

Dimensionality Reduction :

Sid

Data Compression $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\} \quad x^{(i)} \in \mathbb{R}^n$

Data Visualization $\{z^{(1)}, z^{(2)}, \dots, z^{(m)}\} \quad z^{(i)} \in \mathbb{R}^k$

$$k \leq n \quad (2 \text{ or } 3)$$

Principal Component Analysis

- Reduce n-dim \rightarrow k-dim: Find k vectors $u^{(1)} \dots u^{(k)}$ to project data & minimize projection error

\rightarrow avg. of distances of features to projection plane/line

- \neq linear regression (squared error vs shortest distance)

PCA algorithm

- Data preprocessing (feature scaling / mean normalization)

$$\mu_j = \frac{1}{m} \sum_{i=1}^m x_j^{(i)} \Rightarrow \text{replace } x_j^{(i)} \text{ w/ } x_j^{(i)} - \mu_j$$

$$x_j^{(i)} \leftarrow \frac{x_j^{(i)} - \mu_j}{s_j - \text{mean}}$$

- Compute covariance matrix:

$$\Sigma = \frac{1}{m} \sum_{i=1}^m (x^{(i)})(x^{(i)})^T \quad \Sigma = \frac{1}{m} * X' * X ;$$

- Compute eigenvectors of matrix

$$[U, S, V] = \text{svd}(\Sigma)$$

$$U = \begin{bmatrix} | & | & & | \\ u^{(1)} & u^{(2)} & \dots & u^{(n)} \\ | & | & & | \end{bmatrix} \in \mathbb{R}^{n \times n}$$

$$x \in \mathbb{R}^n \rightarrow z \in \mathbb{R}^k \quad z = U_{\text{reduce}}' * x;$$

$$z^{(i)} = \underbrace{\begin{bmatrix} | & | & \dots & | \\ u^{(1)} & u^{(2)} & \dots & u^{(k)} \\ | & | & \dots & | \end{bmatrix}^T}_{\substack{u \times k \\ (U_{\text{reduce}})}} \cdot x^{(i)} = \underbrace{\begin{bmatrix} - (u^{(1)})^T - \\ \vdots \\ - (u^{(k)})^T - \end{bmatrix}}_{\substack{k \times n \\ (U_{\text{reduce}}')}} \underbrace{x^{(i)}}_{n \times 1}$$

```
Sigma = (1/m) * X' * X; % compute the covariance matrix
[U,S,V] = svd(Sigma); % compute our projected directions
Ureduce = U(:,1:k); % take the first k directions
Z = X * Ureduce; % compute the projected data points
```

Reconstruction from compressed representation

$$\underbrace{X_{\text{approx}}}_{\mathbb{R}^n} = \underbrace{U_{\text{reduce}}}_{n \times k} \cdot \underbrace{Z}_{k \times 1}$$

Choosing # of Principal components (k)

$$\text{Avg squared projection error} : \frac{1}{m} \sum_{i=1}^m \|x^{(i)} - x_{\text{approx}}^{(i)}\|^2$$

$$\text{Total variation in data} : \frac{1}{m} \sum_{i=1}^m \|x^{(i)}\|^2$$

Choose k such that:

$$\frac{\frac{1}{m} \sum_{i=1}^m \|x^{(i)} - x_{\text{approx}}^{(i)}\|^2}{\frac{1}{m} \sum_{i=1}^m \|x^{(i)}\|^2} \leq 0.01$$

(99% of variance retained)

$[u, s, v] = \text{svd}(\text{Sigma})$

$S = \begin{bmatrix} s_{11} & s_{12} & \dots & 0 \\ 0 & s_{22} & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & s_{nn} \end{bmatrix}$

$\left| - \frac{\sum_{i=1}^k s_{ii}}{\sum_{i=1}^n s_{ii}} \right|$