# PCA

Thursday, April 25, 2024    8:13 PM
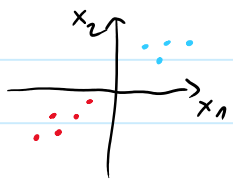
PCA - Principal Component Analysis

If input has too many dimensions:
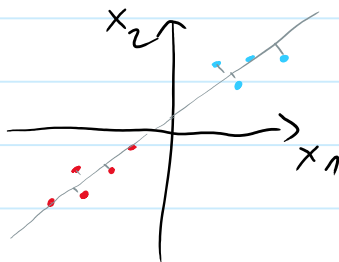- need more data to train
- more computationally difficult
- some ML Algos don't work well in high dimensions

---

Intuitive solution:

Assume we have data for a classification problem



We can project on a line



Now our new data is only



We got rid of 1 dimension.

Formally:

Instead of inputs $(x_1, ..., x_D)$, use $(y_1, ..., y_n)$, where $y_i = \sum a_j x_j + b_i$

(new input feature is just a linear combination of the original ones.
(linear with bias).

(linear with bias).

$y_i$'s are called principal components

$y_1$ such that projecting on $y_1$ gives maximal variance.

To get $y_2$, project onto $\langle y_1 \rangle^\perp$ and again, take $y_2$ in $\langle y_1 \rangle^\perp$ that maximizes variance.

## Algorithm

__Step 1__ Standardizing

  transform each input $x \in \mathbb{R}^D$

  $x \rightsquigarrow \frac{x - \mu}{\sigma}$, where $\mu$ is the mean, $\sigma$ the standard deviation (componentwise)

__Step 2__    Let $A := \Big( Cov(x_i, x_j) \Big)_{i,j \in \{1,..,D\}^2}$

Where $Cov(x_i, x_j) = \frac{1}{N-1} \sum_{\ell=1}^{N} (x_i^\ell - \bar{x}_i)(x_j^\ell - \bar{x}_j)$.

One can show, that the first $k$ principal components are the $k$ eigenvectors corresponding to the $k$ largest eigenvalues of $A$ (without abs. val.).