

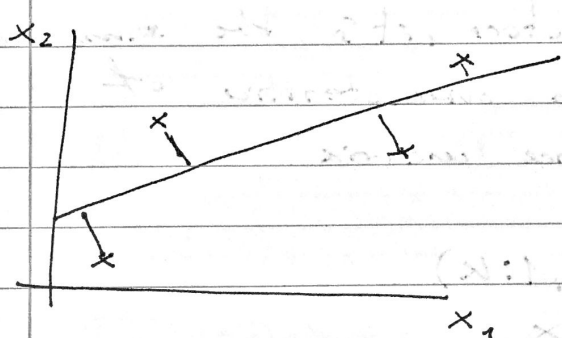
# Dimensionality reduction

- PCA  $\rightarrow$  principal component analysis

It's to make the common shortest projection in ~~an~~ a lower dimension to all your points

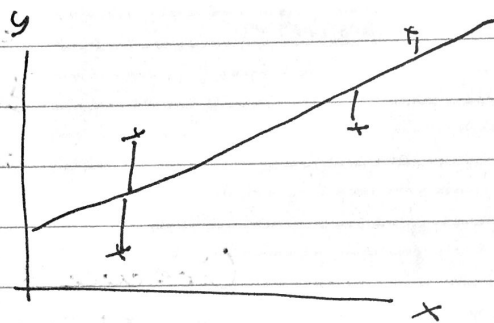
PCA

There's no label



Linear regression

There's a label



The distances to the projection are orthogonal

Parallel to the label

- Data preprocessing

Training set:  $x^1, x^2, \dots, x^m$

Preprocessing (feature scaling / mean normalization)

$$\mu_j = \frac{1}{m} \sum_{i=1}^m x_j^i \leftarrow \text{replace } x_j^i \text{ with } x_j^i - \mu_j$$

Scale features to have comparable range of values

# • PCA algorithm

Reduce data from  $n$ -dimensions to  $k$ -dimensions

- Covariance matrix

$$\Sigma = \frac{1}{m} \sum_{i=1}^m \underbrace{x^i}_{n \times 1} \underbrace{(x^i)^T}_{1 \times n} \rightarrow \bullet$$

- In Matlab  $\rightarrow [U, S, V] = \text{svd}(\Sigma)$   
 $\text{eig}(\Sigma)$   $\rightarrow$

eigenvectors, it's the same here  
 due to characteristics of  
 covariance matrix

$$U_{\text{reduce}} = U(:, 1:k)$$

$$z = U_{\text{reduce}}^T \cdot X$$

$$\Sigma = (1/m) \cdot X^T \cdot X$$

- Reconstruction for compressed representation

$$z = U_{\text{reduce}}^T \cdot X$$

$$\underbrace{\quad}_{n \times k} \quad \underbrace{\quad}_{k \times 1}$$

- Choosing  $k$  (no. of principal components)

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

total variation in  
 the data

avg squared projection  
 error

In Matlab,  $[U, S, V] = \text{svd}(E)$

$$S = \begin{bmatrix} s_{11} & & 0 \\ & s_{22} & \\ 0 & & s_{nn} \end{bmatrix}$$

$$\text{For given } k \quad 1 - \frac{\sum_{i=1}^k s_{ii}}{\sum_{i=1}^n s_{ii}} \leq 0.01 \Leftrightarrow \frac{\sum_{i=1}^k s_{ii}}{\sum_{i=1}^n s_{ii}} \geq 0.99$$

99% of variance  
retained

### • Applications of PCA

#### - Compression

- Reduce memory/disk needed to store data
- Speed up learning algorithm

#### - Visualization

- Don't use to prevent overfitting! For that you have other tools like regularization.