

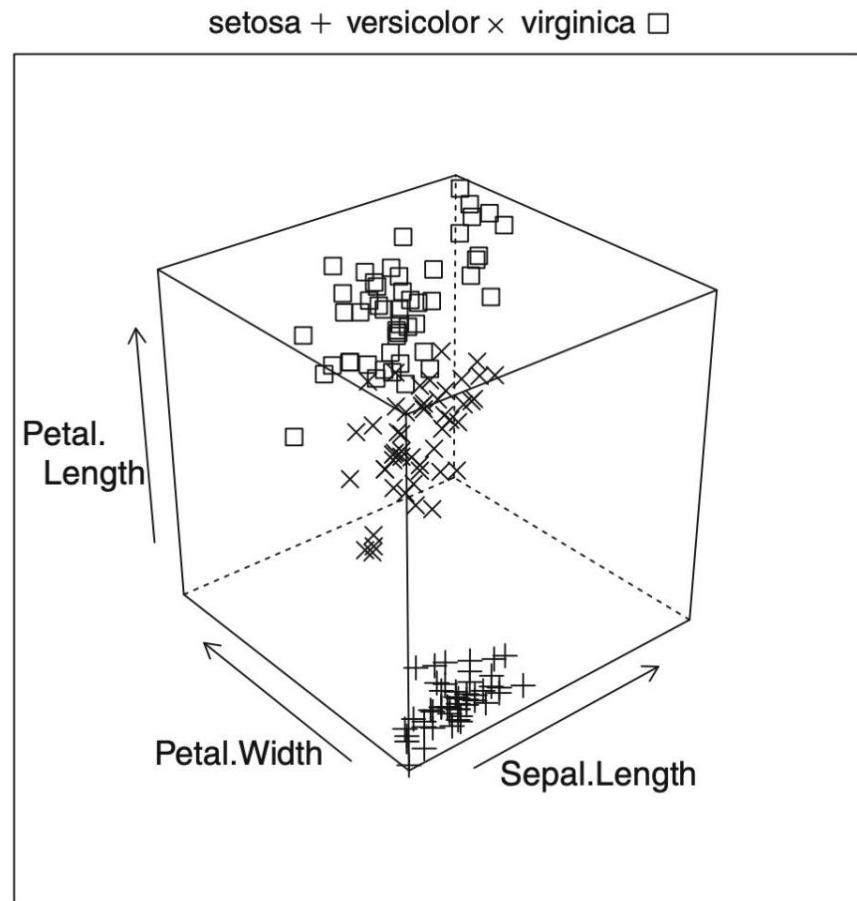
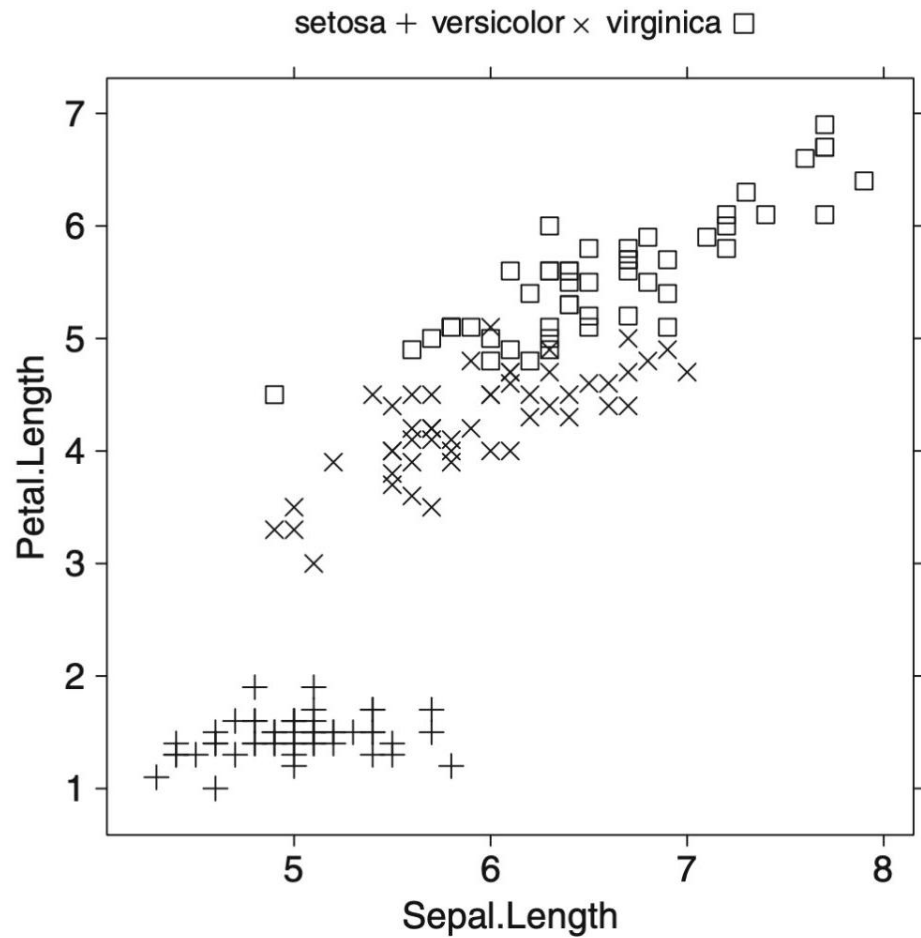
Computing for Medicine

Data Science: PCA and Feature Selection

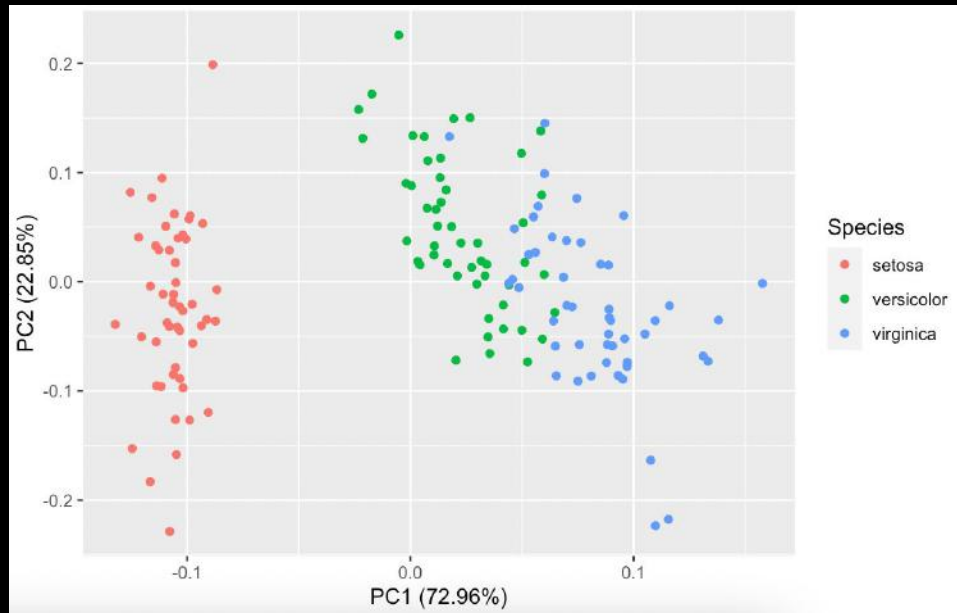
Tavpritesh Sethi



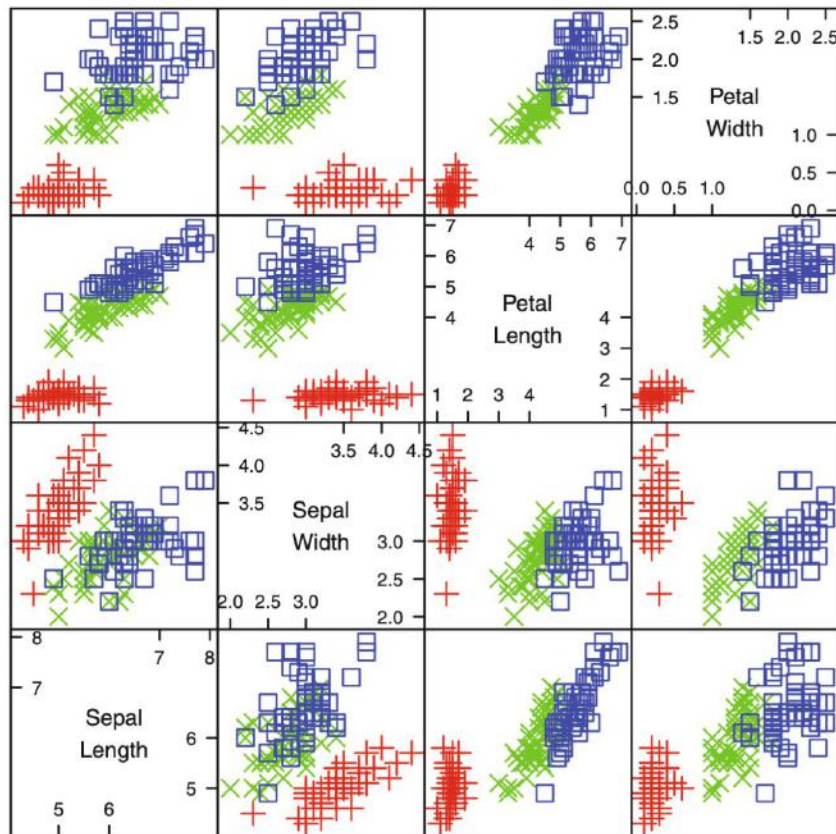
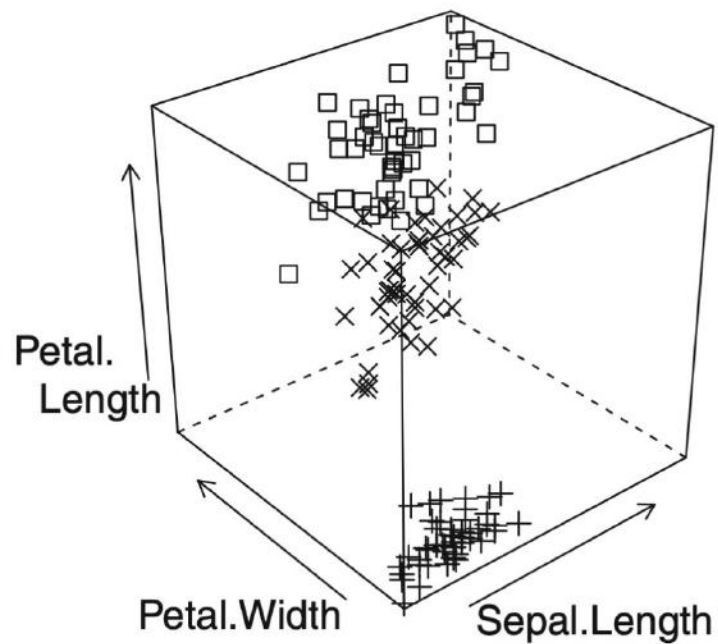
High Dimensional Data Principal Component Analysis



```
1 library(ggfortify)
2 df <- iris[1:4]
3 pca_res <- prcomp(df, scale. = TRUE)
4
5 autoplot(pca_res)
6 autoplot(pca_res, data = iris, colour = 'Species')
```

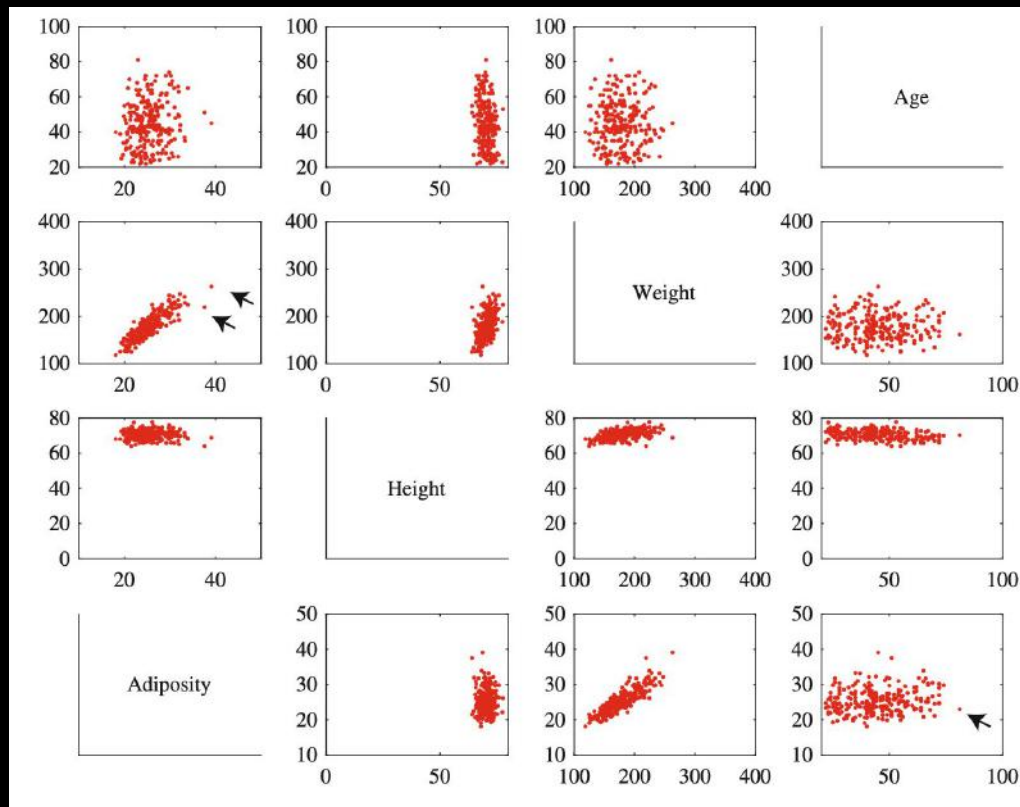


setosa + versicolor × virginica □



Scatter Plot Matrix

Height-Weight Dataset: Original Scatterplots



Challenges of High Dimensional Data

Remember This: *High dimensional data does not behave in a way that is consistent with most people's intuition. Points are always close to the boundary and further apart than you think. This property makes a nuisance of itself in a variety of ways. The most important is that only the simplest models work well in high dimensions.*

Recap: Covariance and Correlation

$$\text{cov}(\{x\}, \{y\}) = \frac{\sum_i (x_i - \text{mean}(\{x\}))(y_i - \text{mean}(\{y\}))}{N}$$

$$\text{corr}(\{(x, y)\}) = \frac{\text{cov}(\{x\}, \{y\})}{\sqrt{\text{cov}(\{x\}, \{x\})} \sqrt{\text{cov}(\{y\}, \{y\})}}.$$

Covariance Matrix

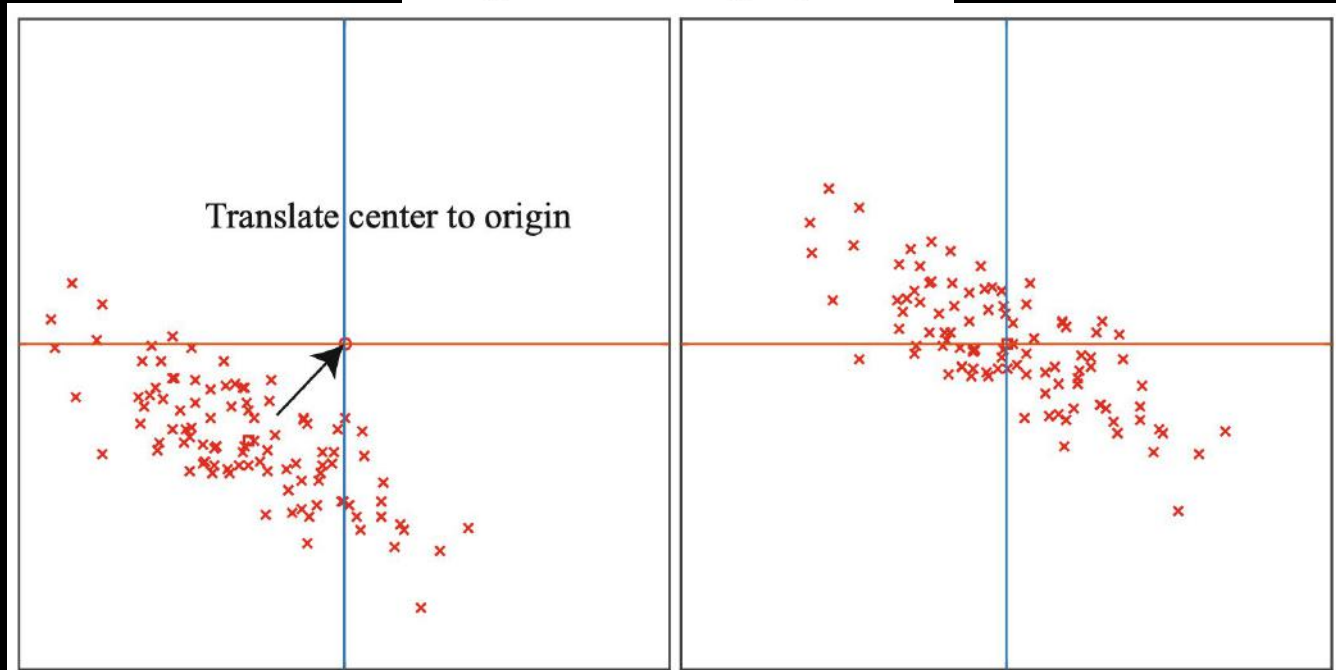
$$\text{Covmat}(\{\mathbf{x}\}) = \frac{\sum_i (\mathbf{x}_i - \text{mean}(\{\mathbf{x}\}))(\mathbf{x}_i - \text{mean}(\{\mathbf{x}\}))^T}{N}$$

Useful Facts: 4.3 *Properties of the Covariance Matrix*

- The j , k 'th entry of the covariance matrix is the covariance of the j 'th and the k 'th components of \mathbf{x} , which we write $\text{cov}(\{x^{(j)}\}, \{x^{(k)}\})$.
- The j , j 'th entry of the covariance matrix is the variance of the j 'th component of \mathbf{x} .
- The covariance matrix is symmetric.
- The covariance matrix is always positive semidefinite; it is positive definite, *unless* there is some vector \mathbf{a} such that $\mathbf{a}^T(\mathbf{x}_i - \text{mean}(\{\mathbf{x}_i\})) = 0$ for all i .

Affine Transformation

$$\mathbf{m}_i = \mathcal{A}\mathbf{x}_i + \mathbf{b}.$$



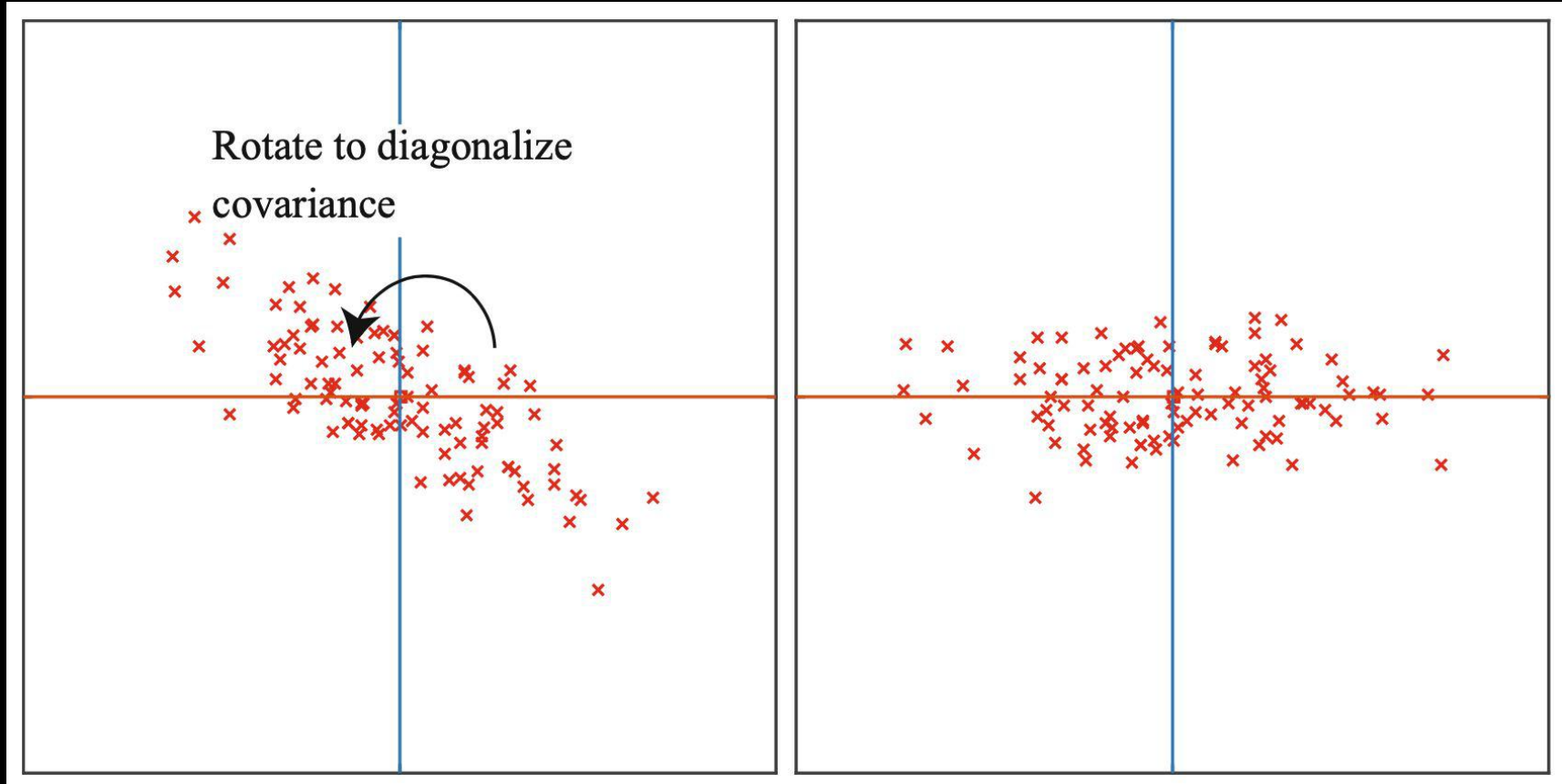
Mean and Covariance matrix under Affine Transformations

$$\begin{aligned}\text{mean}(\{\mathbf{m}\}) &= \text{mean}(\{\mathcal{A}\mathbf{x} + \mathbf{b}\}) \\ &= \mathcal{A}\text{mean}(\{\mathbf{x}\}) + \mathbf{b},\end{aligned}$$

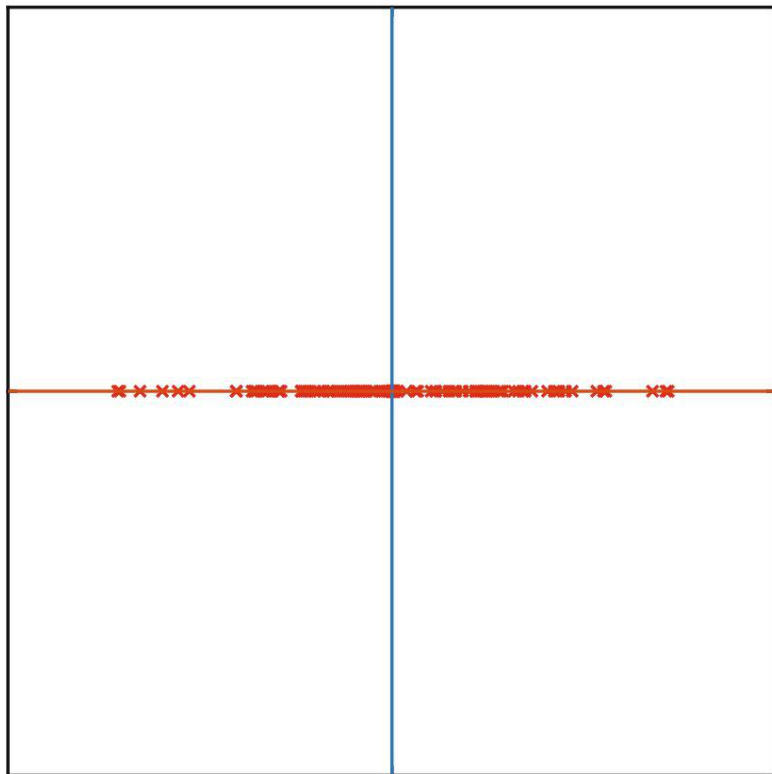
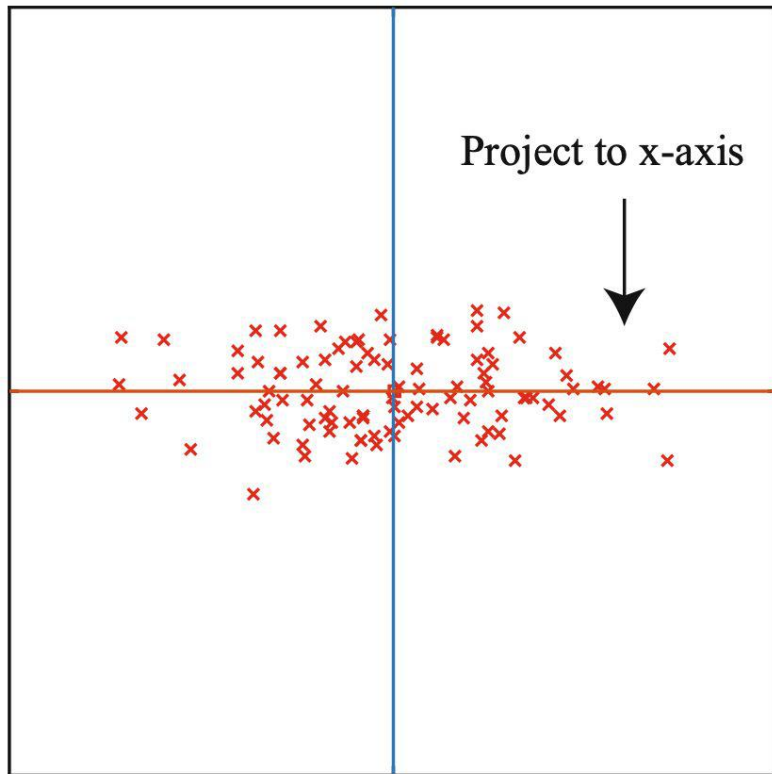
$$\begin{aligned}\text{Covmat}(\{\mathbf{m}\}) &= \text{Covmat}(\{\mathcal{A}\mathbf{x} + \mathbf{b}\}) \\ &= \frac{\sum_i (\mathbf{m}_i - \text{mean}(\{\mathbf{m}\}))(\mathbf{m}_i - \text{mean}(\{\mathbf{m}\}))^T}{N} \\ &= \frac{\sum_i (\mathcal{A}\mathbf{x}_i + \mathbf{b} - \mathcal{A}\text{mean}(\{\mathbf{x}\}) - \mathbf{b})(\mathcal{A}\mathbf{x}_i + \mathbf{b} - \mathcal{A}\text{mean}(\{\mathbf{x}\}) - \mathbf{b})^T}{N} \\ &= \frac{\mathcal{A} \left[\sum_i (\mathbf{x}_i - \text{mean}(\{\mathbf{x}\}))(\mathbf{x}_i - \text{mean}(\{\mathbf{x}\}))^T \right] \mathcal{A}^T}{N} \\ &= \mathcal{A}\text{Covmat}(\{\mathbf{x}\})\mathcal{A}^T.\end{aligned}$$

Principal Component Analysis

Rotated data preserves angles, translations preserve distances



Projections on axes

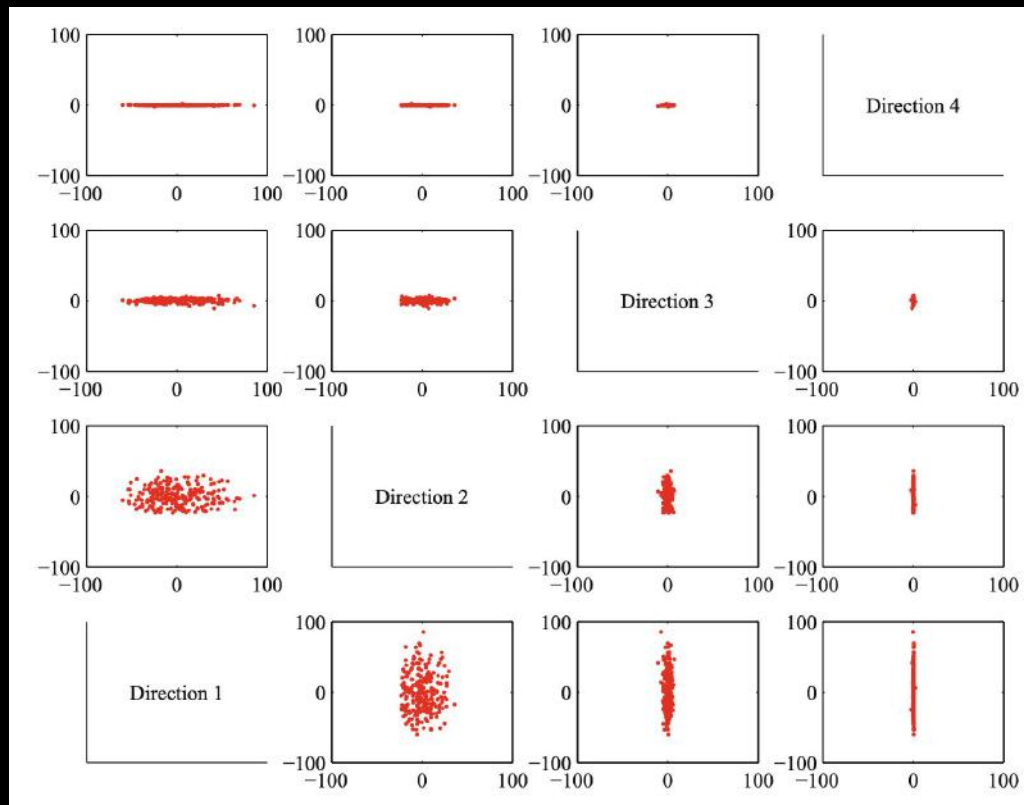


Singular Value Decomposition

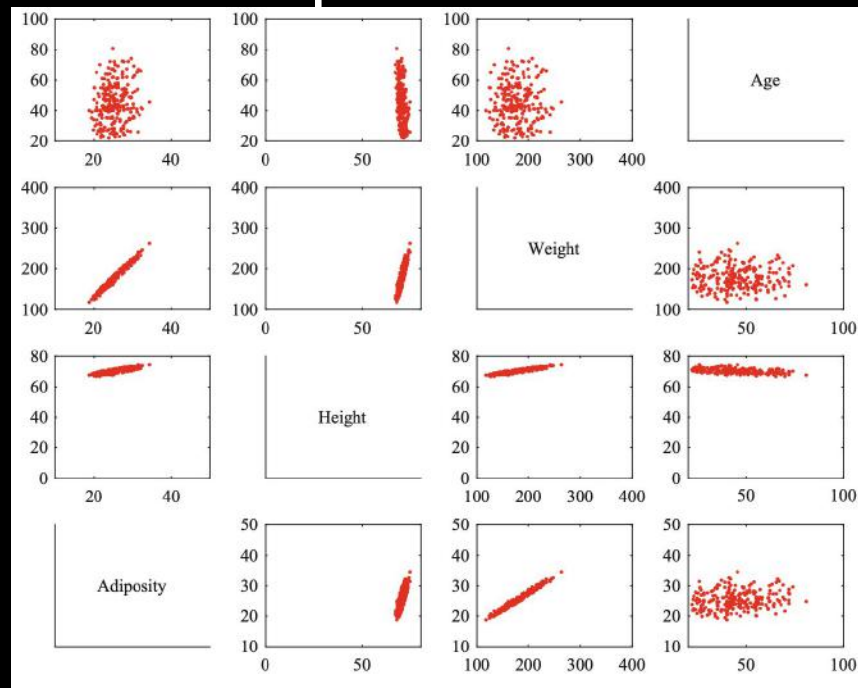
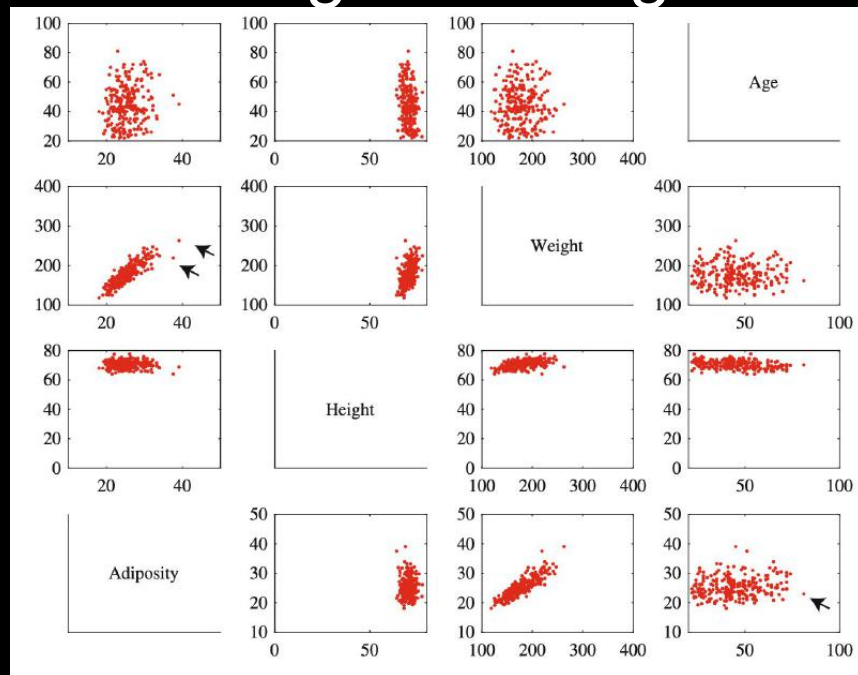
For any $m \times p$ matrix \mathcal{X} , it is possible to obtain a decomposition

$$\mathcal{X} = \mathcal{U}\Sigma\mathcal{V}^T$$

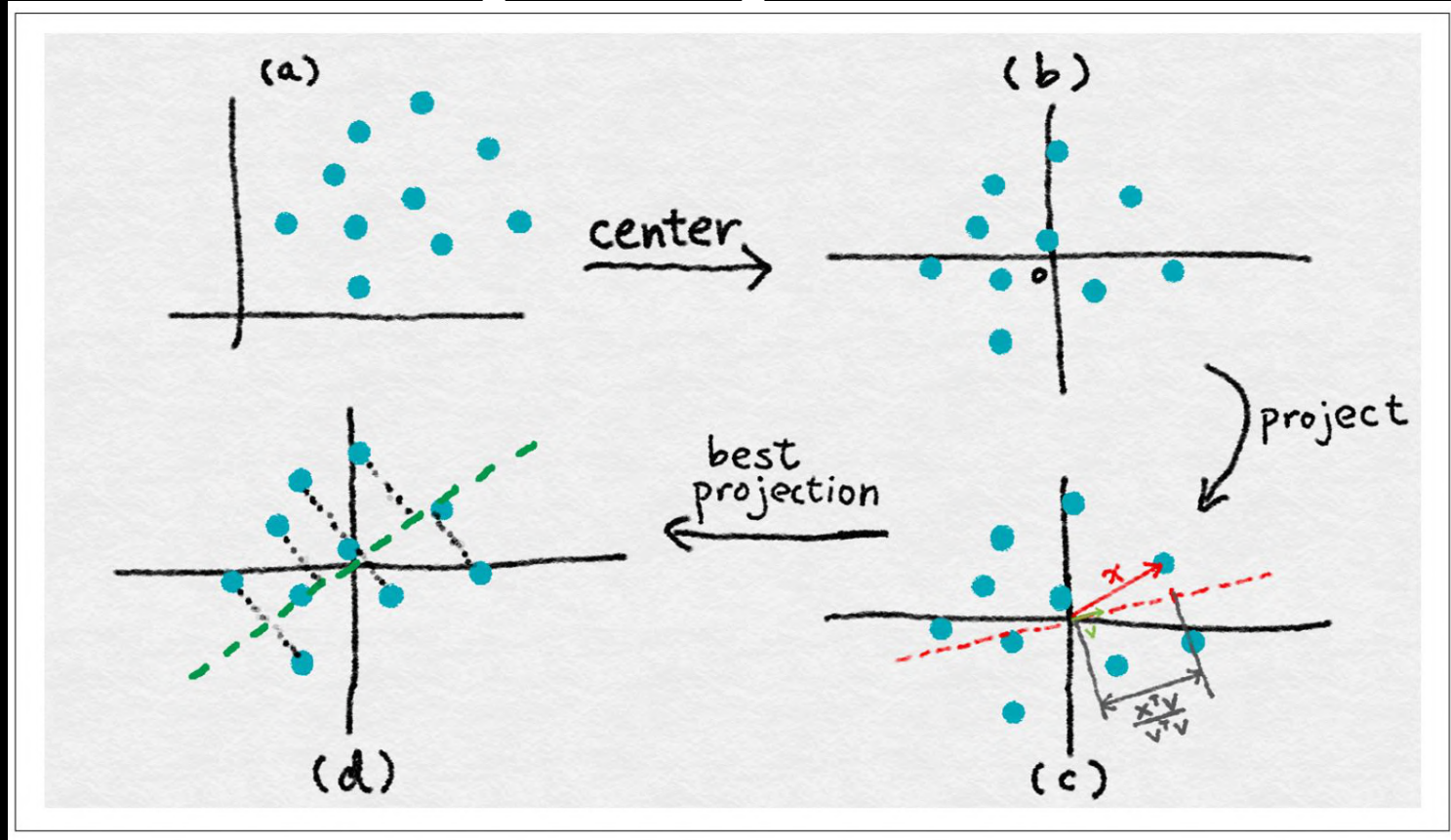
Principal Components are abstract representations



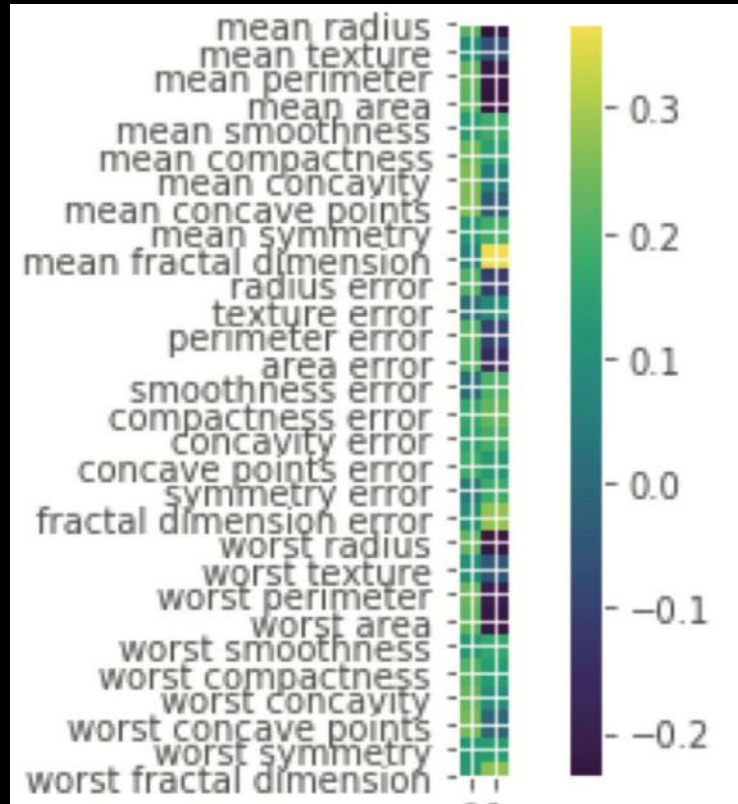
Retrieving back original data from representations



PCA for feature engineering



PCA for feature selection



Example: Approximating Faces

Mean image from Japanese Facial Expression dataset



First sixteen principal components of the Japanese Facial Expression dataset





Sample Face Image

mean

1

5

10

20

50

100



Approaches for High Dimensional Data

Low rank approximations

- Principal Components Analysis
- Multidimensional scaling
- Factor analysis
- Principal Coordinate Analysis
- Latent Semantic Analysis
- TF-IDF

Machine Learning Approaches

- Filter Based Feature Selection
- Wrapper Based Feature Selection
- Embedded approach
- Bayesian Networks