

Principal Component Analysis (PCA)

- $\{\underline{x}_i\}_{i=1}^N$ feature vectors i.i.d. distribution $p_x(\underline{x})$
- \underline{x}_i realization of random vector \underline{X}_i
- \underline{X}_i has mean μ and covariance matrix C
- subtract out μ from all \underline{x}_i 's

Then representative feature vector \underline{x} is a realization of random vector \underline{X} with $E[\underline{X}] = 0$, $\text{Cov}(\underline{X}) = C$

Goal Find orthogonal matrix $W_{d \times d}$ (i.e., $W^T = W^{-1}$)

s.t. (i) $\tilde{\underline{X}} = W \underline{X}$ has uncorrelated features \tilde{X}_i

(ii) The first K features $\tilde{X}_1, \tilde{X}_2, \dots, \tilde{X}_K$

Capture as much variance in data as possible for each $K = 1, 2, \dots, d$.

$$\begin{aligned}\tilde{C} &= \text{Cov}(\tilde{\underline{X}}) = \text{Cov}(W \underline{X}) \stackrel{\text{zero-mean } \underline{X}}{=} E[(W \underline{X})(W \underline{X})^T] \\ &= W E[\underline{X} \underline{X}^T] W^T \\ &= W C W^T\end{aligned}$$

We want \tilde{C} to be diagonal with the diagonal elements (variances of \tilde{X}_i 's) in decreasing order

Solution Eigen-decomposition of C

$$C = U D U^T$$

with $D = \begin{bmatrix} \lambda_1 & & & \\ & \lambda_2 & & 0 \\ & & \ddots & \\ 0 & & & \lambda_d \end{bmatrix}$, $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d$

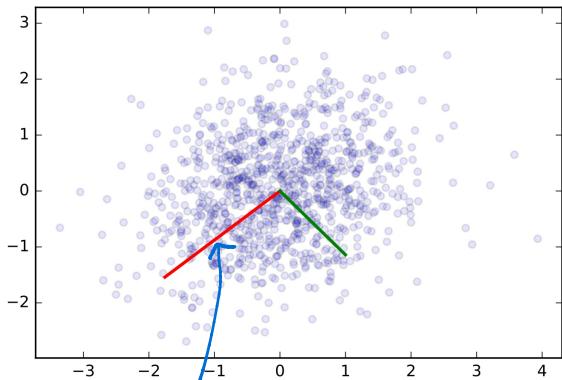
and $U = [\underline{u}_1 \ \underline{u}_2 \ \dots \ \underline{u}_d]$ s.t. $U^T U = I$

Now set $W = U^T$. Then,

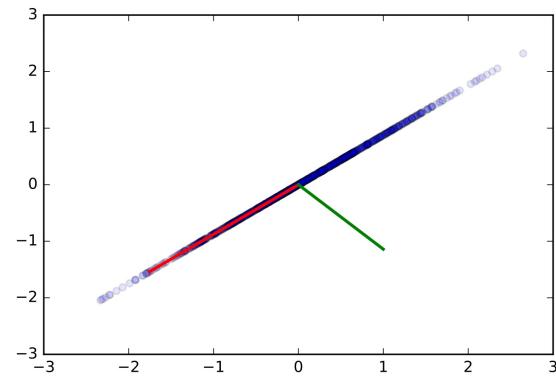
$$\begin{aligned}\tilde{C} &= \text{cov}(\tilde{\underline{X}}) = \text{cov}(U^T \underline{X}) = U^T C U \\ &= U^T U D U^T U = D\end{aligned}$$

$\Rightarrow \tilde{\underline{X}}$ has uncorrelated features with $\text{Var}(\tilde{x}_i) = \lambda_i$

Also, $\text{Var}(\tilde{x}_1) \geq \text{Var}(\tilde{x}_2) \geq \dots \geq \text{Var}(\tilde{x}_d) \geq 0$



$d = 2$
max. variance



Projection to 1-d

PCA can be used for dimensionality reduction by keeping only top k rows ($k < d$) of W

$$W = \begin{bmatrix} \underline{u}_1^T \\ \underline{u}_2^T \\ \vdots \\ \underline{u}_d^T \end{bmatrix}_{d \times d} \quad \xrightarrow{\hspace{1cm}} \quad W_k = \begin{bmatrix} \underline{u}_1^T \\ \vdots \\ \underline{u}_k^T \end{bmatrix}_{k \times d}$$

$$\tilde{\underline{X}}_{k \times 1} = W_k \underline{X}_{d \times 1}$$

To project back into original space

$$\underbrace{W_k^T \tilde{\underline{X}}}_{d \times 1} = W_k^T W_k \underline{X}_{d \times 1} \approx \underline{X}$$

In practice, C is not known a priori but is estimated from data.

data matrix $\rightarrow \mathbf{X}_{N \times d} = \begin{bmatrix} \underline{x}_1^T \\ \underline{x}_2^T \\ \vdots \\ \underline{x}_N^T \end{bmatrix}$ d typically < N

Estimate of Covariance:

$$\hat{C}_{d \times d} = \frac{1}{N} \sum_{i=1}^N \underline{x}_i \underline{x}_i^T = \frac{1}{N} \mathbf{X}^T \mathbf{X}$$

To find PCA transform W , we can compute eigen decomposition of \hat{C} as UDU^T and set $W = U^T$

But since \hat{C} is a scaled version of $\mathbf{X}^T \mathbf{X}$, we can also get W from SVD of \mathbf{X}

$$\mathbf{X} = U S V^T$$

$$\mathbf{X}^T \mathbf{X} = V S^2 V^T$$

$$\hat{C} = \frac{1}{N} \mathbf{X}^T \mathbf{X} = V \frac{S^2}{N} V^T$$

Set $W = V^T$

Variance (estimate) of i^{th} feature = $\frac{\sigma_i^2}{N}$.

Choosing Number of Principal Components $k < d$

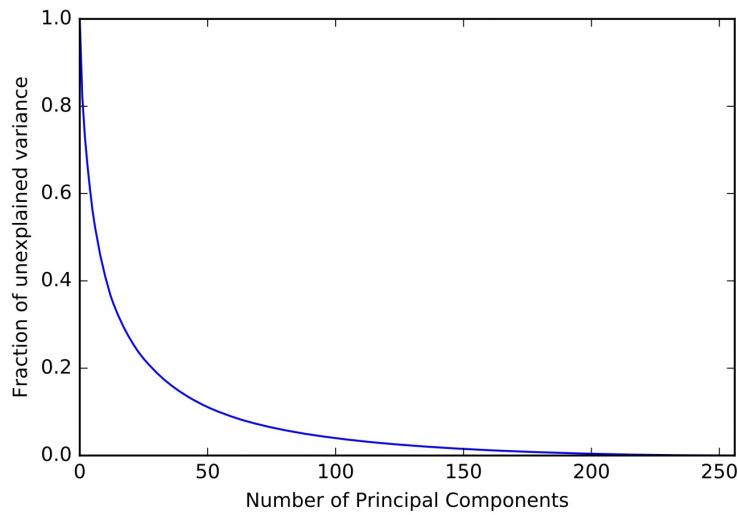
With dim. reduction from d to k ,

$$\underbrace{\underline{x}_i}_{d \times 1} \longrightarrow \underbrace{w_k \underline{x}_i}_{k \times 1}$$

Error in approximating $\underline{x}_i = \|\underline{x}_i - w_k^T w_k \underline{x}_i\|$

Cost function :

$$J(k) = \frac{1}{N} \sum_{i=1}^N \|\underline{x}_i - w_k^T w_k \underline{x}_i\|^2$$



"Scree" Plot

Look for knee in scree plot.