
Dimensionality Reduction (DR) (Feature Extraction)

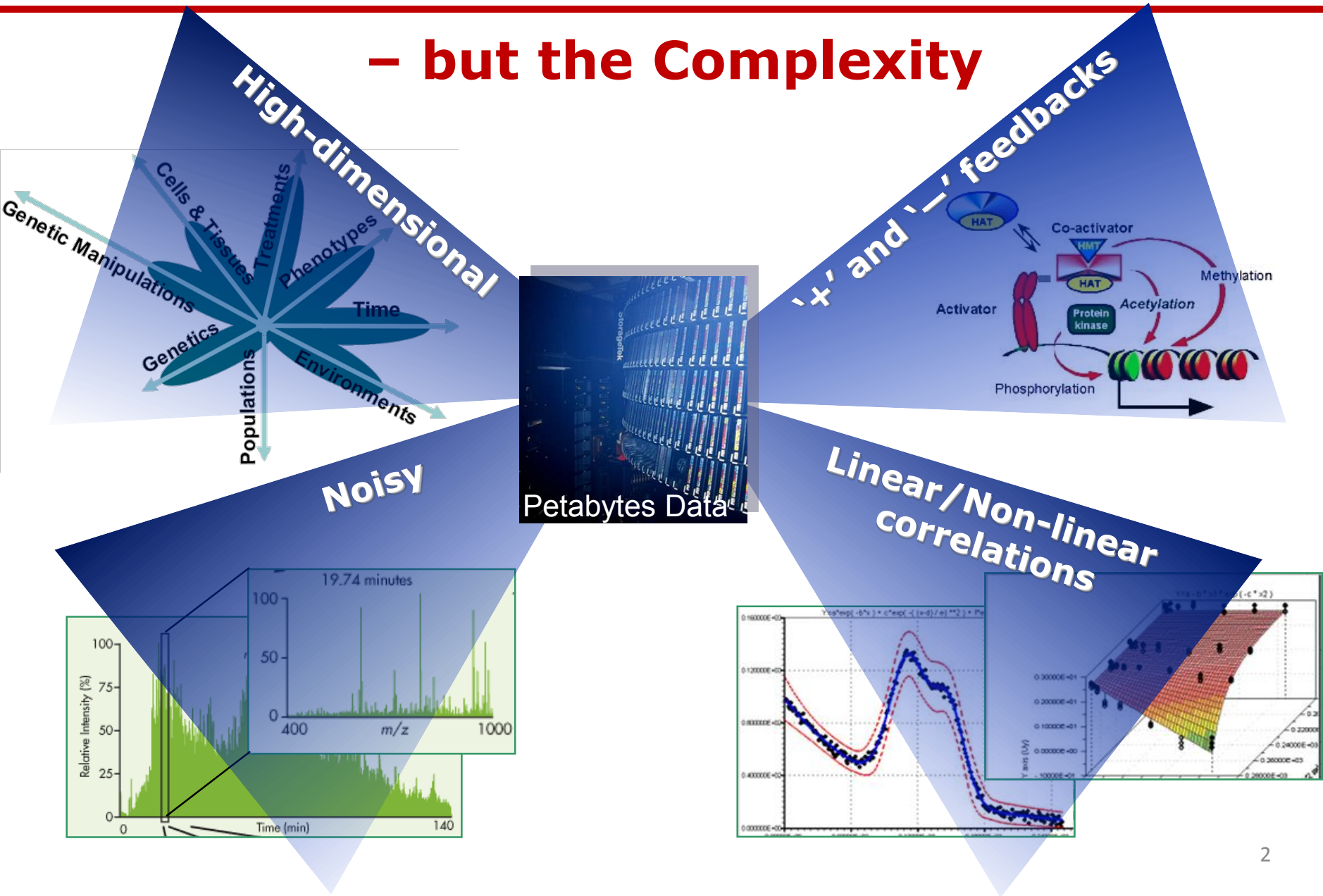
Nagiza F. Samatova, samatova@csc.ncsu.edu

Professor, Department of Computer Science
North Carolina State University

Senior Scientist, Computer Science & Mathematics Division
Oak Ridge National Laboratory

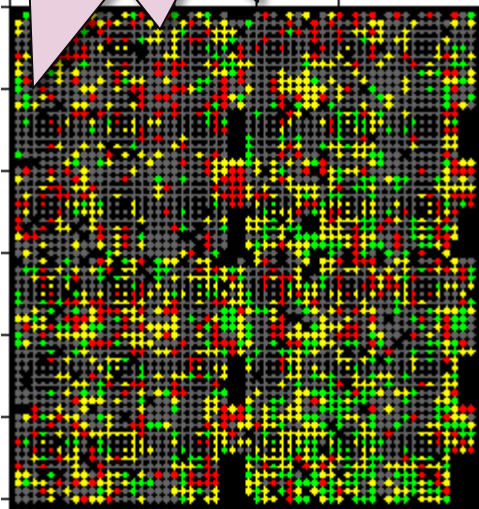
It is not just the Data Size

– but the Complexity



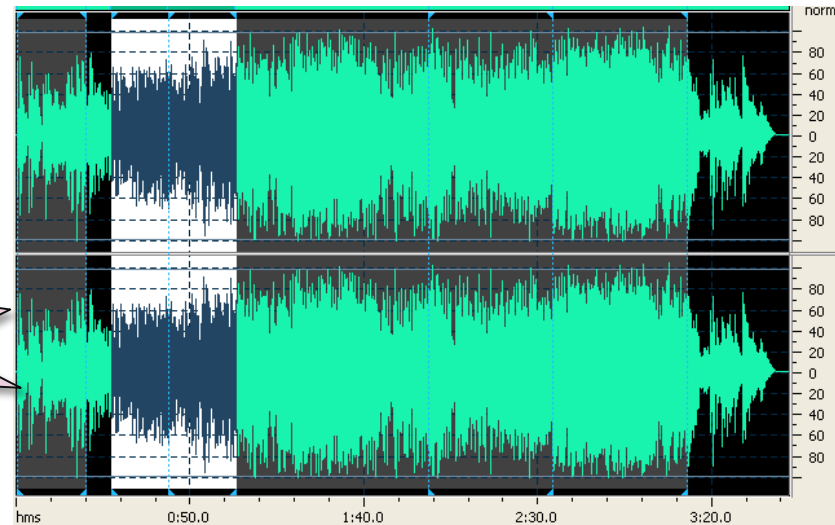
High-dimensional and Noisy Data

**Microarray data:
one sample point
has thousands of
genes and one
gene point has
tens of samples**



**Images contain
a lot of
information in
the form of
pixels**

**Noisy audio
data**



Challenges

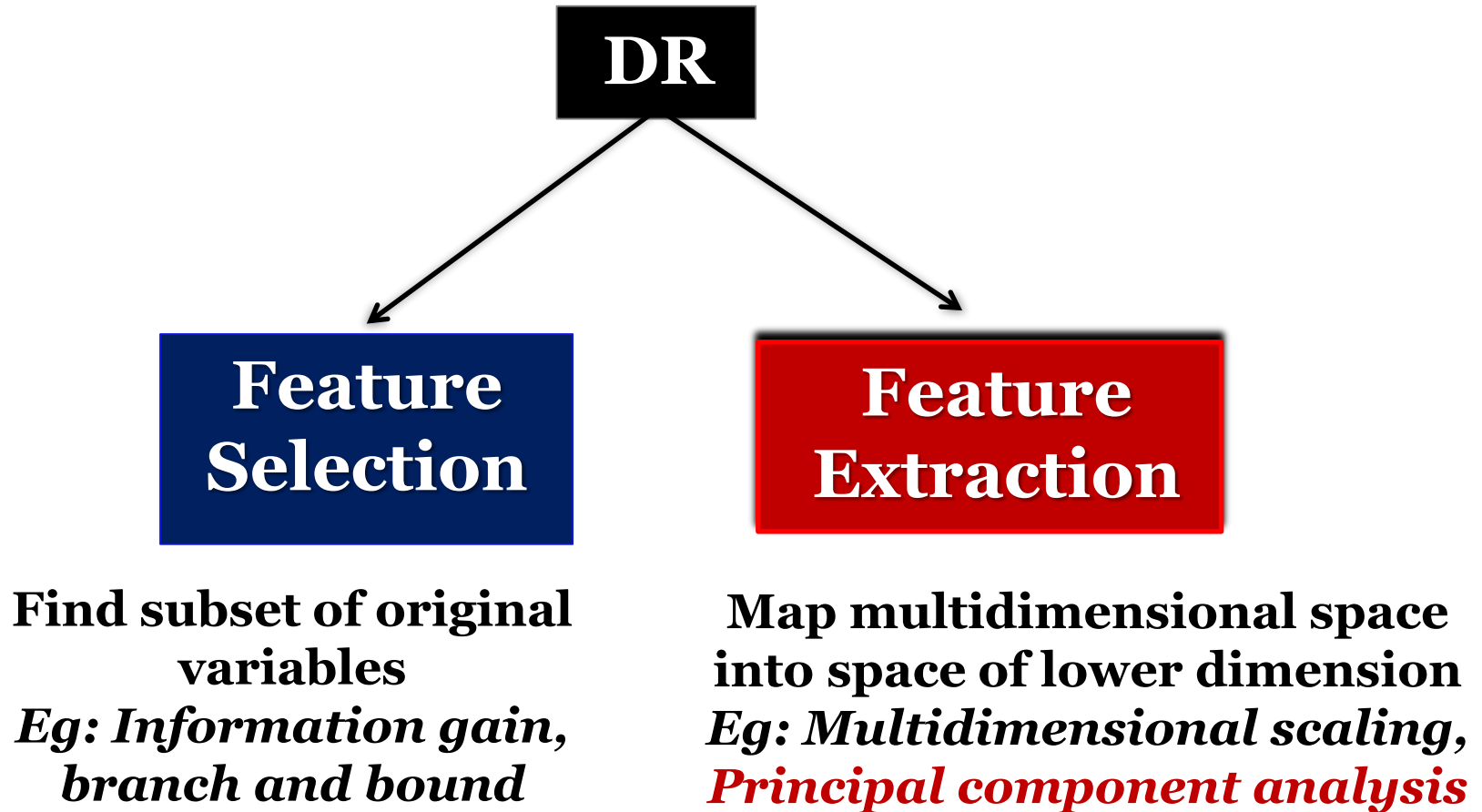
- **Not all the measured variables are important for understanding the underlying “interesting” phenomena – could complicate the process of data analysis**
- **Difficult to visualize**
- **High computational cost (time / space)**

Dimensionality Reduction as a Solution

- **Objective:** To transform data in high-dimensional space to a corresponding representation in some low-dimensional space, while “best” preserving the information
- Given dataset with n objects,

$$\begin{array}{ccc} X \in \mathbb{R}^{n \times m} & \xrightarrow{\text{DR}} & Y \in \mathbb{R}^{n \times p} \\ m \text{ dimensions} & & p \text{ dimensions} \\ & & \text{where } \mathbf{p} \ll \mathbf{m} \end{array}$$

Classification of DR techniques



Dimension Reduction, DR (informal)

- **Given**: a collection of records in d -dimensional space
 - Each record contains a set of d attributes
- **Find**: a lower-dimensional ($k < d$) representation of this data that:
 - Optimizes some objective function and
 - Meets a given set of constraints.
- **Assumption**: In a high-dimensional space, the data is often intrinsically low-dimensional. We can exploit this to avoid the curse of high dimensionality.

Key Questions for Different DR Methods

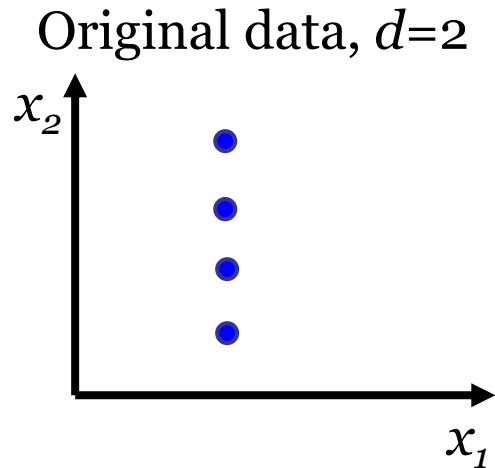
- **What is the OBJECTIVE FUNCTION?**
- **What are the CONSTRAINTS?**
- **How to solve the OPTIMIZATION problem?**

- **Different Dimension Reduction methods answer these questions differently**

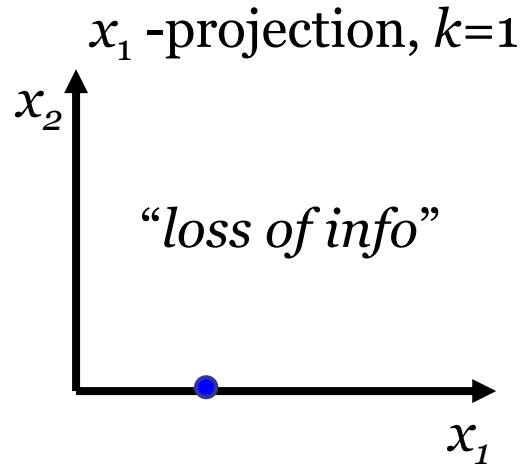
Motivation for Dimension Reduction

- Decrease the **computational cost** for other data mining tasks:
 - Proximity measure calculations: $O(d) \rightarrow O(k)$
- Reduce the **noise** in the data
- Improve the **accuracy** of predictive models
- Reduce **collinearity** among variables/features

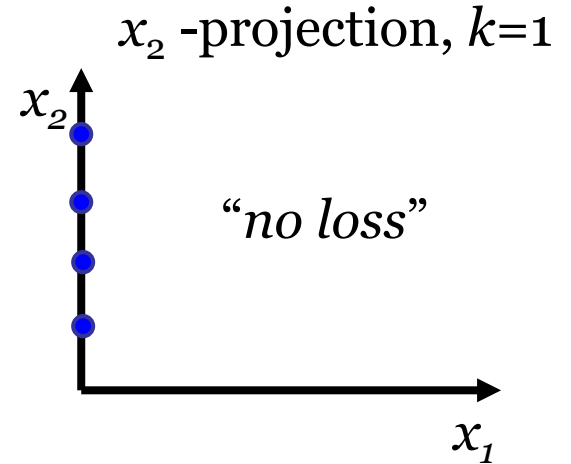
Example 1: $d=2 \rightarrow k=1$



$$X_{4 \times 2}$$



$$X'_{4 \times 2} = X_{4 \times 2} \cdot \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix}$$



$$X'_{4 \times 2} = X_{4 \times 2} \cdot \begin{bmatrix} 0 & 0 \\ 0 & 1 \end{bmatrix}$$

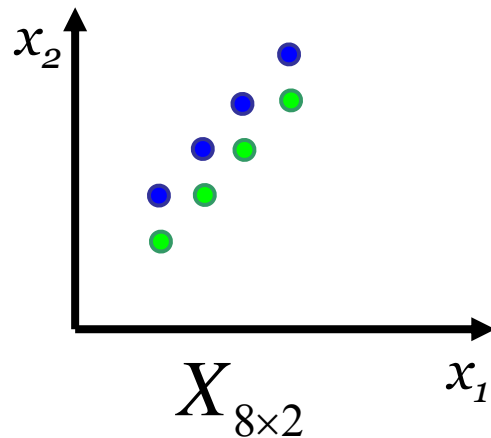
$$X'_{m \times d} = X_{m \times d} \cdot P_{d \times d}$$

Which projection is "better"?

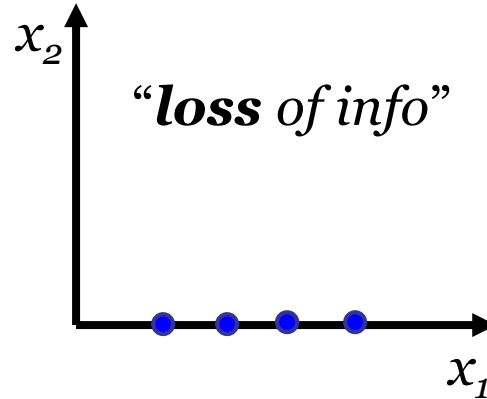
$$P_{d \times d} = \begin{bmatrix} 1 & 0 & \dots & 0 \\ 0 & 0 & \dots & 0 \\ \dots & \dots & \dots & \dots \\ 0 & 0 & 0 & 1 \end{bmatrix} \text{ -- projection matrix; some diagonal elements are 0}$$

Example 2: Linear, Orthogonal Projection

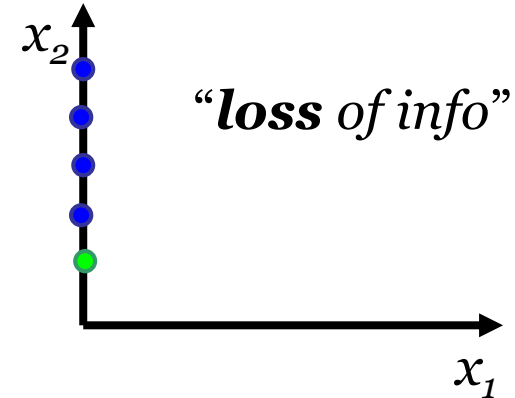
Original data, $d=2$



x_1 -projection, $k=1$



x_2 -projection, $k=1$

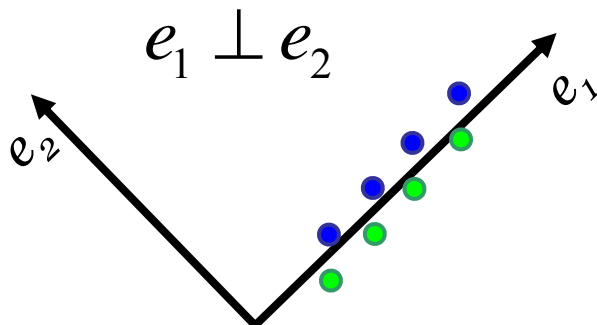


Is there a “better” projection?

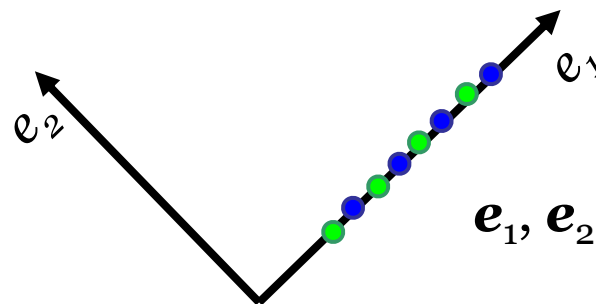
Projection:

- **Linear**, e_1 – line
- **Orthogonal**

Another **basis**, $d=2$
(rotate coord. system)

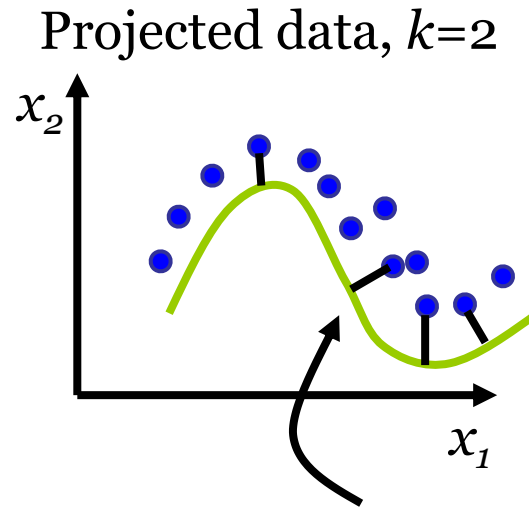
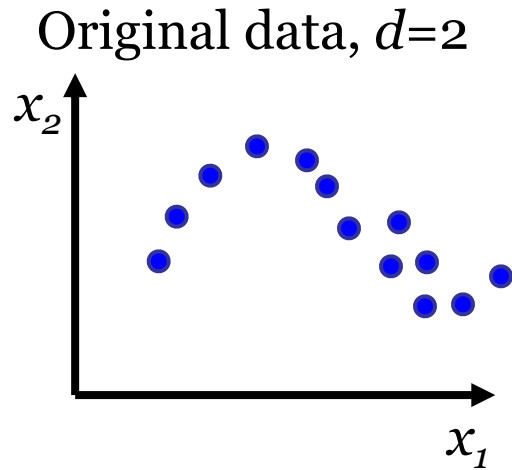


e_1 -projection,
 $k=1$



e_1, e_2 – eigenvectors of X

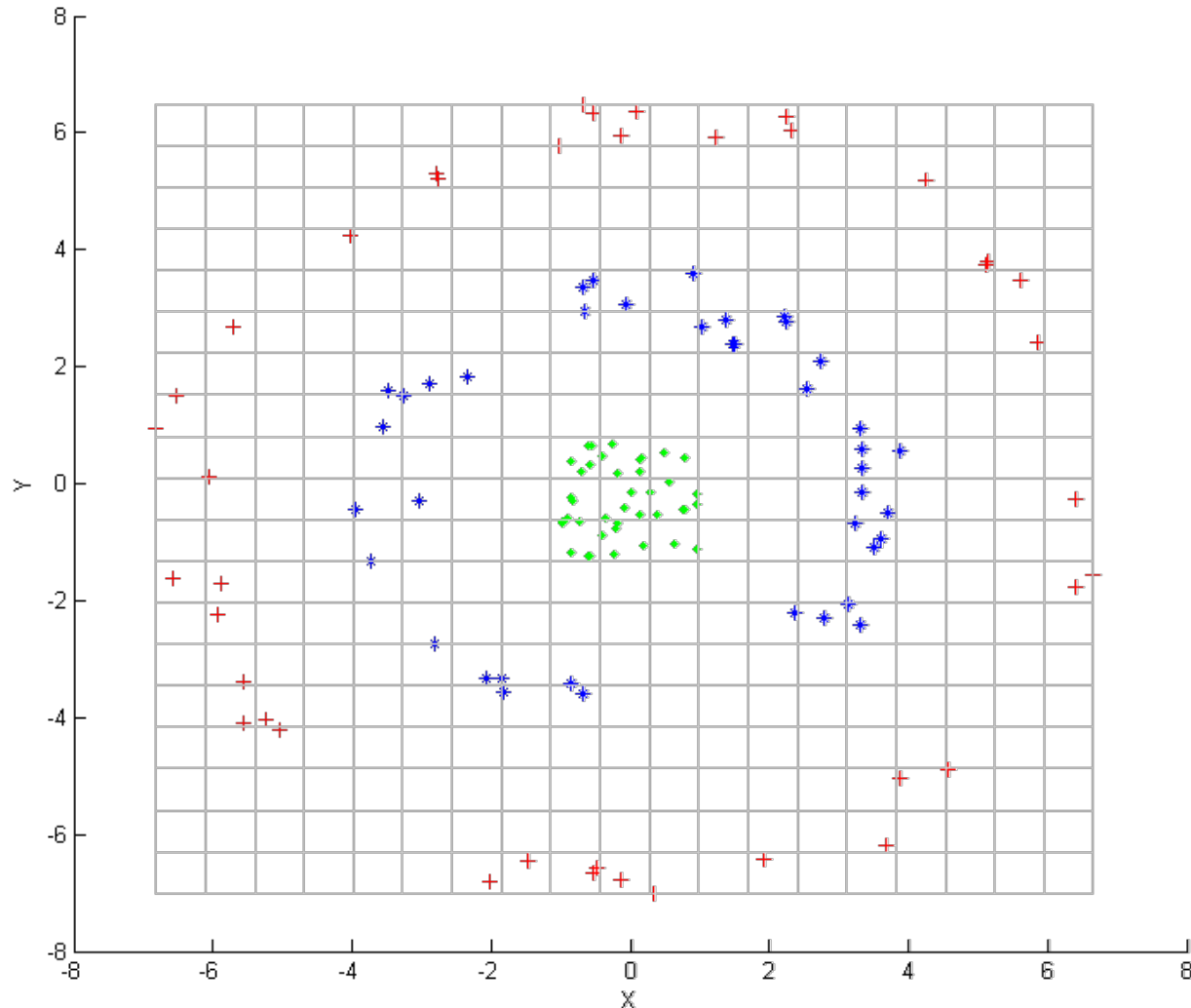
Example 3: Non-linear projection



NON-linear projection

Example 4: DR for Labeled Data

What is a “better” projection?



Types of Projection

- **Linear vs. NON-linear**
- **Orthogonal vs. non-orthogonal**
- **Unsupervised (unlabeled data) vs. supervised (labeled data)**

What is “better” projection: $d \rightarrow k$ ($k < d$)?

- Many definitions are possible
- **Definition 1:**
 - Projection that **maximizes the VARIANCE** of the data in its target **k** -dimensional projection

R Example

```
data (iris)
```

```
iris
```

```
X = iris[, 1:4]
```

```
pca = princomp(X, center=TRUE)
```

```
pca
```

```
plot(pca) # screeplot
```

```
loadings(pca) # matrix of eigenvectors
```

```
summary(pca) # check proportion of variance
```

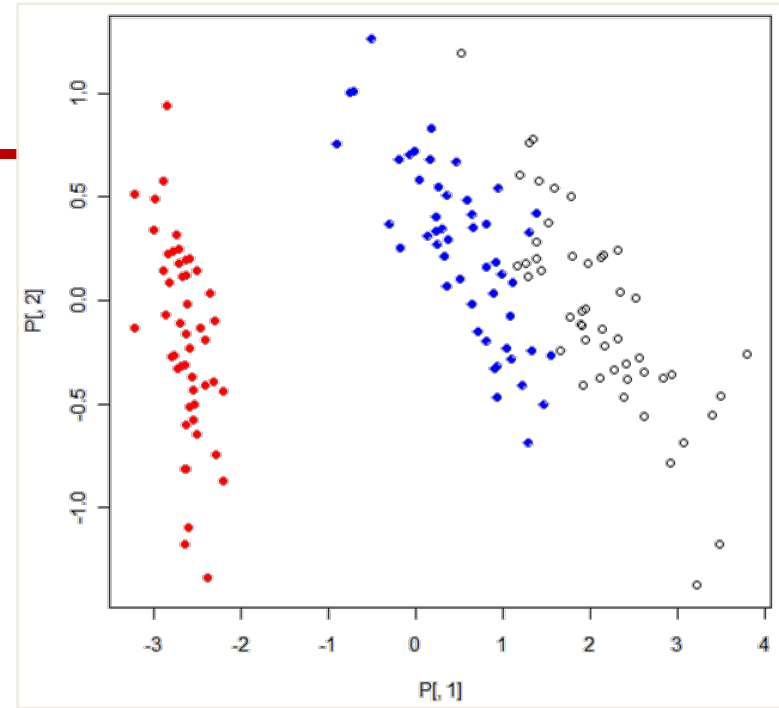
```
P=pca$scores # projection of X onto eigenvectors
```

```
plot(P[,1], P[,2])
```

```
points(P[1:50,1], P[1:50,2], col="red")
```

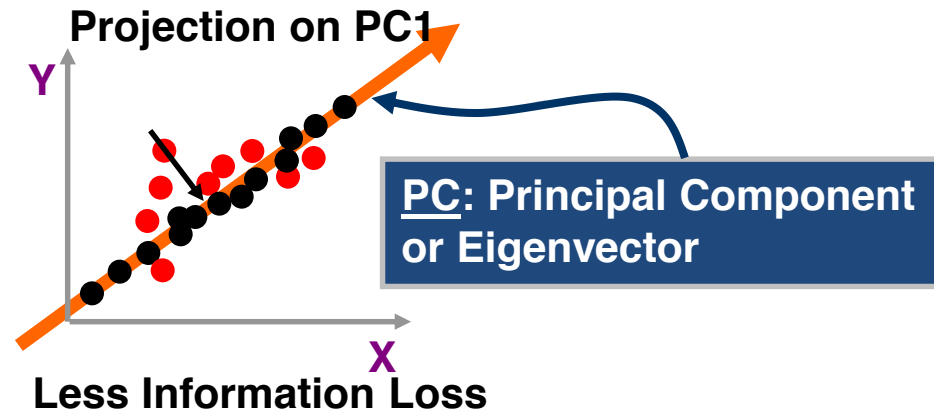
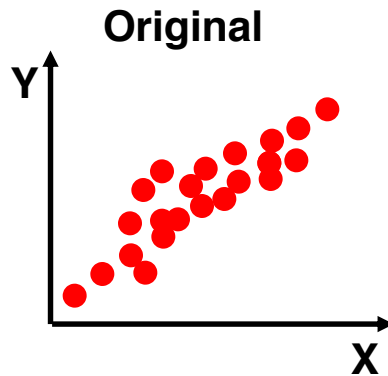
```
points(P[51:100,1], P[51:100,2], col="blue")
```

```
# plot the same for X
```



PCA: Linear Orthogonal DR

Principal Component Analysis (**PCA**) finds **intrinsic** dimensionality and allows for low-dimensional representation.



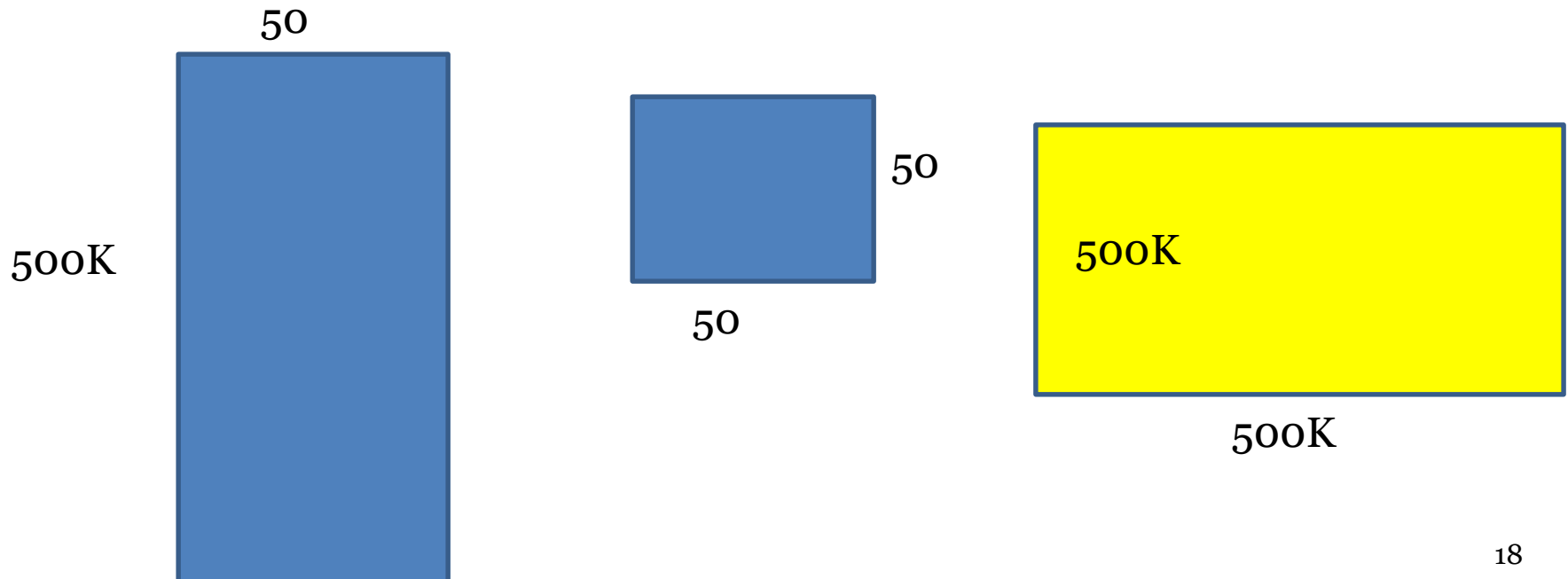
Covariance Matrix

$$\Sigma_{d \times d} = \text{cov}(X_{m \times d}) = \frac{1}{m-1} (X_{m \times d}^T X_{m \times d} - m \cdot \bar{\mathbf{x}} \cdot \bar{\mathbf{x}}^T)$$

Covariance matrix

Centroid: Column means

Exercise: Go over the script **covariance.R**



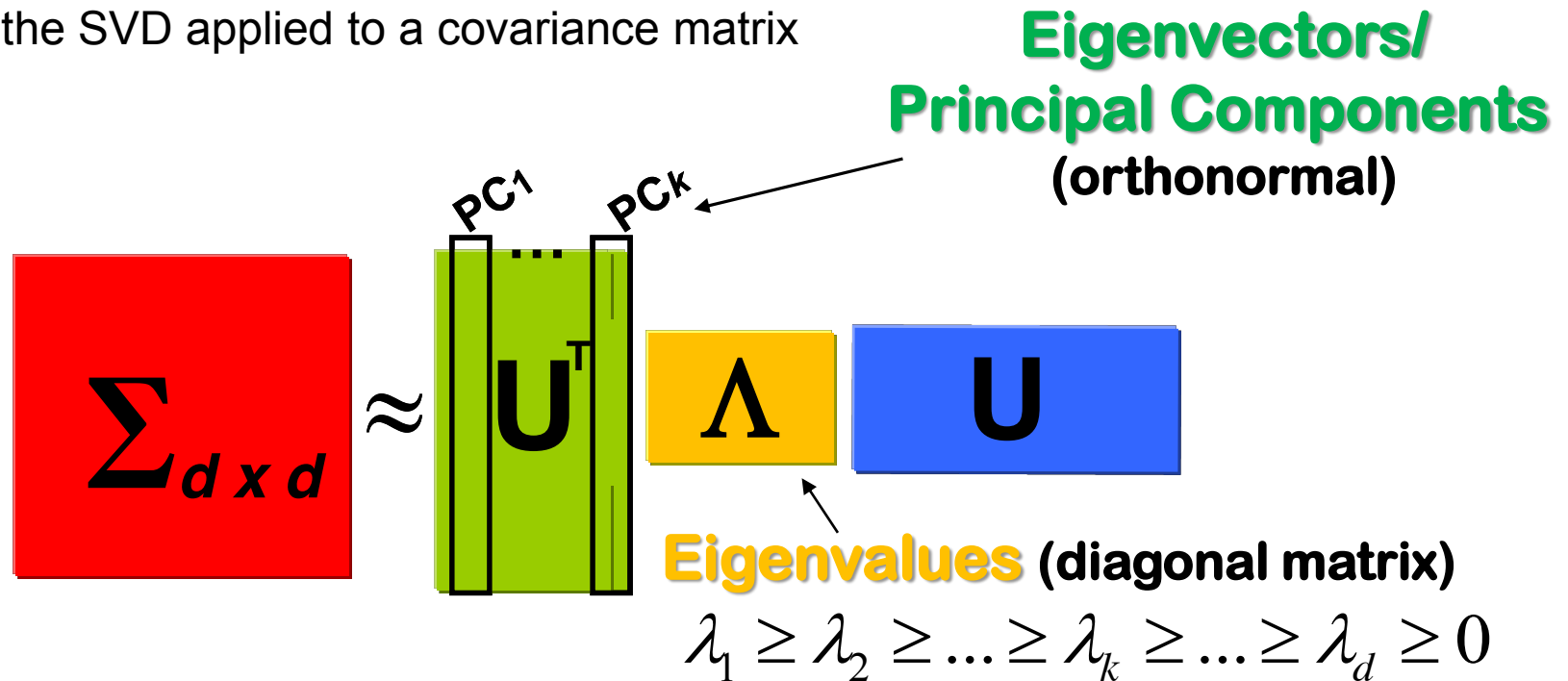
PCA is the SVD of Covariance Matrix

Singular Value Decomposition (SVD)

The technique underlying PCA analysis

PCA = SVD (Covariance Matrix)

PCA is the SVD applied to a covariance matrix



Details

`princomp` is a generic function with "formula" and "default" methods.

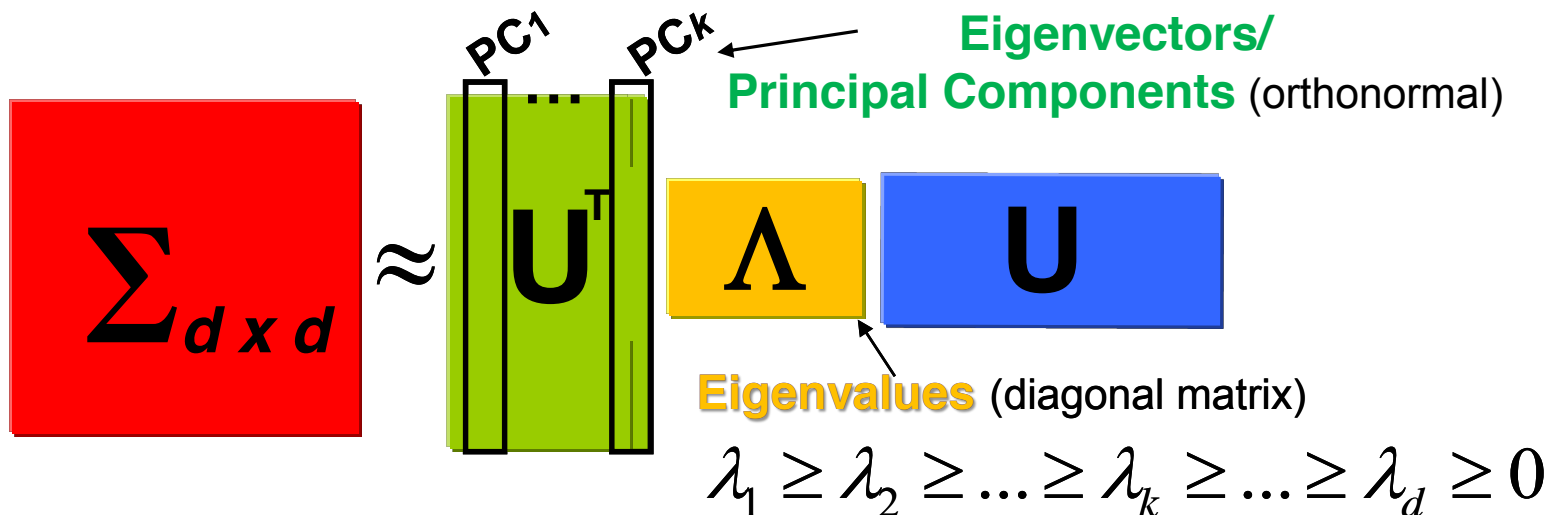
The calculation is done using [eigen](#) on the correlation or covariance matrix, as determined by [cor](#).

PC is a weighted linear sum of original features

Extracted Feature → PCA-based Feature Extraction:

$$PC = w_1 * f_1 + w_2 * f_2 + \dots + w_d * f_d$$

The magnitude of each weight indicates how important the corresponding feature is
→ it could be used as a **feature selection** technique!



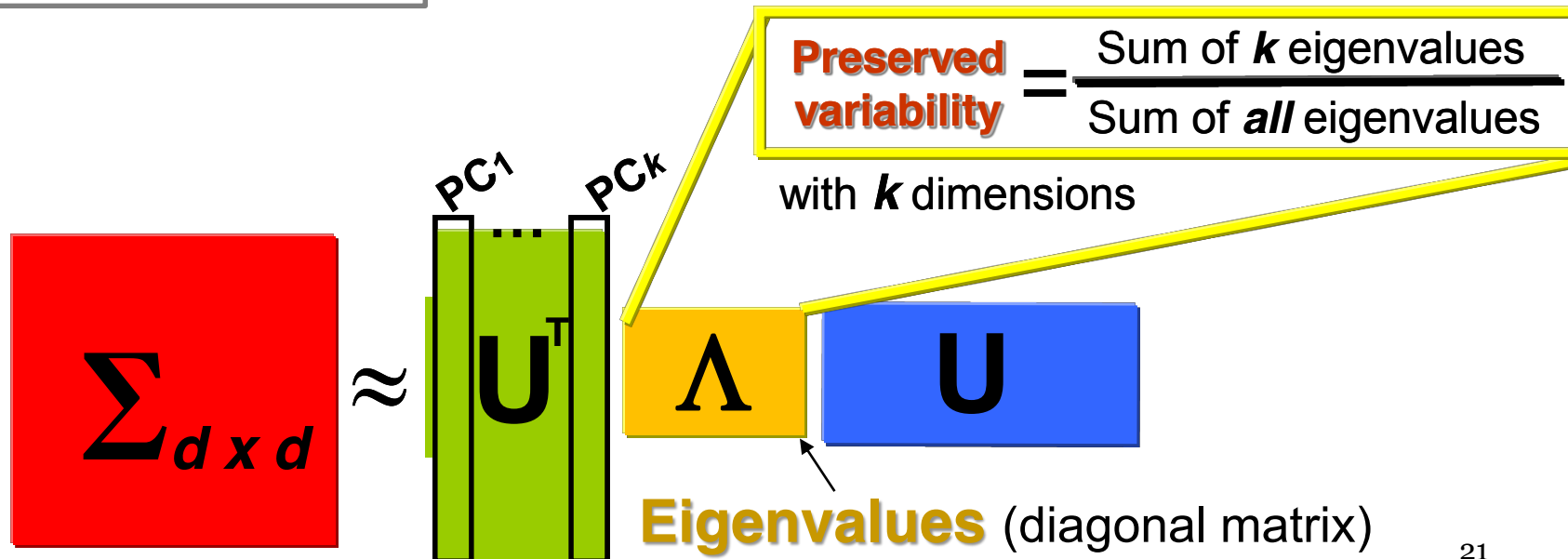
Preserved Variability for top-k PCs

Percentage of variability preserved
if the first **k** PCs are used for **projection**:

SORTED Eigenvalues

$$\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_k \geq \dots \geq \lambda_d \geq 0$$

$$\frac{\sum_{i=1}^k \lambda_i}{\sum_{i=1}^d \lambda_i} = \frac{\sum_{i=1}^k \lambda_i}{\text{trace}(\Sigma)}$$

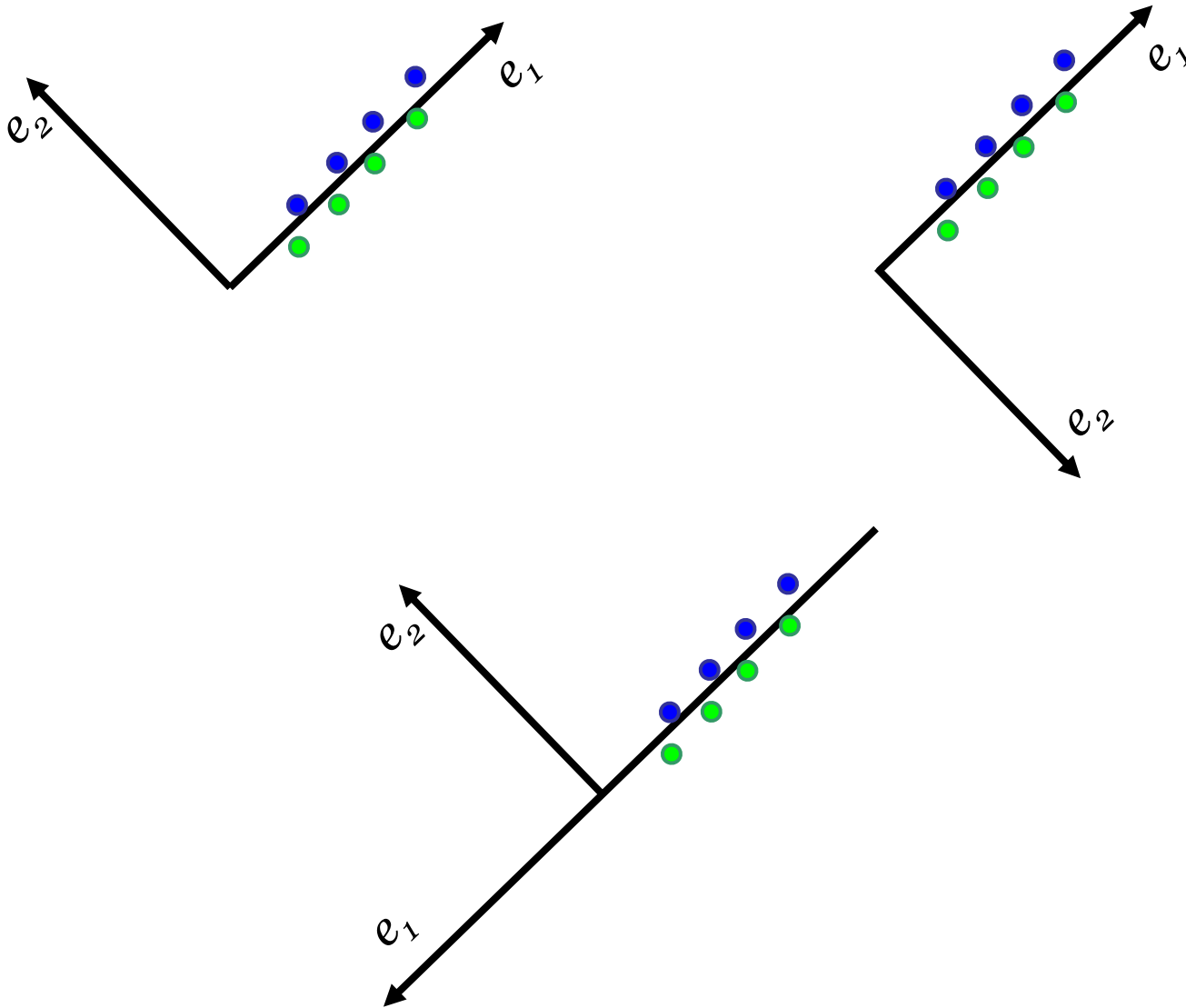


Key Points about PCA

- **PCA = SVD (Covariance Matrix, Σ)**
- **Linear, orthogonal projection:**
 - “Best” linear orthogonal k -dimensional ($k < d$) view of data
- **How “good” the k -dimensional view is:**

$$\frac{\sum_{i=1}^k \lambda_i}{\sum_{i=1}^d \lambda_i} = \frac{\sum_{i=1}^k \lambda_i}{\text{trace}(\Sigma)}$$

Eigenvectors are NOT unique



Eigenvalues: importance of eigenvectors

$$\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_k \geq \dots \geq \lambda_d \geq 0$$

Importance of Principal Components (PC)/Eigenvectors/Loadings

PC1 preserves more variance than PC2

PC2 preserves more variance than PC3

....

Proportion of Variance Preserved if only k PCs are used:

$$\frac{\sum_{i=1}^k \lambda_i}{\sum_{i=1}^d \lambda_i} = \frac{\sum_{i=1}^k \lambda_i}{\text{trace}(\Sigma)}$$