# Machine Learning from Data

Lecture 27: Spring 2021

# Today's Lecture

- Input Preprocessing ✓
- Dimensionality Reduction and Feature Selection
- Principal Components Analysis (PCA)
- Hints, Data Cleaning, Validation

# Learning Aides

Additional tools that can be applied to all techniques

Preprocess data to account for arbitrary choices during data collection (input normalization)

Remove irrelevant dimensions that can mislead learning (PCA)

Incorporate known properties of the target function (hints and invariances)

Remove detrimental data (deterministic and stochastic noise)

Better ways to validate (estimate $E_{out}$) for model selection

# Nearest Neighbor

Mr. Good and Mr. Bad were both given credit cards by the Bank of Learning (BoL).

|  | Mr. Good | Mr. Bad |
| --- | --- | --- |
| (Age in years, Income in $ \times 1,000$) | (47,35) | (22,40) |

Mr. Unknown who has "coordinates" (21yrs,$36K) applies for credit. Should the BoL give him credit, according to the nearest neighbor algorithm?
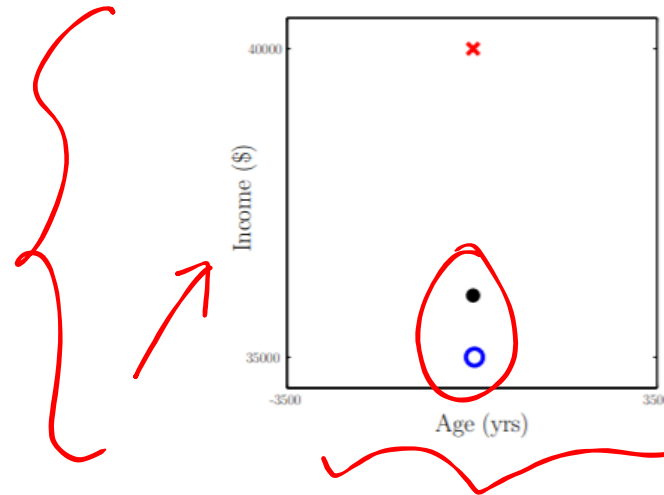
# Nearest Neighbor Uses Euclidean Distance

Mr. Good and Mr. Bad were both given credit cards by the Bank of Learning (BoL).

|  | Mr. Good | Mr. Bad |
|---|---|---|
| (Age in years, Income in $) | (47,35000) | (22,40000) |

Mr. Unknown who has "coordinates" (21yrs,$36000) applies for credit. Should the BoL give him credit, according to the nearest neighbor algorithm?

What if, income is measured in dollars instead of "K" (thousands of dollars)?



Learning →

Approve Similarity

Built-in Scale

# Uniform Treatment of Dimensions

Most learning algorithms treat each dimension equally

Nearear neighbor: $d(\mathbf{x}, \mathbf{x}') = \|\mathbf{x} - \mathbf{x}'\|$

Weight Decay: $\Omega(\mathbf{w}) = \lambda \mathbf{w}^{\mathsf{T}} \mathbf{w}$

SVM: margin defined using Euclidean distance

RBF: bump function decays with Euclidean distance

**Input Preprocessing**

Unless you want to emphasize certain dimensions, the data should be *preprocessed* to present each dimension on an equal footing
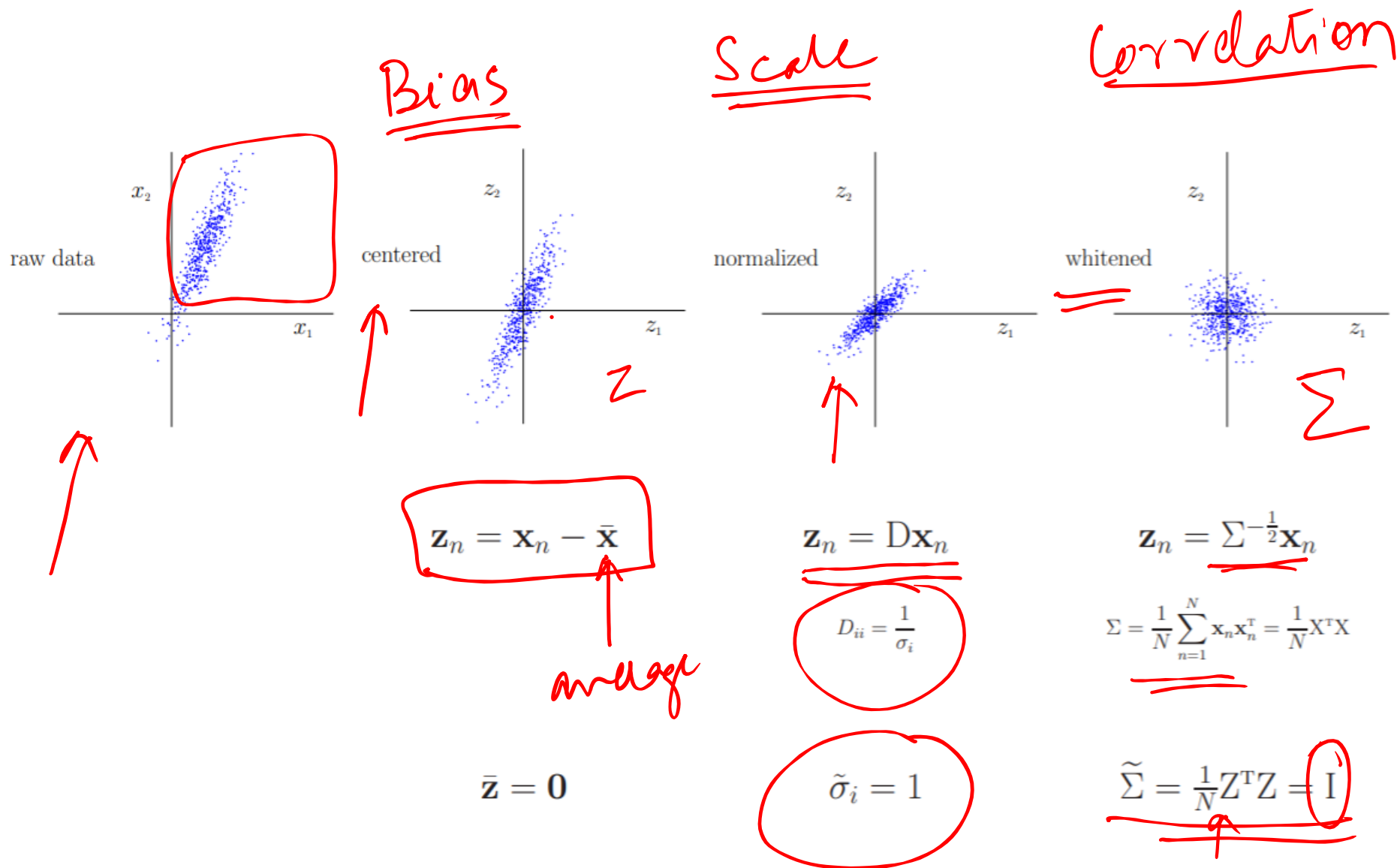
# Input Preprocessing is a Data Transform

$$X = \begin{bmatrix} - & \mathbf{x}_1^{\mathrm{T}} & - \\ - & \mathbf{x}_2^{\mathrm{T}} & - \\ & \vdots & \\ - & \mathbf{x}_N^{\mathrm{T}} & - \end{bmatrix} \qquad\qquad \mathbf{x}_n \mapsto \mathbf{z}_n$$

$$g(\mathbf{x}) = \tilde{g}(\Phi(\mathbf{x})).$$

Raw $\{\mathbf{x}_n\}$ have (for example) arbitrary scalings in each dimension, and $\{\mathbf{z}_n\}$ will not.

**Bias**

**Scale**

**Correlation**

raw data

centered

$z$

normalized

whitened

$\Sigma$

$$\mathbf{z}_n = \mathbf{x}_n - \bar{\mathbf{x}}$$

average

$$\bar{\mathbf{z}} = \mathbf{0}$$

$$\mathbf{z}_n = \mathrm{D}\mathbf{x}_n$$

$$D_{ii} = \frac{1}{\sigma_i}$$

$$\tilde{\sigma}_i = 1$$

$$\mathbf{z}_n = \Sigma^{-\frac{1}{2}}\mathbf{x}_n$$

$$\Sigma = \frac{1}{N}\sum_{n=1}^{N}\mathbf{x}_n\mathbf{x}_n^{\mathsf{T}} = \frac{1}{N}\mathrm{X}^{\mathsf{T}}\mathrm{X}$$

$$\widetilde{\Sigma} = \frac{1}{N}\mathrm{Z}^{\mathsf{T}}\mathrm{Z} = \mathrm{I}$$

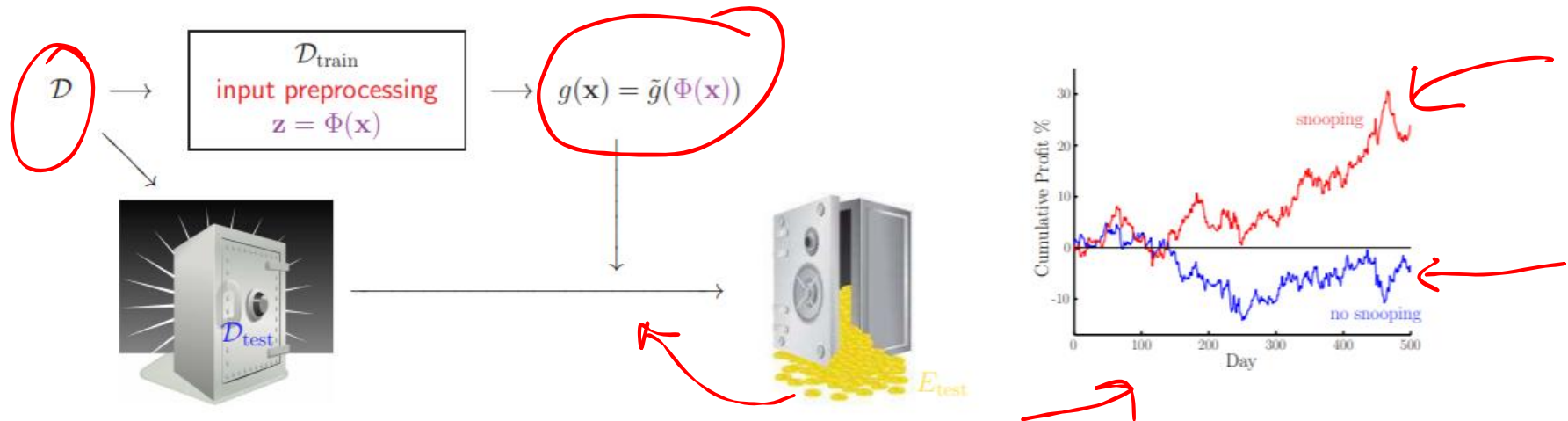# Only Use Training Data For Preprocessing



**WARNING!**

Transforming data into a more convenient format has a hidden trap which leads to data snooping.
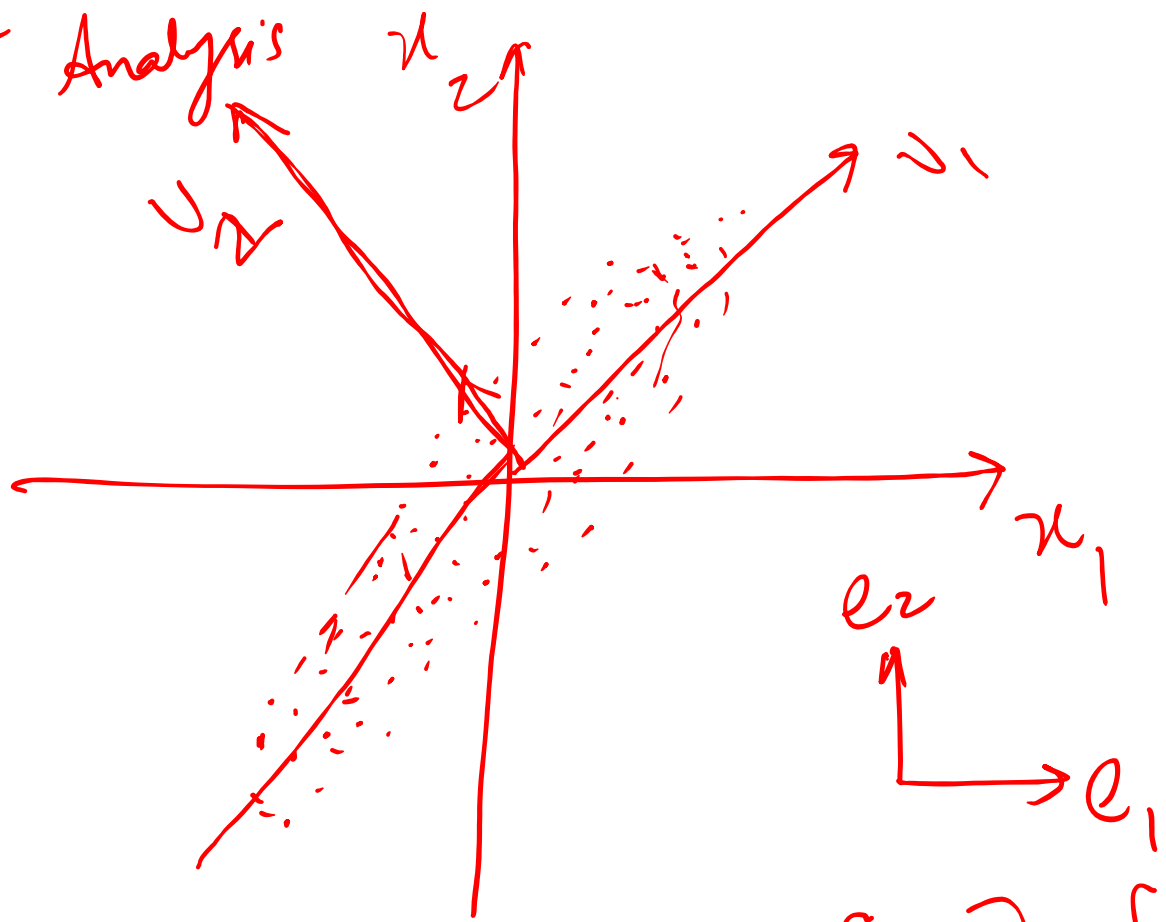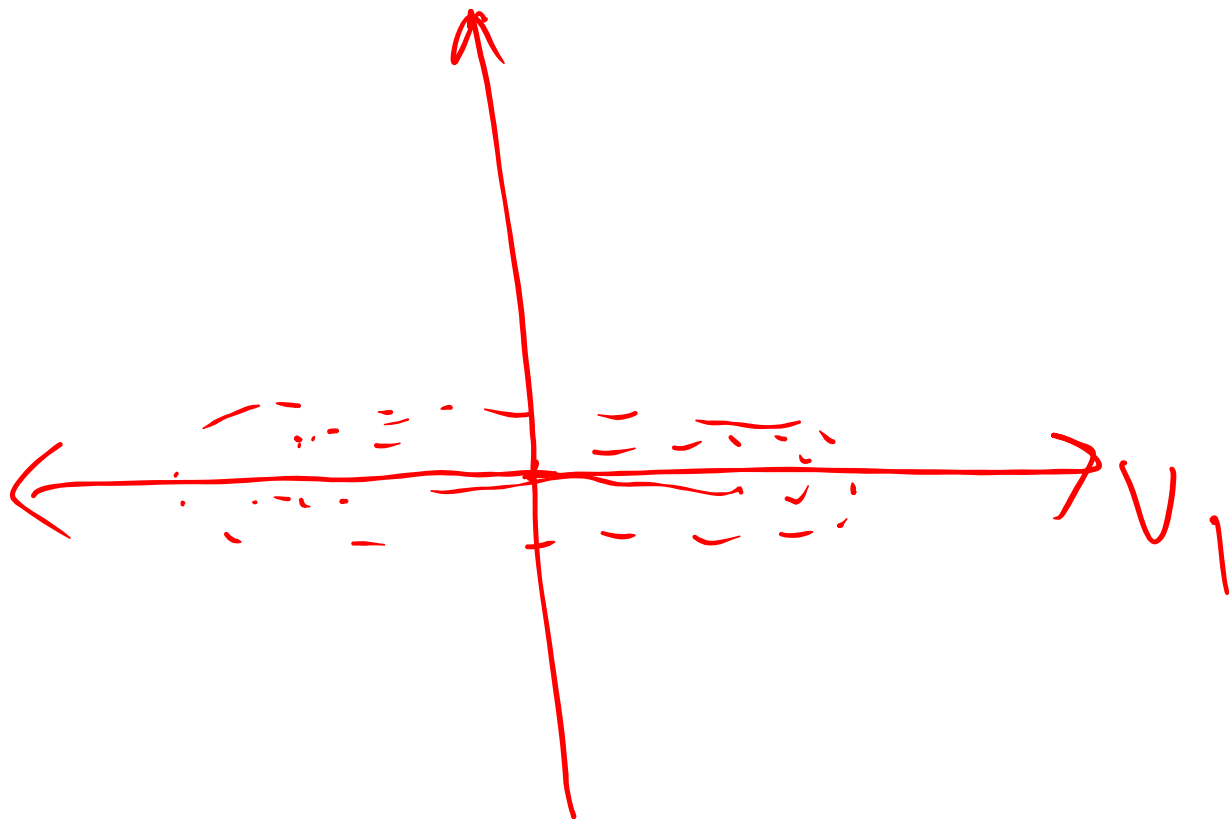
When using a test set, determine the input transformation from training data *only*.
**Rule:** lock away the test data until you have your final hypothesis.



$\mathcal{D}$

$\mathcal{D}_{\text{train}}$
input preprocessing
$\mathbf{z} = \Phi(\mathbf{x})$

$g(\mathbf{x}) = \tilde{g}(\Phi(\mathbf{x}))$

$\mathcal{D}_{\text{test}}$

$E_{\text{test}}$

Cumulative Profit %

snooping

no snooping

Day

# Principal Component Analysis

# (PCA)



$V_2$ (Noise)

$V_1$

$$x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} x^T e_1 \\ x^T e_2 \end{bmatrix}$$

Co-ordinates in natural basis:

$$x = \begin{bmatrix} x^T U_1 \\ x^T U_2 \end{bmatrix} = \begin{bmatrix} Z_1 \\ Z_2 \end{bmatrix} = Z$$
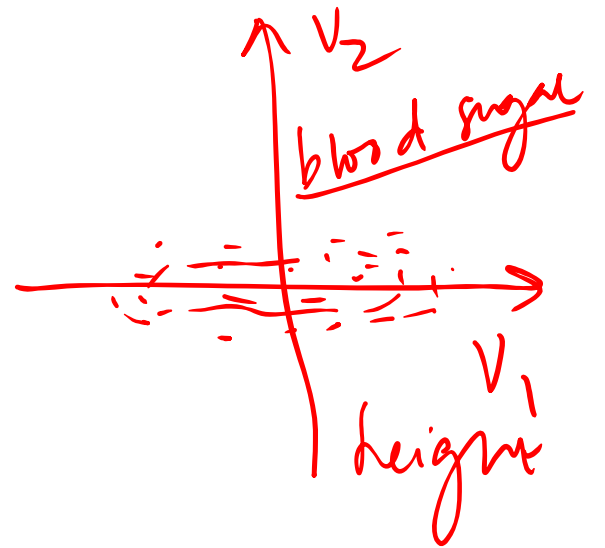
ORTHOGONAL          Rotation

$$x = x_1 e_1 + x_2 e_2 = Z_1 U_1 + Z_2 V_2$$

NOISE

$$x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \longrightarrow \begin{bmatrix} Z_1 \\ \cancel{Z_2} \end{bmatrix} \rightarrow \begin{bmatrix} Z_1 \end{bmatrix}$$

$$x \simeq Z_1 V_1 \quad (Z_2 V_2)$$

Step1 : Identify dominant directions.

Step2 : Throw away small dimensions.

i) Dimensionality Reduction (Always True)

ii) Irrelevant ~~Beta~~ Info. (Hope)

Drawback:

Dominant direction is $v$

$$Z = x^T v$$

$$\frac{1}{N} \sum_{n=1}^{N} Z_n^2$$

(Largest squared co-ordinals on avg.).

$$Z_n^2 = (x_n^T v)^2 = v^T x_n x_n^T v$$

$$\frac{1}{N} \sum_{n=1}^{N} v^T x_n x_n^T v = \frac{1}{N} v^T \left( \sum x_n x_n^T \right) v$$

$$= v^T \left( \frac{1}{N} \sum x_n x_n^T \right) v$$

$$= v^T \Sigma v \leftarrow$$

$$\Sigma \rightarrow v_1 \, v_2 \cdots v_d$$

eigenval $\rightarrow \lambda_1 \geq \lambda_2 \geq \cdots \lambda_d \geq 0$

$$v^* = v_1$$

1) Find $\Sigma$ and eigenvectors $\begin{rcases} U_1, U_2 \cdots V_d \\ \lambda_1, \lambda_2 \cdots \lambda_d \end{rcases}$ orthogonal basis

$$Z = \begin{bmatrix} x^T U_1 \\ x^T V_2 \\ \vdots \\ x^T V_d \\ \cdots \end{bmatrix}$$

$\longrightarrow$ max

direction of decreasing variance

$\longrightarrow$ max variance $\perp$ to $U_1$

2) $K$ - coordinates : Top $K$ PCA directions.

$$x \simeq Z_1 V_1 + Z_2 V_2 \cdots Z_k V_k$$

Reconstruction error
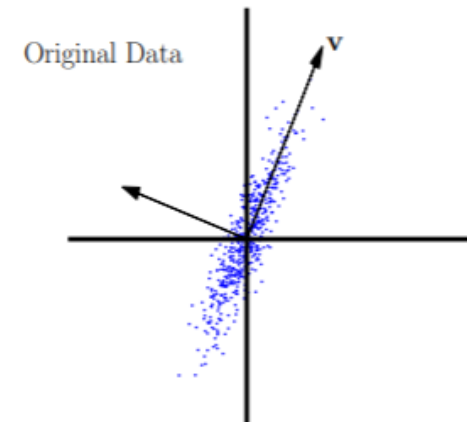
# The Principal Components

$$z_1 = \mathbf{x}^{\mathrm{T}} \mathbf{v}_1$$

$$z_2 = \mathbf{x}^{\mathrm{T}} \mathbf{v}_2$$

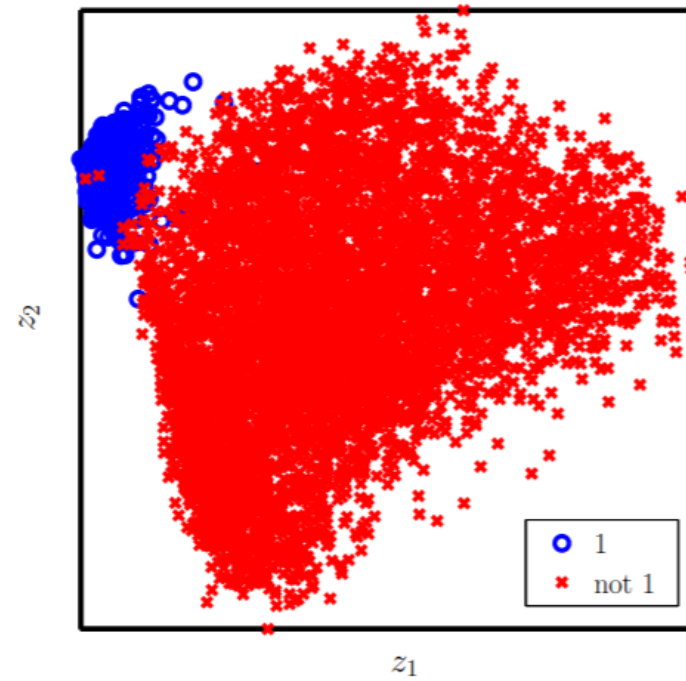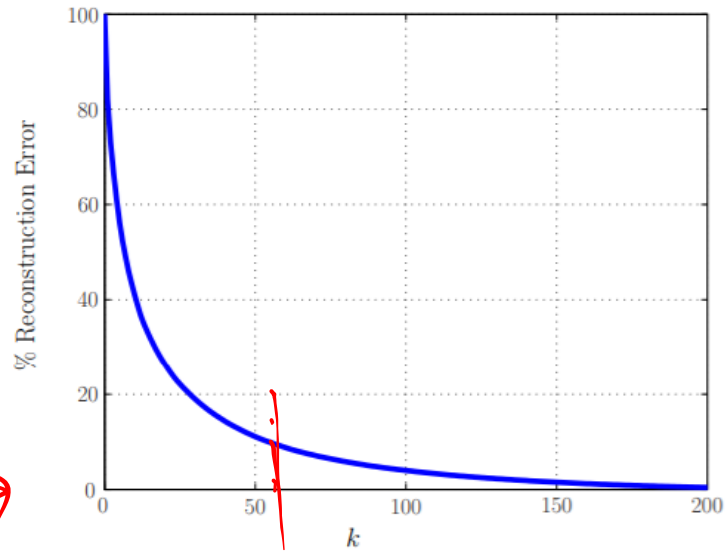$$z_3 = \mathbf{x}^{\mathrm{T}} \mathbf{v}_3$$

$$\vdots$$

$\mathbf{v}_1, \mathbf{v}_2, \cdots, \mathbf{v}_d$ are the eigenvectors of $\Sigma$ with eigenvalues $\lambda_1 \geq \lambda_2 \geq \cdots \geq \lambda_d$

**Theorem [Eckart-Young].** These directions give best reconstruction of data; also capture maximum variance.

Original Data

# PCA Features for Digits Data



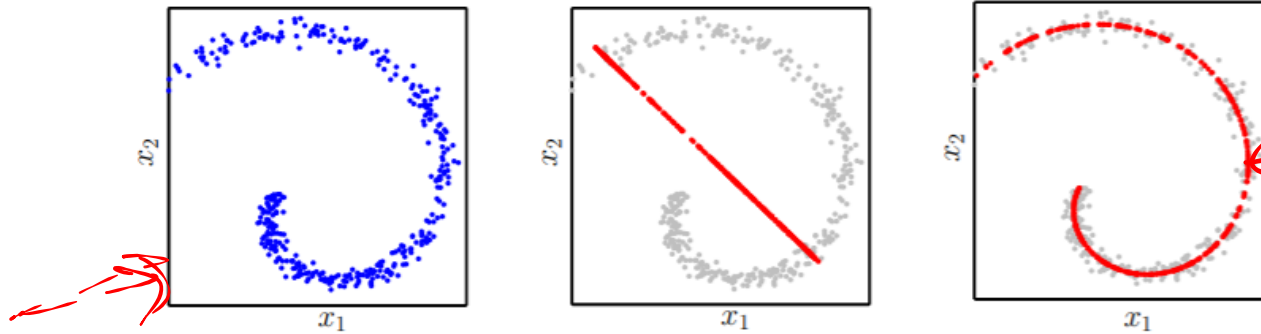Principal components are **automated**

Captures dominant directions of the data.

May not capture dominant dimensions for $f$.
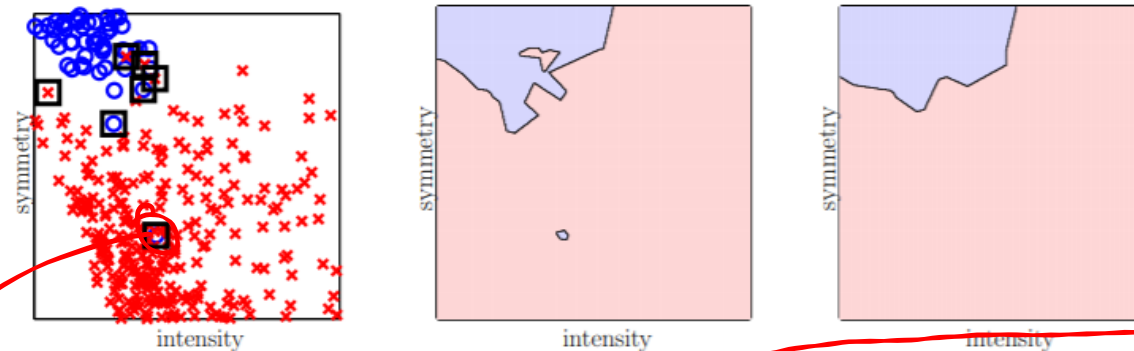
# Other Learning Aides

### 1. **Nonlinear dimension reduction:**



### 2. **Hints (invariances and prior information):**

rotational invariance, monotonicity, symmetry, ....

### 3. **Removing noisy data:**



### 4. **Advanced validation techniques:** Rademacher and Permutation penalties

More efficient than CV, more convenient and accurate than VC.

# Thanks!