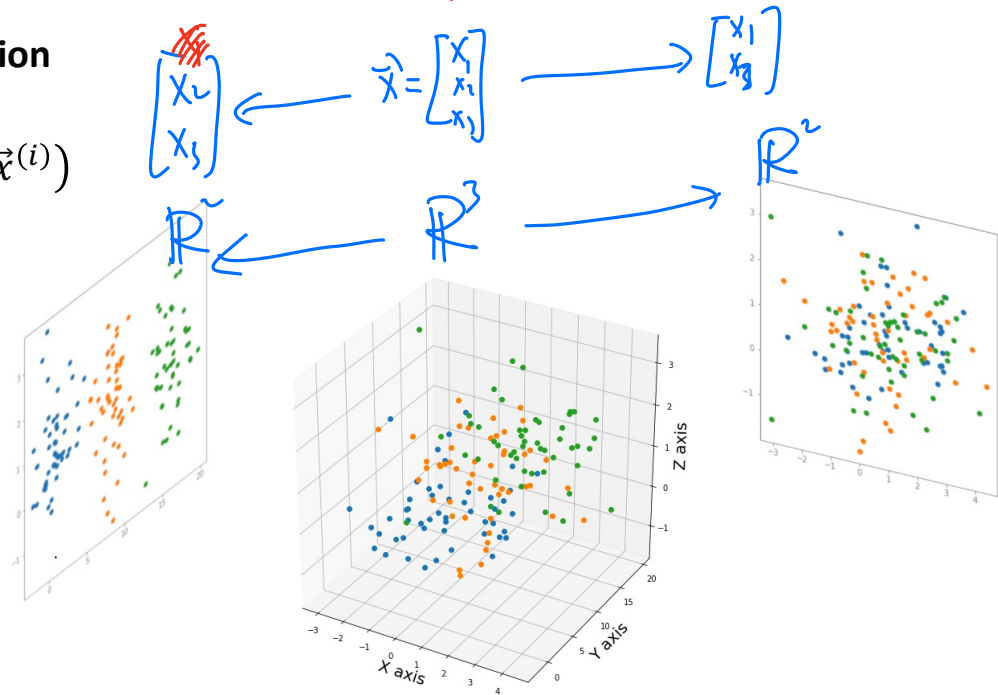**Section 19. Principal Components Analysis**

- Dimension Reduction
- Singular value decomposition
- Principal Components
- PCA regression
- PCA classification
- Non-linear PCA

➢ **Dimensional Reduction**

**Data:** $\left(\vec{x}^{(i)}, y^{(i)}\right)$ or $\left(\vec{x}^{(i)}\right)$

$$\vec{x} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} \longrightarrow \begin{bmatrix} x_1 \\ x_3 \end{bmatrix}$$

$$\begin{bmatrix} x_2 \\ x_3 \end{bmatrix} \longleftarrow$$

$\tilde{R} \longleftarrow R^3 \longrightarrow R^2$



**Dimensional Reduction** is the process of combining **high** dimensional information into **low** dimensional representations. Low dimensional representations can be computationally easier to work with, while hopefully throwing away only extraneous information.

Low dimensional projections are meaningful. In fact, all of data visualization can be said to be the result careful low dimensional projections. One goal is to discover the small set of variables that controls the larger phenomena.
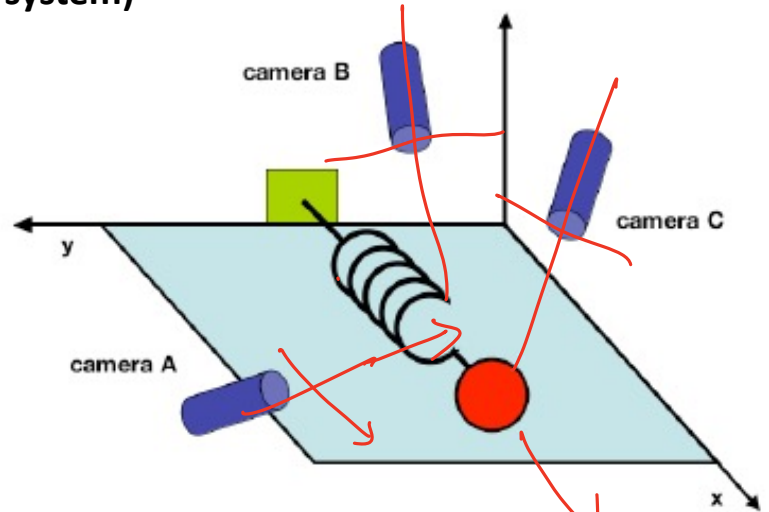
## ➢ Feature engineering

Constructed features play an important role in terms of subset selection. Using our naive subset selection for a linear model, we would quickly throw away some features. However, when we pick the right **combination of the features** suddenly we're given one our best predictors.

The process of **constructing meaningful features** is called **feature engineering**. Feature engineering is useful both in terms of constructing better models, and communicating results in a human readable format. One of the project of statistics, and by extension machine learning, is to construct a minimal set of mean rich features that explain all of the **variance** in the data.

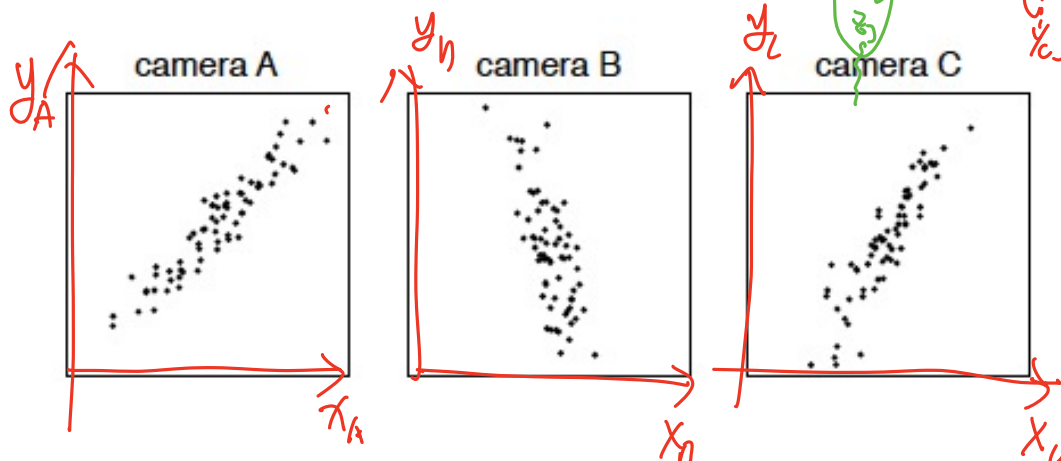To this end, feature engineering often uses the techniques of dimensional reduction, and visa versa.

## Another motivation example (mass–spring system)

Suppose we do not know what are the real x, y and z axes, so we choose three camera positions A, B, C at some arbitrary angles with respect to the system. The angles between our measurements might not even be $90^o$. Now, we record with the cameras for several minutes. Each camera produces a two-dimensional representation of the data.
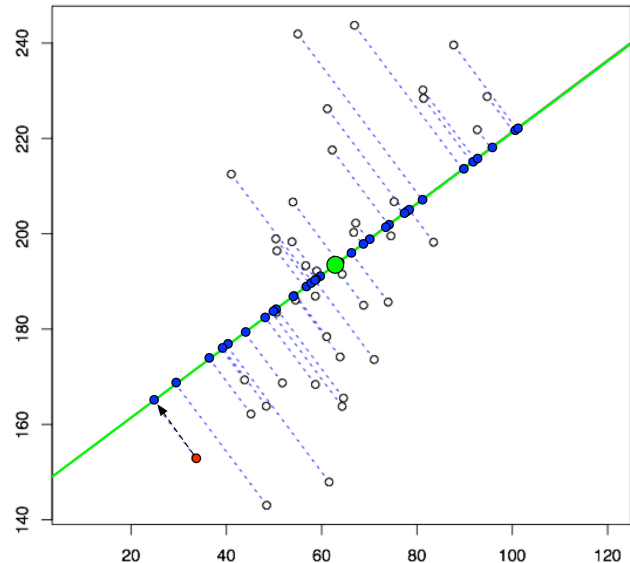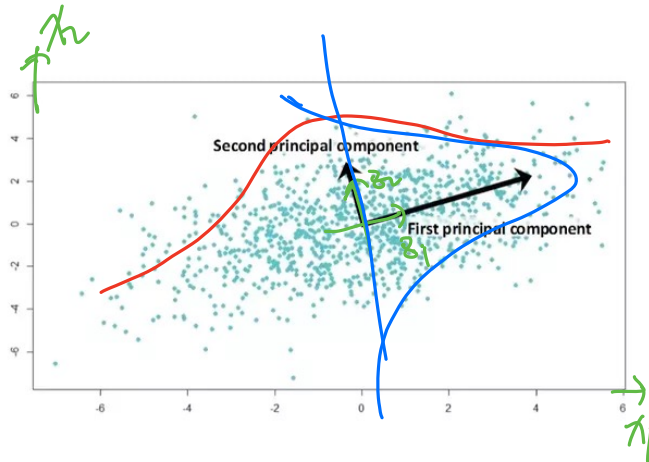
$t=0 \quad t=0.1 \quad t=0.2 \cdots$

$\vec{x}^{(0)}, \quad \vec{x}^{(1)}, \quad \vec{x}^{(2)}$

$$\vec{X} = \begin{bmatrix} x_A \\ y_A \\ x_B \\ y_B \\ x_C \\ y_C \end{bmatrix} \in \mathbb{R}^6$$



$\vec{x} = C \begin{bmatrix} \cdots \end{bmatrix} \begin{bmatrix} x_A \\ y_A \\ \vdots \\ y_C \end{bmatrix}$
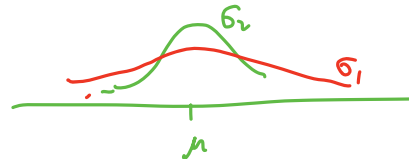
camera A

camera B

camera C

## ➢ Principal Components Analysis



There are two main techniques in dimensional reduction: **projection onto linear subspaces and manifold learning**. In projection onto linear subspaces, we try to discover the **linear combination of features** the provide the clearest coordinate system for the dataset, i.e., the **component factors**.

In an alternative view, factor analysis tries to fit a **linear subspace** to a data set in such a way as to **maximize the projected variance**, that is maximize the amount of data variance captured in orthogonal project to the new space.

$$\sigma_2 < \sigma_1$$

$$\mu = E(x)$$

- **Linearly combine** the variables to create the new variables, called **principal components.**

$$X : r.v.$$

- The first few explain most of the **variation**.

$$Var(X) = E\left[(X - \mu)^2\right]$$

- Achieve data reduction, without much loss of information.

- It's important to note though that PCA doesn't know anything about the labels $y$. This can be an **advantage**. Unlabeled data can import the PCA (an unlabeled data is often cheap).

PCA is a greedy algorithm with a beautiful mathematical interpretation.

The idea is to proceed iteratively:

- Find principal component (direction) that accounts for the largest possible variance.
- Project onto the subspace orthogonal to that vector.
- Repeat, collecting the orthogonal vectors into the rows of U until U is a $p \times p$ orthogonal matrix.

$$A \in \mathbb{R}^{n \times n}$$

$\star$ • A symmetric $\overset{\text{def}}{\Longleftrightarrow}$ $A = A^T$

Spectral Decomposition

$$\Longleftrightarrow A = PDP^{-1} \text{ and } P^T = P^{-1}$$

$$D = \begin{bmatrix} \lambda_1 & & & \\ & \lambda_2 & & \\ & & \ddots & \\ & & & \lambda_n \end{bmatrix}$$

$\boxed{P \text{ is orthogonal}}$

$E_{\lambda_1} \mid E_{\lambda_2} \cdots \mid E_{\lambda_s}$

• Ex: $f(x_1, x_2) = x_1^2 + 3x_2^2 + 4x_1 x_2 = \begin{bmatrix} x_1 & x_2 \end{bmatrix} \begin{bmatrix} 1 & 2 \\ 2 & 3 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$

• quadratic form $\boxed{P(\vec{x}) = \vec{x}^T A \vec{x}}$     $A = A^T$.

• $\underline{P(\vec{x}) \text{ or } A}$ is $\underline{\text{positive (semi)definite}}$ $\overset{\text{def}}{\Longleftrightarrow}$

$$P(\vec{x}) = \vec{x}^T A \vec{x} \geq 0 \text{ for any } \vec{x} \neq \vec{0} \in \mathbb{R}^n$$

$$\Longleftrightarrow \left( \text{all eigenvalues of } A \right) \geq 0$$

$\star$ • Let $X$ be any $m \times n$ matrix, then

  • $X^T X$ is positive semi-definite.
  
    $n \times n$

  • If rank $X = n$, then $X^T X$ is positive definite.

$A = a_1, a_2, \cdots, a_n$

➢ **Singular Value Decomposition**

$$\text{Data Matrix } X = \begin{bmatrix} x^{(1)T} \\ \vdots \end{bmatrix}_{n \times d}$$

Suppose $E(X) = \vec{0}$

$\cdot$ **Def :** Sample Covariance

$$\text{Var}(A) := \frac{1}{n-1} \sum_{i=1}^{n} (a_i - \bar{a})^2$$

$$\text{Cov}(X) := \left(\frac{1}{n-1}\right) X^T X = \begin{bmatrix} \text{Var}(X_1) & \text{Cov}(X_1, X_2) & \cdots \\ \vdots \\ \vdots \end{bmatrix}_{d \times d}$$

ages  height.

$$\cdot \text{Ex: } X = \begin{bmatrix} 3 & 7 \\ -4 & -6 \\ 7 & 8 \\ -1 & -1 \\ -4 & -1 \\ -3 & -7 \end{bmatrix}$$

$$\text{Cov}(X) = \frac{1}{5} X^T X = \begin{bmatrix} 20 & 25 \\ 25 & 46 \end{bmatrix}$$

$$\cdot \boxed{A = X^T X}$$
$d \times d.$

$$r = \text{rank } A = \text{rank}(X)$$

① A is symmetric $\implies$ A has $\underline{d \text{ real eigenvalues}}$

① A is $\underline{\text{positive-semidefinite}} \implies \lambda_1 \geqslant \lambda_2 \geqslant \lambda_3 \cdots \geqslant \lambda_r > 0$

$$X^T X = A = VDV^{-1} = \boxed{VDV^T}$$

$$D = \begin{bmatrix} \lambda_1 & & & & & 0 \\ & \lambda_2 & & & & \\ & & \ddots & \lambda_r & & \\ & & & & 0 & \\ 0 & & & & & 0 \end{bmatrix}$$

$$= \lambda_1 \vec{v_1} \vec{v_1}^T + \lambda_2 \vec{v_2} \vec{v_2}^T + \cdots + \lambda_r \vec{v_r} \vec{v_r}^T$$

$$V = \begin{bmatrix} \vec{v_1} & \cdots & \vec{v_d} \end{bmatrix}$$

$$U = [u_1, \cdots, u_r, \cdots, u_n]$$

$\underline{\text{Goal 2:}}$

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

$$= \sigma_1 \vec{u_1} \vec{v_1}^T + \cdots + \sigma_r \vec{u_r} \vec{v_r}^T$$

E.g. $A^T A =$
$$\begin{bmatrix} \vec{a_1}^T \\ \vec{a_v}^T \end{bmatrix} \begin{bmatrix} \vec{a_1} & \vec{a_v} \end{bmatrix}$$
$$= \begin{bmatrix} a_1 \cdot a_1 & a_1 \cdot a_v \\ a_v \cdot a_1 & a_v \cdot a_v \end{bmatrix}$$

$$\downarrow \quad X^TX = VDV^T \iff V^TX^TXV = D \iff (XV)^TXV = D$$

- $XV = \begin{bmatrix} X\vec{v_1} & \cdots & X\vec{v_d} \end{bmatrix} \longrightarrow$

$$\iff \begin{bmatrix} \cdots & X\vec{v_i} \cdot X\vec{v_j} & \cdots \\ & \cdots & \end{bmatrix} = \begin{bmatrix} \lambda_1 \lambda_2 & & 0 \\ & \cdots & \lambda_r \\ 0 & & \cdots & 0 \end{bmatrix}$$

$$\iff \begin{cases} \|X\vec{v_i}\|^2 = \lambda_i & \text{for } i = 1, 2, \cdots r \\ \\ X\vec{v_i} \cdot X\vec{v_j} = 0 & \text{for } i \neq j \iff \text{all orthogonal} \\ \\ X\vec{v_i} = 0 & \text{for } i > r \end{cases}$$

<span style="color:red">Singular values</span>

· Denote: $\vec{u_i} := X\vec{v_i}$ for $i=1,\dots,r$.    Denote: $\sigma_i = \|X\vec{v_i}\|$

$$X V = \underbrace{[\underbrace{X\vec{v_1}}_{} \cdots X\vec{v_r} \; \underbrace{X\vec{v_{r+1}}}_{} \cdots X\vec{v_d}]}_{} = [\sigma_1\vec{u_1} \; \sigma_2\vec{u_2} \cdots \sigma_r\vec{u_r} \; 0 \; 0 \cdots 0] \quad = \sqrt{\lambda_i}$$

$\|X\vec{v_i}\|$

$n\times d$  $d\times d$   $n\times n$  $n\times d$

$$\boxed{X V = U \Sigma}$$

orthogonal

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

$$= [\vec{u_1} \; \vec{u_2} \cdots \vec{u_r} \; ?\,?\cdots?] \begin{bmatrix} \sigma_1 \sigma_2 \cdots \sigma_r \; 0 \cdots 0 \\ & & O \\ O & & \\ 0 \cdots & 0 \end{bmatrix}$$

orthogonal complement

$n\times d$



$V^T = V^{-1}$

$\Sigma \quad \begin{bmatrix} \sigma_1 & 0 \\ 0 & \sigma_2 \end{bmatrix}$

$U$

$v_2$  $v_1$

$e_2$  $e_1$

$\sigma_2 e_2$  $\sigma_1 e_1$

$\sigma_2 u_2$  $\sigma_1 u_1$

$V$

## ➢ Principal Components

$$\sigma_i = \sqrt{\lambda_i}$$

Data $X_{n \times d}$

$$X = U \Sigma V^T = \sigma_1 \vec{u}_1 \vec{v}_1^T + \cdots + \sigma_r \vec{u}_r \vec{v}_r^T$$

Def: $\vec{v}_1 \cdots \vec{v}_r$ are called the principal components of $X$.

flip

Def New Data Matrix

$$X^T X = V D V^T$$

$$Z = X V$$

$$\quad n \times d \qquad n \times d \qquad d \times d$$

$$\text{COV}(Z) = \frac{1}{n-1} Z^T Z = \frac{1}{n-1} V^T X^T X V$$

$$= \frac{1}{n-1} V^T V D V^T V$$

$$= \frac{1}{n-1} \begin{bmatrix} \lambda_1 & & & & 0 \\ & \lambda_2 & \ddots & & \\ & & & \lambda_{r0} & \\ 0 & & & & 0 \end{bmatrix}$$

$$\begin{bmatrix} z_1 \\ z_2 \\ \vdots \\ z_d \end{bmatrix} = \vec{z} = V^T \vec{x} = \begin{bmatrix} \vec{v_1}^T \\ \vec{v_2}^T \\ \vdots \\ \vec{v_d}^T \end{bmatrix} \vec{x} = \begin{bmatrix} \vec{v_1}^T \vec{x} \\ \vdots \\ \vec{v_d}^T \vec{x} \end{bmatrix}$$

• The new features $z_i, z_j$ has correlation 0.

$$X = VDV^T$$

$$D = \begin{bmatrix} 100 & & & & \\ & 90 & & & \\ & & 10 & 1 & \\ & & & 0 \cdots 0 \end{bmatrix}$$

$$\boxed{Z = XV}$$

$$Z^T Z = \frac{1}{n-1} \begin{bmatrix} 100 & & & & 0 \\ & 90 & & & \\ & & 1 \cdot 1 & & \\ 0 & & & k \cdots 0 \end{bmatrix}$$

$$\tilde{Z}^T \tilde{Z} = \frac{1}{n-1} \begin{bmatrix} 100 & & \\ & 90 & 1000 \\ & & \end{bmatrix}$$

It is often desired to reduce the number of variables, especially when the number of variables in concern are too many. But the reduction must be done without much loss of information.

Principal components provide an ideal way of such reduction. One may retain the first $k$ principal components , which altogether explains of the total variation.

$$\frac{\lambda_1 + \cdots + \lambda_k}{\lambda_1 + \cdots + \lambda_p}$$

## Principal Components Regression

*(handwritten:)* Old data $X$ (nxd)    New data $Z = XV = \begin{bmatrix} | & | & | \\ & & \\ | & | & | \end{bmatrix}$ (nxd)

1. perform principal components analysis (PCA) on the original data,
2. perform dimension reduction by selecting the number of principal components ($m$) using cross-validation or test set error,
3. Conduct regression using the first $m$ dimension reduced principal components.

Principal components regression forms the derived input columns $z_m = X v_m$ and then regresses *on* $z_1, z_2, \cdots, z_m$ for some $m \leq p$.

Principal components regression discards the $p - m$ smallest eigenvalue components.

## Multiple linear regression



$X$   ($d$, $N$)   $\xrightarrow{\text{MLR}}$   y ($1$, $N$)

## Principal component regression



$X$   ($d$, $N$)   $\xrightarrow{\text{PCA}}$   $z_1, z_2, z_3$ ($d$)   $\xrightarrow{\text{MLR}}$   y ($1$, $N$)

PCA v.s. least square regression

## ➢ Principal Components for classification

Take as an example projection onto the first two principal components in the Iris data set.



**Iris Versicolor**  **Iris Setosa**  **Iris Virginica**

| Sepal Length | Sepal Width | Petal Length | Petal Width | Species |
|---|---|---|---|---|
| 5.1 | 3.5 | 1.4 | 0.2 | Setosa |
| 7.0 | 3.2 | 4.7 | 1.4 | Versicolor |

**Iris Data (red=setosa,green=versicolor,blue=virginica)**

There are four variables, with correlations given above.

PCA of IRIS dataset

Using PCA to project onto the first two principal components yields a projection onto

$$\vec{v_1} = w(1) = (0.36, -0.08, 0.86, 0.36) \; ; \; w(2) = (0.66, 0.73, -0.18, -0.07)$$

$$V\vec{v} =$$

We see the species are separable knowing only their features, and furthermore we have concrete measurement ratios to determine species.

Compare for a moment with the LDA projection. LDA requires knowing the labels and requires significantly more computational time for a similar fit.



LDA of IRIS dataset

**PCA MNIST**



The MNIST data set may provide one of the most startling examples. A simple PCA can be computed quite quickly and yields a wealth of information about the structure of the data set. Projecting onto the first two components shows a lot of structure on the MNIST dataset.

**Key assumption:** a small number of principal components suffice to explain most of the variability in the data, as well as the relationship with the response.

Choose $z_1, \ldots, z_M$ as the first M principal components.

This assumption may not hold!

For example, performing PCA on the fashion MNIST data set doesn't yield as strong of forms as on the digits MNIST.



$$\text{Data matrix} \rightarrow (X^T X) = V D V^T$$

"PCA"

$$\text{New Data} \quad Z = X V \quad \text{"RISE"}$$
$$\qquad\qquad\quad n \times d \quad n \times d \; d \times d$$

Example.



Projection of Admissions Data

In the context of regression, PCA can perform subset selection for us. In the graph above, we used PCA to project the features of the dataset onto the first two principal components.
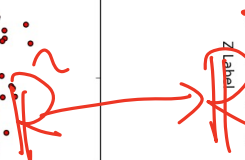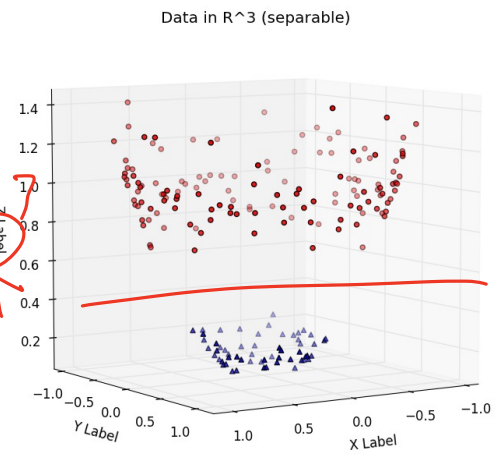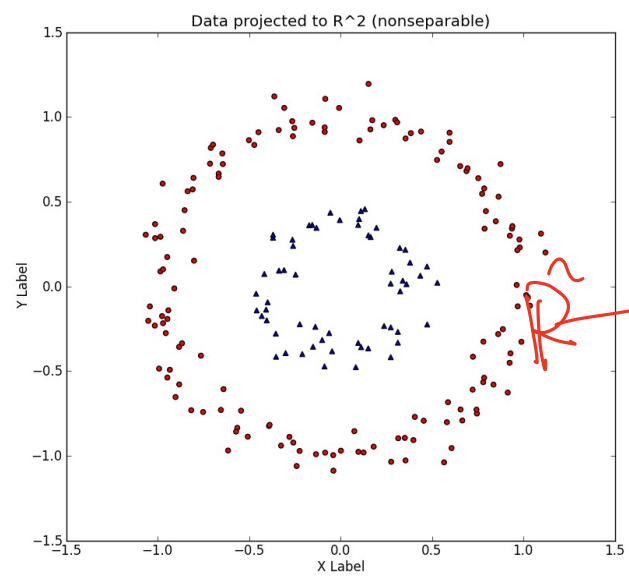
Of course, linear methods may be almost useless if the underlying structure isn't linear. For truly unknown structure we must construct smoothing maps using splines, random graphs, or other high complexity techniques. For example, in the following distribution, linear PCA will do nothing to reduce the complexity of the data set.
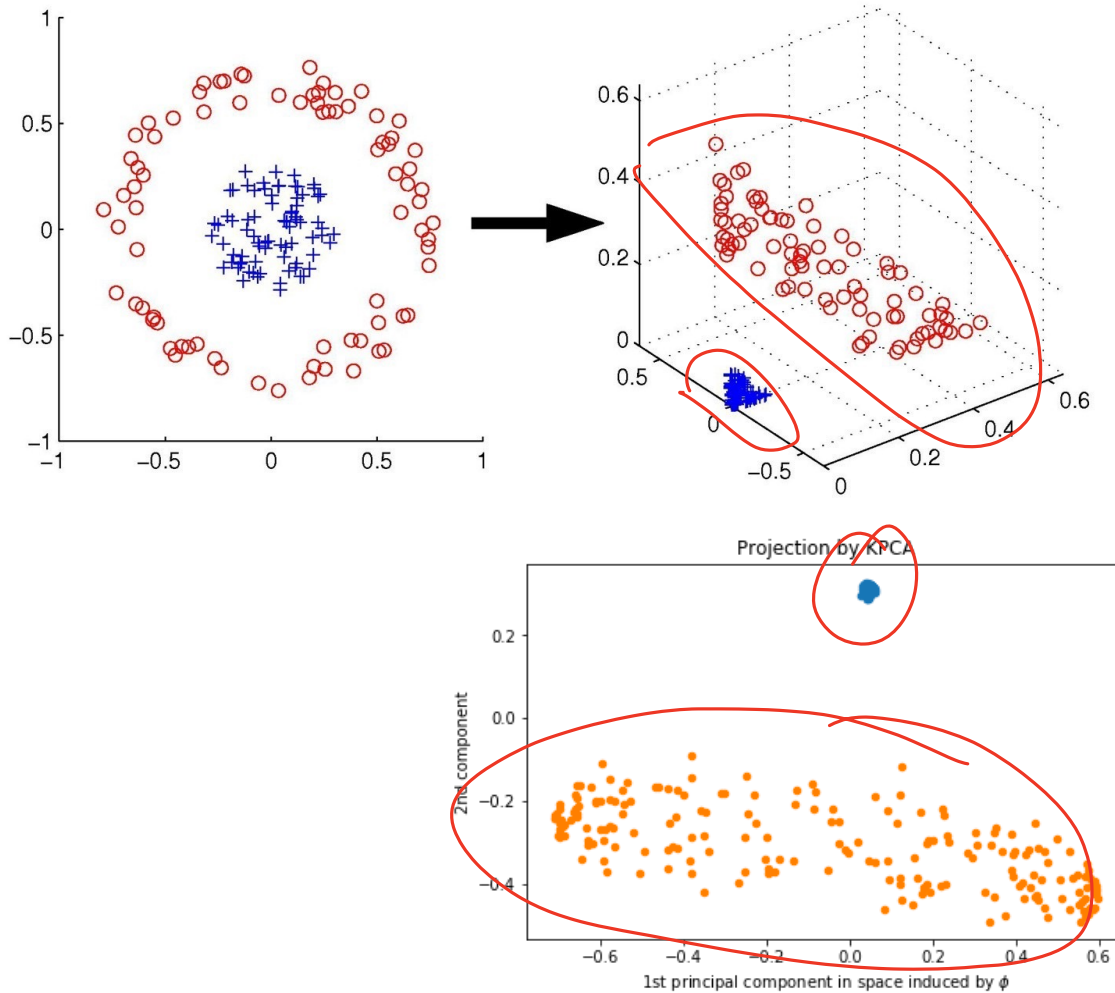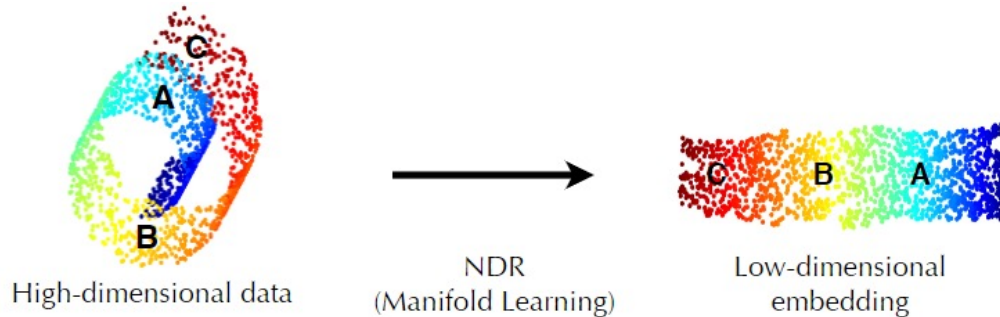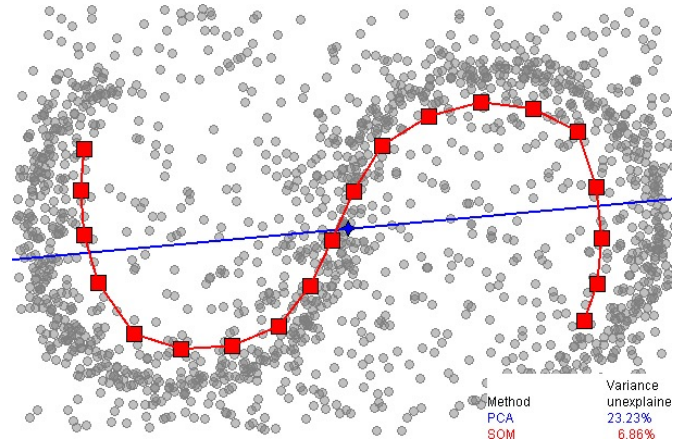


Original space

However, whenever we have linear methods there is the hope that we can apply the linear methods to generated features and get a a better fit. For nonlinear PCA, this means that we generate a high dimensional feature space, say $(1, \vec{x}, \vec{x}^2)$, and then project down onto a low dimensional subspace of this feature set, that is a normalized polynomial feature.



Data projected to R^2 (nonseparable)

Data in R^3 (separable)

Applying the **radial basis function** (rbf) kernel PCA to the concentric circle data set results in a dramatic clustering.



Projection by KPCA

➢ **Manifold learning**



In manifold learning, we try to fit a more complicated manifold to the underlying data. Manifold learning is of course much more complicated than simple linear dimensional reduction, but there do exist good algorithms.

Scikit-Learn:

https://scikit-learn.org/stable/modules/manifold.html

https://scikit-learn.org/stable/modules/decomposition.html#pca