PCA minimizes the **square ( $\ell_2$ )** reconstruction error of a **linear** projection  $z_n = x_n^\top U$

$$J(U) = \frac{1}{N} \sum_{n=1}^N \|x_n - \tilde{x}_n\|^2 = \frac{1}{N} \sum_{n=1}^N \|x_n - x_n U U^\top\|^2 = \frac{1}{N} \sum_{n=1}^N \ell_2(x_n, \psi(\phi(x_n)))$$

with *encoder*  $\phi(x) = U^\top(x - \bar{x})^\top$  and *decoder*  $\psi(z) = Uz + \bar{x} = (x')^\top$



- ▶ Exchange the **linear** projections by **nonlinear** ones to get an **autoencoder**.

# Nonlinear Dimensionality Reduction

beyond PCA

- ▶ PCA minimizes the **square ( $\ell_2$ ) reconstruction error** of a **linear** projection  $z_n = x_n^\top U$

$$J(U) = \frac{1}{N} \sum_{n=1}^N \|x_n - \tilde{x}_n\|^2 = \frac{1}{N} \sum_{n=1}^N \|x_n - x_n U U^\top\|^2 = \frac{1}{N} \sum_{n=1}^N \ell_2(x_n, \psi(\phi(x_n)))$$

with *encoder*  $\phi(x) = U^\top x^\top$  and *decoder*  $\psi(z) = Uz = (x')^\top$

- ▶ Exchange the **linear** projections by **nonlinear** ones to get an **autoencoder**.
- ▶ replace the  $\ell_2$  with the log loss (KL-divergence) to get a *variational autoencoder* (VAE)

$$J(\phi, \psi) = D_{\text{KL}}(q_\phi(z | x) \| p_\psi(z | x)) = \frac{1}{N} \sum_{n=1}^N \log q_\phi(z | x) - \log p_\psi(z | x)$$

- ▶ this requires parametrization choices for  $q_\phi, p_\psi$ . Typically Gaussian.

This can help with the problems above, but only if you know what you're doing!



## Generative Dimensionality Reduction

- ▶ amounts to *projecting* the data  $X$  onto a low-dimensional *latent representation*  $Z = \phi(X)$  such that the *reconstruction*

$$X' = \psi(Z) = (\psi \circ \phi)X$$

"captures structure" of  $X$ , as described by some loss function  $J$

- ▶ The simplest case is PCA, where  $\phi, \psi$  are linear and  $J$  is quadratic.
- ▶ Note: "Capturing structure" does not automatically do what you may want
  - ▶ the latent representation may not be a good visualization
  - ▶ if there is structure the loss does not capture, it will not generally be preserved by the encoding

But there are also *non-generative* ways to reduce dimensionality. For example, we may search for a low-dimensional encoding that preserves "structure" of the high-dimensional data, but without requiring an invertible map.

# Embedding without Encoding

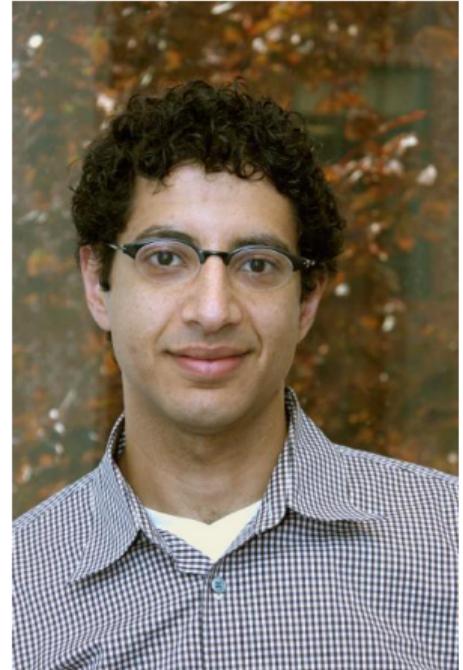
neighborhood embedding



- ▶ Instead of explicitly constructing the encoding/decoding  $\phi, \psi$ , define affinity graph / neighborhood probabilities  $p_{ij}$  and  $q_{ij}$  on input  $X = [x_1, \dots, x_n] \in \mathbb{R}^{n \times d}$  and visualization space  $Y = [y_1, \dots, y_n] \in \mathbb{R}^{n \times 2}$
- ▶ tune  $y$  (thus  $q_{ij}$ ) to minimize some divergence, e.g.  $D_{\text{KL}}(p||q)$

Some immediate observations

- ▶ this means the embedded representation  $y$  can not be interpreted as an encoding. There are no explicit maps  $\phi : x \mapsto y$  or  $\psi : y \mapsto x$ .
- ▶ building blocks: affinity  $p$ , affinity  $q$ , divergence, optimizer. All of which come with assumptions and caveats.



Sam Roweis  
1972–2010

# *t*-SNE

Roweis & Hinton, NeurIPS 2002. v.d. Maaten & Hinton, JMLR 9 (11/2008), pp. 2579–2605

- ▶ affinity  $p$  as ‘probability for  $i$  to pick  $j$  as neighbor if picked by isotropic Gaussian, symmetrised’

$$p_{j|i} = \frac{\exp(-d_{ij}^2)}{\sum_{k \neq i} \exp(-d_{ik}^2)} \quad \text{with} \quad d_{ij}^2 = \frac{\|x_i - x_j\|^2}{2\sigma_i^2} \quad \text{and set} \quad p_{jj} = \frac{1}{2n}(p_{j|i} + p_{i|j}), p_{ii} = 0$$

- ▶ set hyperparameters  $\sigma_i$  to achieve ad-hoc *perplexity*  $2^{\mathbb{H}(p_i)} = 2^{-\sum_j p_{ij} \log_2 p_{ij}}$
- ▶ affinity  $q$  as Cauchy distribution (‘Student- $t$  distribution with one degree of freedom’)

$$q_{ij} = \frac{(1 + \|y_i - y_j\|^2)^{-1}}{\sum_k \sum_{\ell \neq k} (1 + \|y_k - y_\ell\|^2)^{-1}} \quad \text{and} \quad q_{ii} = 0$$

- ▶ divergence: Kullback-Leibler  $D_{KL}(p||q) = \sum_{ij, i \neq j} p_{ij} \log \frac{p_{ij}}{q_{ij}}$
- ▶ optimization: random initialization, momentum stochastic gradient descent

$$\frac{\partial D_{KL}(p||q)}{\partial y_i} = 4 \sum_{j \neq i} \frac{(p_{ij} - q_{ij})(y_i - y_j)}{1 + \|y_i - y_j\|^2}$$



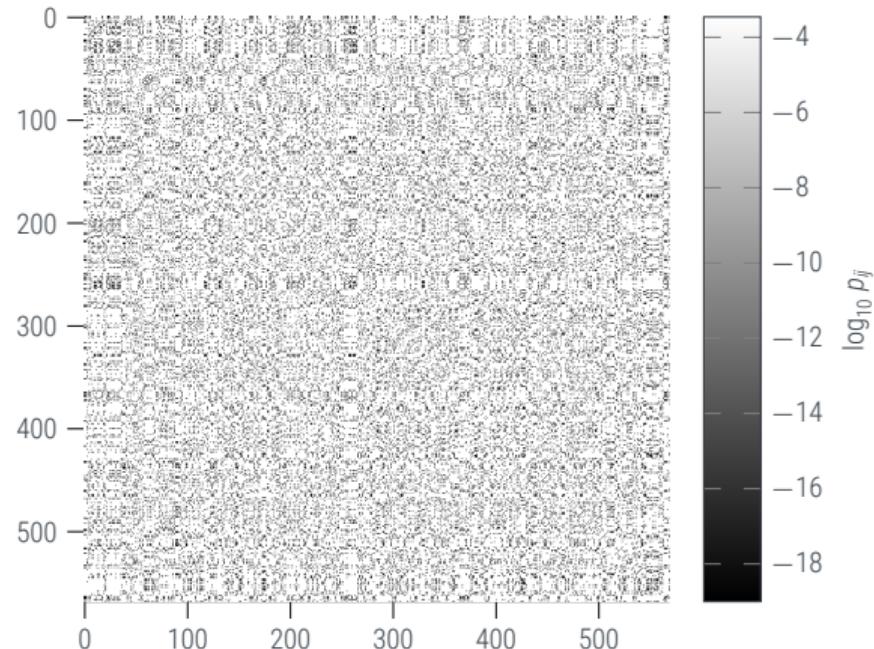
# Affinity $p$

WBC dataset

$$p_{j|i} = \frac{\exp(-d_{ij}^2)}{\sum_{k \neq i} \exp(-d_{ik})}$$

where  $d_{ij}^2 = \frac{\|x_i - x_j\|^2}{2\sigma_i^2}$

and set  $p_{ij} = \frac{1}{2n}(p_{j|i} + p_{i|j}),$   
 $p_{ii} = 0$





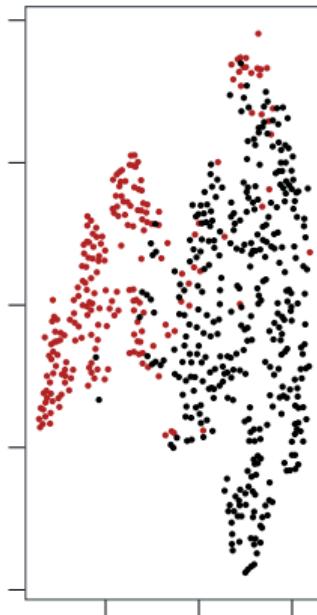
# What do your dimensions mean?

re-scaled inputs

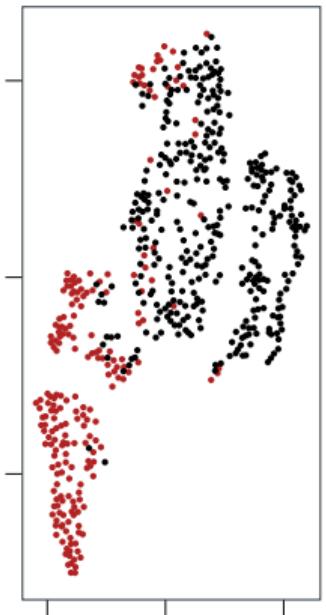
- ▶ Remember that our (WBC) dataset has 30 dimensions of differing meaning
- ▶ in particular,  $X[:, : 10]$  contains means,  $X[:, 10 : 20]$  contains variances,  $X[:, 20 : 30]$  maxima

```
X2 = np.hstack([X[:, :10], np.sqrt(X[:, 10:20]), X[:, 20:30]])
```

perplexity 25



perplexity 25, transformed





## Summary $t$ -stochastic neighborhood embedding ( $t$ -SNE)

- ▶ an unsupervised way to make low-dimensional representations of high-dimensional data that retain as much *local* similarity as possible
- ▶ lots of free parameters, assumptions and algorithmic tricks needed. Not a *black box* like PCA
- ▶ SNE is not a formal way to find scientific “insights”, but a way to explore a dataset and find leads to interesting properties

Powerful tools can be dangerous or useless if you do not understand them.

So far, we have ignored that our dataset is labeled.

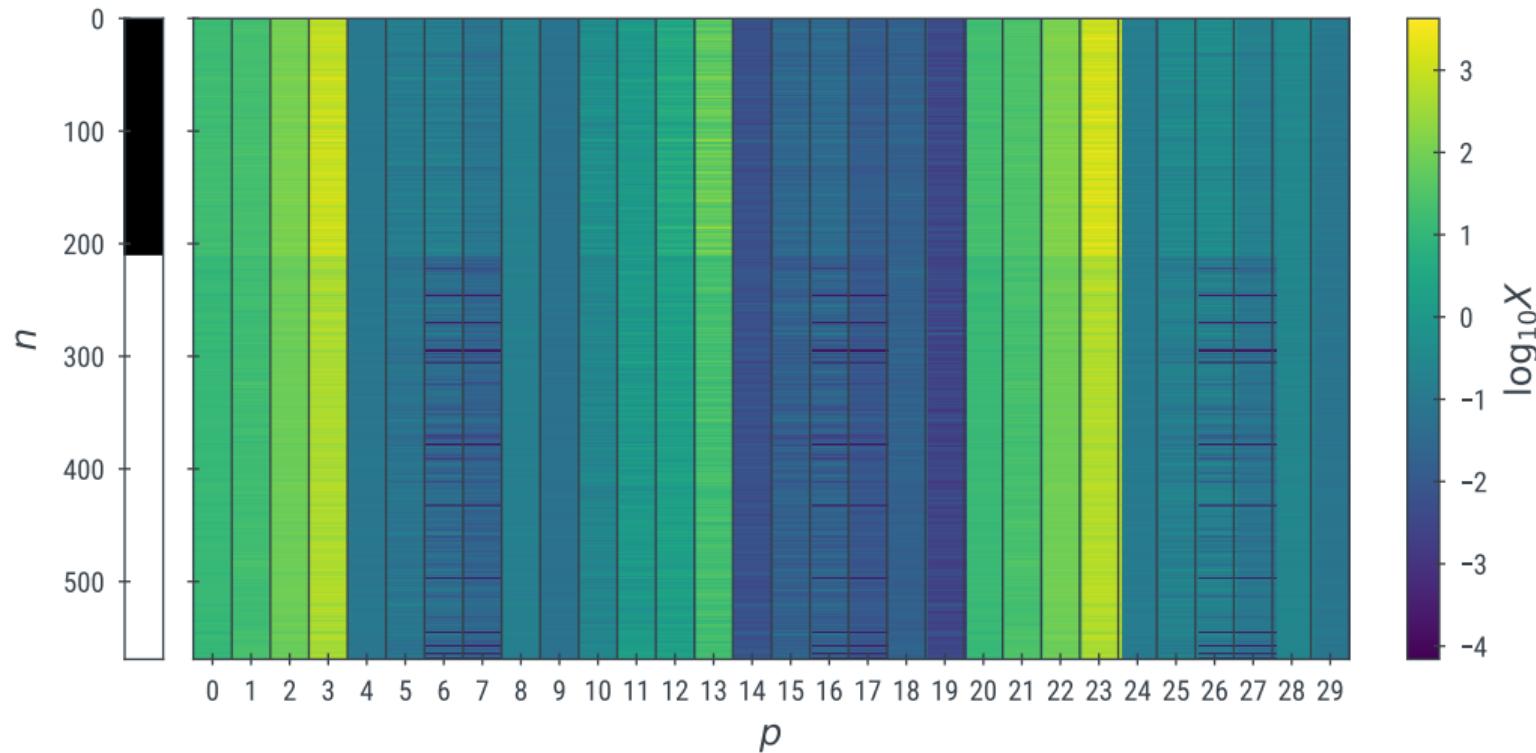
Can we use the labels to learn a *supervised* embedding?

# Can we separate the two classes?

What can we do to *separate* the classes from each other? To *discriminate* among them?



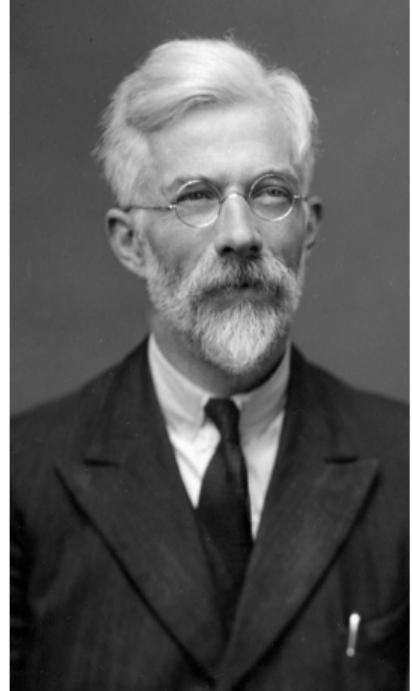
WBC dataset, sorted by class label



# Can we separate the two classes?

Latent Discriminant Analysis

R.A. Fisher, *The Use of Multiple Measurements in Taxonomic Problems*. Annals of Eugenics, 1936



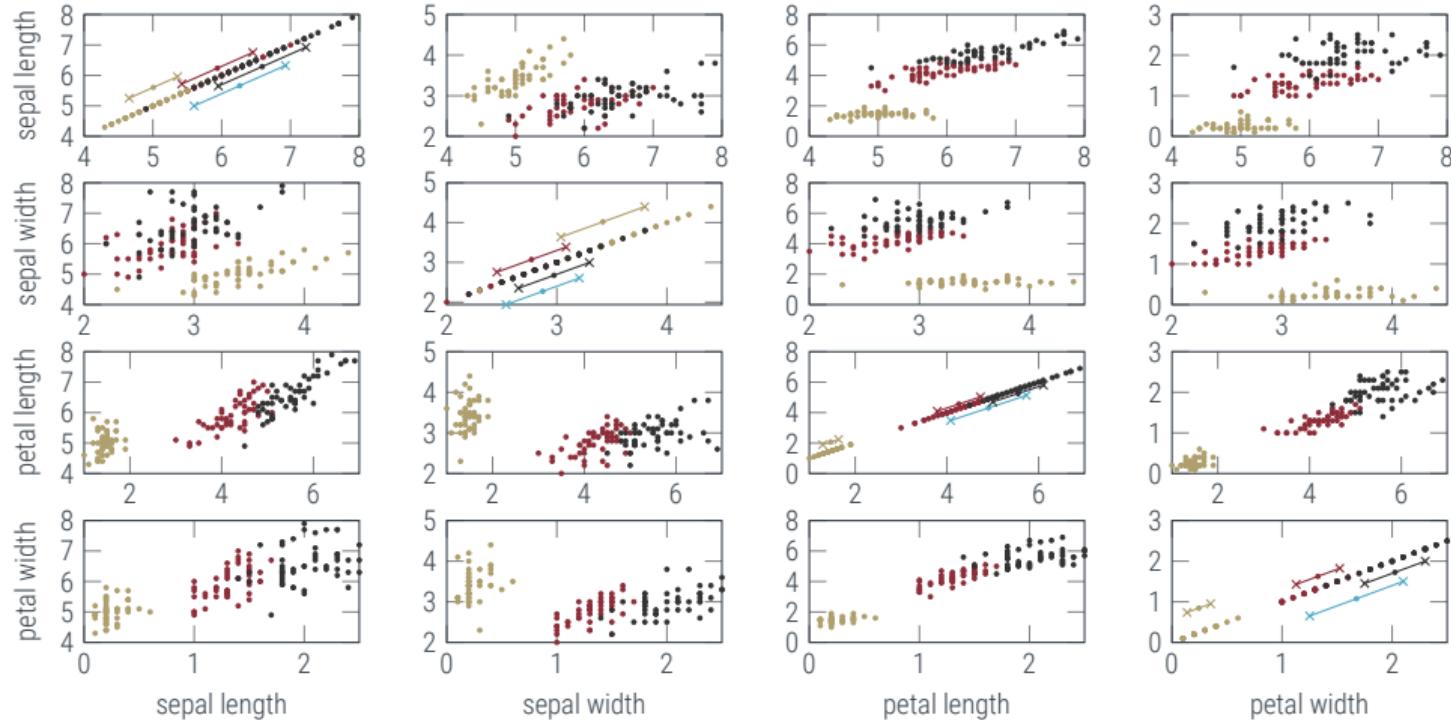
<i>Iris setosa</i>				<i>Iris versicolor</i>				<i>Iris virginica</i>			
Sepal length	Sepal width	Petal length	Petal width	Sepal length	Sepal width	Petal length	Petal width	Sepal length	Sepal width	Petal length	Petal width
5.1	3.5	1.4	0.2	7.0	3.2	4.7	1.4	6.3	3.3	6.0	2.5
4.9	3.0	1.4	0.2	6.4	3.2	4.5	1.5	5.8	2.7	5.1	1.9
4.7	3.2	1.3	0.2	6.9	3.1	4.9	1.5	7.1	3.0	5.9	2.1
4.6	3.1	1.5	0.2	5.5	2.3	4.0	1.3	6.3	2.9	5.6	1.8
5.0	3.6	1.4	0.2	6.5	2.8	4.6	1.5	6.5	3.0	5.8	2.2
5.4	3.9	1.7	0.4	5.7	2.8	4.5	1.3	7.6	3.0	6.6	2.1
4.6	3.4	1.4	0.3	6.3	3.3	4.7	1.6	4.9	2.5	4.5	1.7
5.0	3.4	1.5	0.2	6.0	3.0	4.6	1.3	7.3	2.9	6.2	1.2

Ronald A. Fisher  
1890–1962



# Which projection discriminates best?

The Iris dataset



- *versicolor*
- *setosa*
- *virginica*
- *virginica* ∪ *versicolor*



# Linear Discriminant Analysis

finding the best separator of Gaussian distributions

Consider Data  $\mathbf{x}_n \in \mathbb{R}^D$ ,  $n = 1, \dots, N$  from two classes  $\mathcal{C}_1, \mathcal{C}_2$  (i.e. index sets with  $\mathcal{C}_1 \cup \mathcal{C}_2 = [1, \dots, N]$  and  $\mathcal{C}_1 \cap \mathcal{C}_2 = \emptyset$ ) with class means and un-normalised covariances

$$\mathbf{m}_i := \frac{1}{N_i} \sum_{n \in \mathcal{C}_i} \mathbf{x}_n \quad S_i := \sum_{n \in \mathcal{C}_i} (\mathbf{x}_n - \mathbf{m}_i)(\mathbf{x}_n - \mathbf{m}_i)^\top \quad \text{where } N_i = |\mathcal{C}_i|.$$

Try to find a projection  $\mathbf{w}$  such that the scalar variables  $y_i = \mathbf{w}^\top \mathbf{x}_i$  maximize the ratio between inter-class variance  $(m_2 - m_1)^2 = (\mathbf{w}^\top (\mathbf{m}_2 - \mathbf{m}_1))^2$  and in-class square distances

$$s_i^2 = \sum_{n \in \mathcal{C}_i} (y_i - m_i)^2 = \sum_{n \in \mathcal{C}_i} \mathbf{w}^\top (\mathbf{x}_n - \mathbf{m}_i)(\mathbf{x}_n - \mathbf{m}_i)^\top \mathbf{w} = \mathbf{w}^\top S_i \mathbf{w}.$$

That is, maximise

$$J(\mathbf{w}) = \frac{(m_2 - m_1)^2}{s_1^2 + s_2^2}$$



# Linear Discriminant Analysis

finding the best separator of Gaussian distributions

Maximise (note we can scale  $w$  arbitrarily, so only its direction matters)

$$J(w) = \frac{(m_2 - m_1)^2}{s_1^2 + s_2^2} = \frac{w^\top S_B w}{w^\top S_W w} = \frac{w^\top \overbrace{(m_2 - m_1)(m_2 - m_1)^\top}^{=:S_B} w}{w^\top \underbrace{\left( \sum_{n \in C_1} (x_n - m_1)(x_n - m_1)^\top + \sum_{n \in C_2} (x_n - m_2)(x_n - m_2)^\top \right)}_{S_W} w}$$

$$\frac{\partial J}{\partial w} = 2 \frac{S_B w}{w^\top S_W w} - 2 \frac{w^\top S_B w}{(w^\top S_W w)^2} S_W w = C_1 \cdot (m_2 - m_1) - C_2 \cdot S_W w$$

$$w \propto S_W^{-1} (m_2 - m_1)$$



# Linear Discriminant Analysis

finding the best separator of Gaussian distributions

When should we predict class 1 over 2? The  $(m_i, S_i)$  define *generative* Gaussian distributions  $p(y_n \mid c = i) = \mathcal{N}(y_n; m_i, s_i)$ , and the marginal frequencies give a “prior” class probability  $p(c_i) = N_i/N$ .

**Exercise:** Find the value of the *best discriminator*  $y_0$  set by considering *posterior log odds*

$$r(y) = \log \frac{p(c_1 \mid y)}{p(c_2 \mid y)} = \log \frac{p(y \mid c_1)p(c_1)}{p(y \mid c_2)p(c_2)} = \log \frac{\mathcal{N}(y; m_1, s_1^2)N_1}{\mathcal{N}(y; m_2, s_2^2)N_2}$$

and solving for  $r(y_0) = 0$ .

# LDA for the Iris dataset

Directly from Fisher's paper

R.A. Fisher, *The Use of Multiple Measurements in Taxonomic Problems*. Annals of Eugenics, 1936

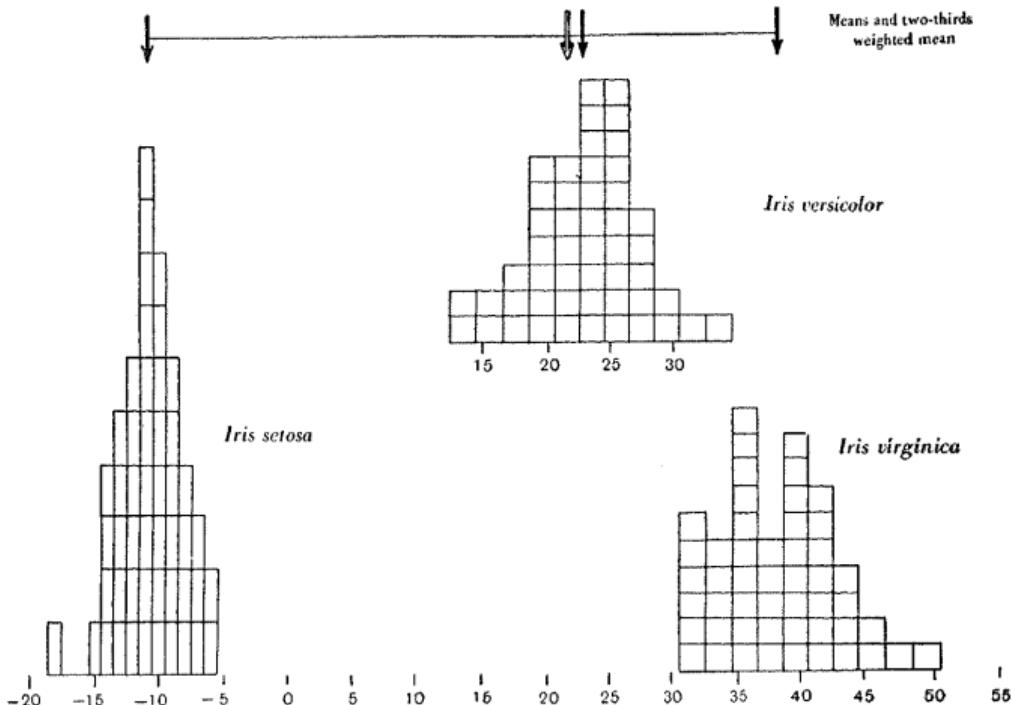


Fig. 1. Frequency histograms of the discriminating linear function, for three species of *Iris*.



## Summary: Dimensionality Reduction

- ▶ There are different motivations for dimensionality reduction
  - ▶ encoding/decoding paradigm (save memory)
  - ▶ neighborhood embedding (preserve similarity)
  - ▶ discrimination (separate classes from each other)
- ▶ which one to choose depends on your application. But you should know *why* you are using one technique over the other, because each dimensionality reduction creates some artifacts

Please provide feedback:





# A date for your calendar?

Methods of Machine Learning Research Seminar

Wednesday, 24 November, 15:00 (st)  
Sarah Bechtle (MPI IS / NYU)

## Lifelong learning in the real world

I'm going to present work centred around the fundamental research question of how a robot can learn in the real world. The robot gathers the data it needs to learn through interactions with - and perception of - the environment. In pursuing answering this research question, I work on two interconnected directions: 1) Model Based Learning in the Action Perception Loop, where a model of the environment and a policy are learned iteratively from data collected on the robot and 2) Learning to Learn in the Action Perception Loop, where the robot learns how to learn a task. These directions fundamentally are concerned with learning representations that generalize quickly to new tasks and scenarios, which is a key capability that allows us humans to continuously learn over a lifetime.



details at <https://talks.tue.ai/talks/talk/id=50>