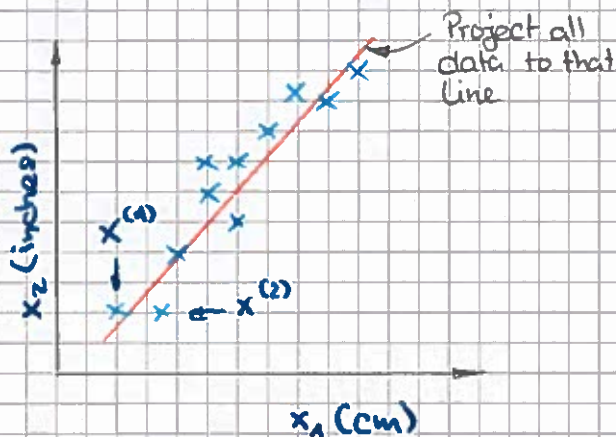


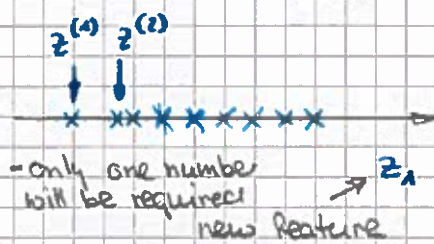
14. Dimensionality Reduction (PCA)

Data Compression:



Reduce Data from
2D to 1D

- Redundant data set - different units for same attributes
⇒ can happen when different teams are working independently

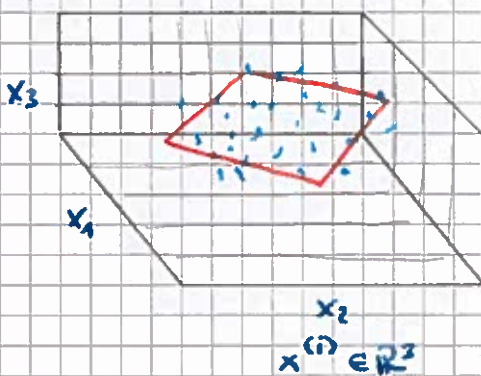


⇒ compression lead to faster algorithm

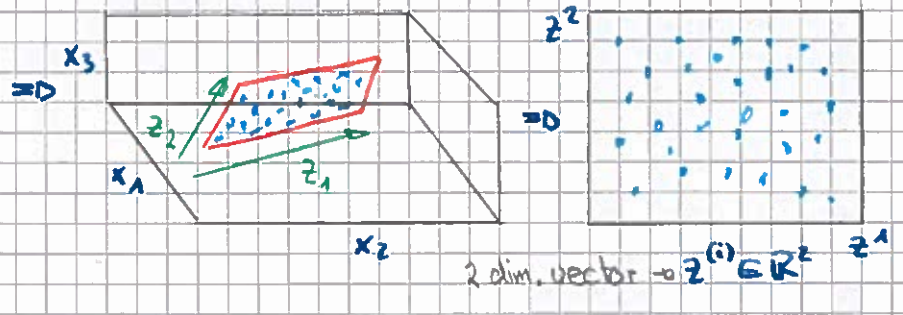
$$\begin{aligned} x^{(1)} &\in \mathbb{R}^2 \rightarrow z^{(1)} \in \mathbb{R} \\ x^{(2)} &\in \mathbb{R}^2 \rightarrow z^{(2)} \in \mathbb{R} \end{aligned}$$

⋮

$$x^{(m)} \in \mathbb{R}^2 \rightarrow z^{(m)} \in \mathbb{R}$$



Reduce 3D data to 2D



$$z = \begin{bmatrix} z_1 \\ z_2 \end{bmatrix} \quad z^{(i)} = \begin{bmatrix} z_1^{(i)} \\ z_2^{(i)} \end{bmatrix}$$

Dimensionality reduction of a dataset $(x^{(1)}, x^{(2)}, \dots, x^{(m)})$
where $x^{(i)} \in \mathbb{R}^n$

⇒ A lower dimensional dataset $(z^{(1)}, z^{(2)}, \dots, z^{(m)})$ of m examples where $z^{(i)} \in \mathbb{R}^k$ for some value of k and $k \leq n$

Data Visualization:

Country	x_1 GPD	x_2 per capita GPD	x_3 Human dev.	x_4 Life	x_5 pow. Index
Canada	1.577	39.17	0.908	80.7	32.6
China	:	:	:	:	:
India	:	:	:	:	:
Singapore	0.223	56.59	0.866	80	42.5
USA	14.527	46.86	0.91	78.3	40.8

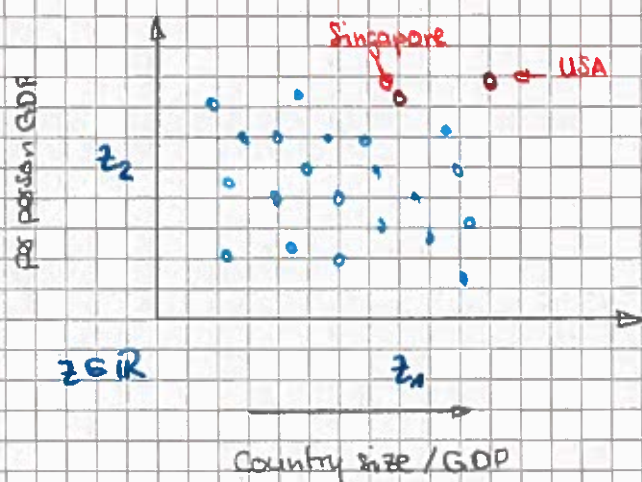
$$\Rightarrow x^{(i)} \in \mathbb{R}^{50}$$



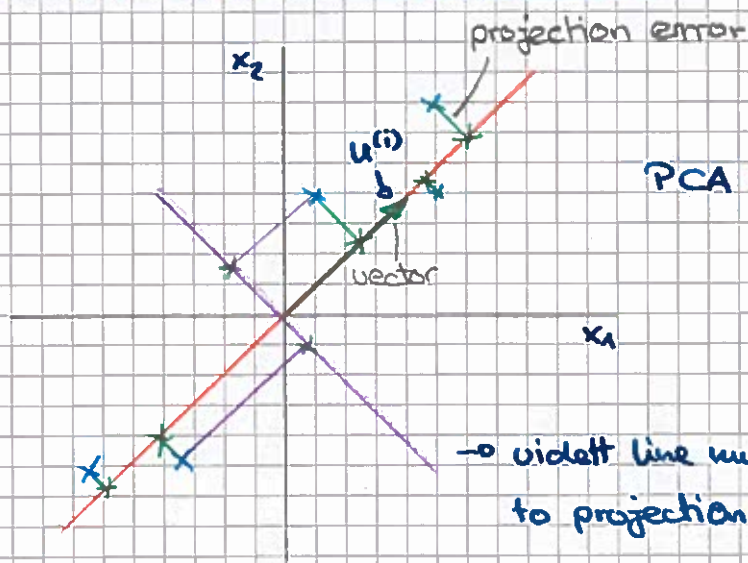
Country	z_1	z_2 ← no particular meaning
Canada	1.6	1.2
China	:	:
India	:	:
Singapore	0.5	1.7
USA	2	1.5

$$z^{(i)} \in \mathbb{R}^2$$

Reduce the data from 50D to 2D



$$k \leq n$$

Principal Component Analysis (PCA)

$$x \in \mathbb{R}^2$$

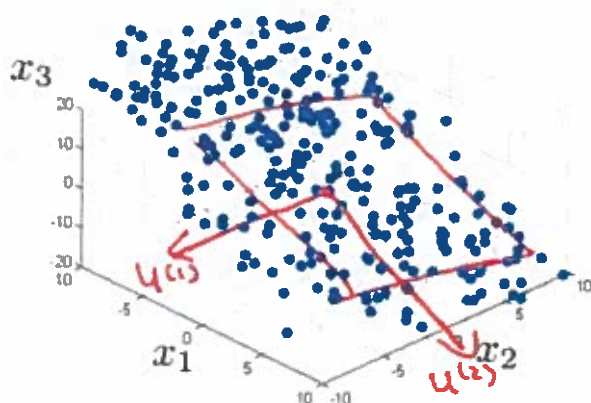
PCA tries to find lower dimensional projection \Rightarrow projection error to be min.

\Rightarrow violet line much worse than red line in respect to projection error.

Reduce from 2-dimension to 1-dimension: Find a direction (a vector $u^{(1)} \in \mathbb{R}^n$) onto which to project the data so as to minimize the projection error.

\Rightarrow it doesn't matter if vector is $u^{(1)}$ or $-u^{(1)}$

Reduce from n -dimensional to k -dimensional: Find vectors $u^{(1)}, u^{(2)}, \dots, u^{(k)}$ onto which to project the data, so to minimize the projection error.



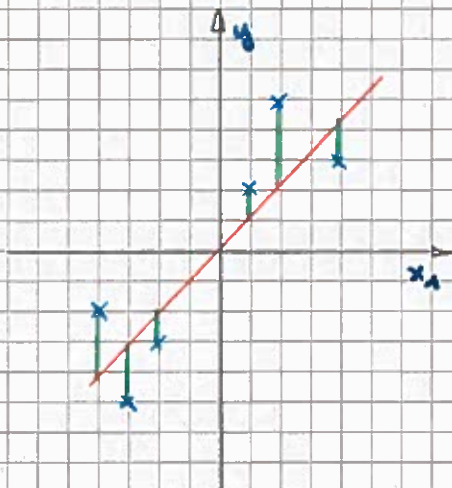
3D \rightarrow 2D

$k=2$

- tries to find a line to minimize the projection error

- with PCA all features are treated equally

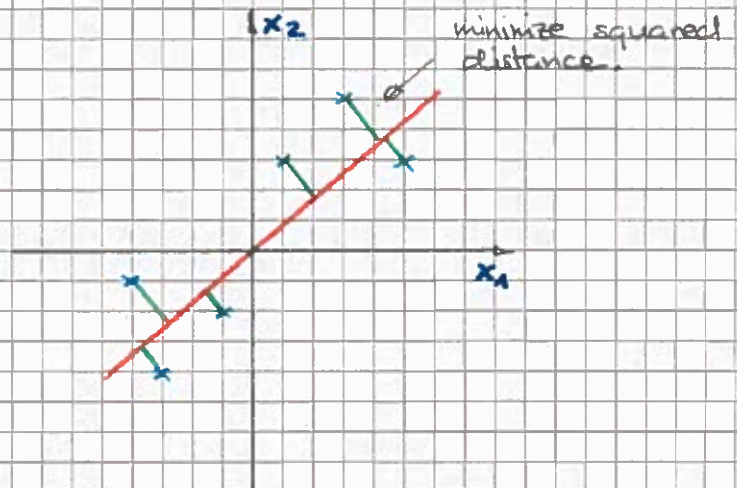
PCA is not Linear Regression



prediction of y based on
 x_1 value (Feature)

$$x \rightarrow y$$

both are very different algorithms.



\Rightarrow projection, each feature is
treated equally

$$x_1, x_2, \dots, x_n$$

Data PreprocessingTraining set: $x^{(1)}, x^{(2)}, \dots, x^{(n)}$

Processing (Feature scaling / mean normalization):

$$\mu_j = \frac{1}{n} \sum_{i=1}^n x_j^{(i)}$$

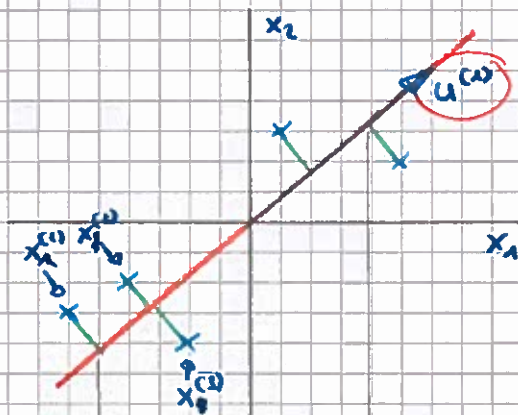
↳ for unlabeled data

Replace each $x_j^{(i)}$ with $x_j - \mu_j$ If different features on different scales (e.g. x_1 = size of house \Rightarrow x_2 = # of bedrooms), scale features to have comparable ranges of values.

$$x_j = \frac{x_j^{(i)} - \mu_j}{s_j}$$

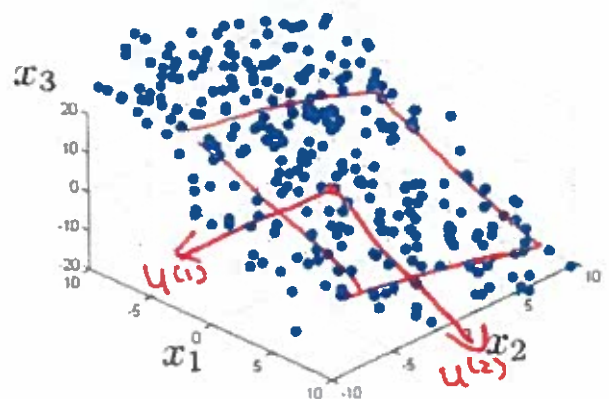
 \Rightarrow resp. standard dev.

or max - min Value

Principal Component Analysis Algorithm

Reduce data from 2D to 1D

$$x^{(i)} \in \mathbb{R}^2 \Rightarrow z^{(i)} \in \mathbb{R}$$



Reduce data from 3D to 2D

$$x^{(i)} \in \mathbb{R}^3 \Rightarrow z^{(i)} \in \mathbb{R}^2$$

$$z^{(i)} = \begin{bmatrix} z_1 \\ z_2 \end{bmatrix}$$

Reduce data from n -dimensions to k -dimensions
 Compute the "covariance matrix"

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

(Sigma)

Compute "eigenvectors" of matrix Σ :

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

\downarrow
 required matrix

\downarrow
 matrix
 $n \times n$ vector

singular value decomposition
 $\text{eig}(\text{Sigma}) \Rightarrow$ less stable as svd

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

$\underbrace{\hspace{10em}}_k$

$$x \in \mathbb{R}^n \rightarrow z \in \mathbb{R}^k$$

$$z^{(i)} = \begin{bmatrix} | & | & & | \\ u^{(1)} & u^{(2)} & \dots & u^{(k)} \\ | & | & & | \end{bmatrix}^T \cdot x^{(i)} = \begin{bmatrix} - (u^{(1)})^T - \\ \vdots \\ - (u^{(k)})^T - \end{bmatrix} \cdot x^{(i)}$$

$\underbrace{\hspace{10em}}_{n \times k}$ $\underbrace{\hspace{10em}}_{k \times n}$ $\underbrace{\hspace{10em}}_{n \times 1}$

$z \in \mathbb{R}^k$ Unreduced $k \times 1$

\Rightarrow After mean normalization (ensure every feature has zero mean) and optionally feature scaling.

$$\text{Sigma} = \frac{1}{m} \sum_{i=1}^m (x^{(i)}) (x^{(i)})^T \Rightarrow X = \begin{bmatrix} -x^{(1)T} - \\ \vdots \\ -x^{(m)T} - \end{bmatrix}$$

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

$$\text{Unreduce} = U(:, 1:k);$$

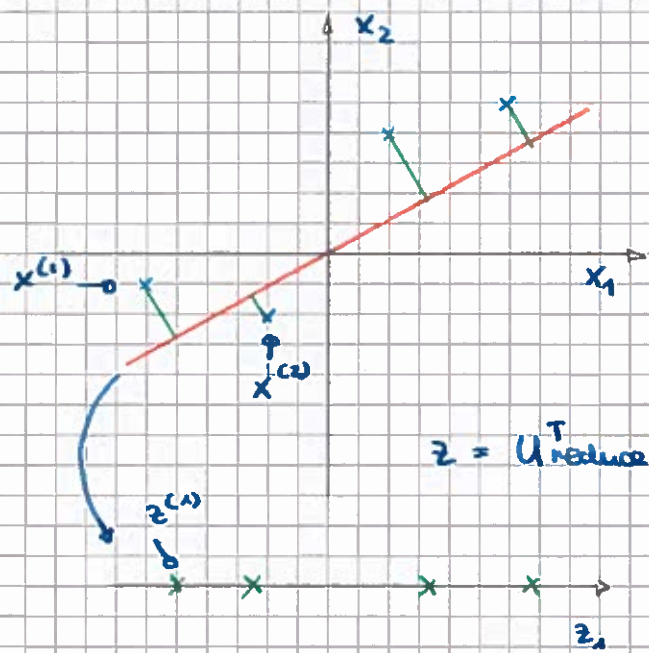
$$z = \text{Unreduce}^T \cdot x;$$

$$\hookrightarrow x \in \mathbb{R}^n \Rightarrow x_0 = 1$$

vectorized $\text{Sigma} = \left(\frac{1}{m}\right) X^T X;$

Reconstruction from Compressed Representation

7



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

$$x_{\text{approx}}^{(1)} \in \mathbb{R}^2$$

$$x_{\text{approx}}^{(2)}$$

$$z \in \mathbb{R} \Rightarrow x \in \mathbb{R}^2$$

$$x_{\text{approx}}^{(1)} = \underbrace{U_{\text{reduce}}}_{n \times k} \cdot \underbrace{z}_{k \times 1}$$

\Downarrow
 \mathbb{R}^n $n \times 1$

Choosing K (number of Principal Components)

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

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

Typically choose k to be smallest value so that

$$K \rightarrow \left\{ \frac{\frac{1}{m} \sum_{i=1}^m \|x^{(i)} - x_{\text{approx}}^{(i)}\|^2}{\frac{1}{m} \sum_{i=1}^m \|x^{(i)}\|^2} \leq 0.01 \text{ (1\%)} \right.$$

validation

=> "99% of variance retained"

=> 95% - 98% are usual

Algorithm:

Try PCA with $k=1$

Compute: $U_{\text{reduce}}, z^{(1)}, z^{(2)}, \dots, z^{(m)}, x_{\text{approx}}^{(1)}, \dots, x_{\text{approx}}^{(m)}$

then check if

=> Formula => K-Validation ≤ 0.01

if not start over again with $k=2, \dots$

↑

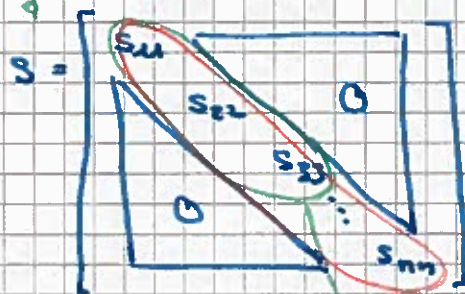
Approach inefficient!

Pick smallest value of k for which the result is:

"99% of variance retained"



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



For a given k

$$1 - \frac{\sum_{i=1}^k s_{ii}}{\sum_{i=1}^n s_{ii}} \Rightarrow \leq 0.01$$

K=3

equivalent to K-Validation Formula

$$\left\{ \frac{\sum_{i=1}^k s_{ii}}{\sum_{i=1}^n s_{ii}} \geq 0.99 \right\} \begin{array}{l} \text{calculate} \\ \text{this for different} \\ \text{k values} \\ \text{smallest value of k} \end{array}$$

svd(Sigma) must only be called once

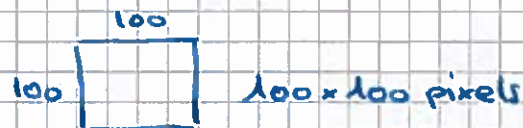
Supervised learning speedup

Very high Feature vector
assume computer vision problem

$$\Rightarrow (x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)}) \Rightarrow x^{(i)} \in \mathbb{R}^{10'000}$$

Extract inputs:

$$\text{Unlabeled data set: } x^{(1)}, x^{(2)}, \dots, x^{(m)} \in \mathbb{R}^{10'000}$$



PCA

$$z^{(1)}, z^{(2)}, \dots, z^{(m)} \in \mathbb{R}^{1000}$$

for prediction x
reduced by 10x

New training set:

$$(z^{(1)}, y^{(1)}), (z^{(2)}, y^{(2)}), \dots, (z^{(m)}, y^{(m)}) \Rightarrow h_{\theta}(z) = \frac{1}{1 + e^{-\theta^T z}}$$

Note: Mapping $x^{(i)} \rightarrow z^{(i)}$ should be defined by running PCA only on the training set. This mapping can be applied as well to the examples $x_{\text{cv}}^{(i)}$ and $x_{\text{test}}^{(i)}$, in the cross validation and test set.

\Rightarrow is calculating U_{reduce} \Rightarrow parameter learned by PCA as well as mean normalization

\Rightarrow PCA makes a mapping between $x \rightarrow z$

Application of PCA

- Compression

- Reduce memory/disk needed to store data

- Speed up learning algorithm

- Choose k by % of variance retained.

- Visualization

only $k=2$ or $k=3$ could be used to plot data.

Bad use of PCA: To prevent overfitting

10

Use $z^{(i)}$ instead $x^{(i)}$ to reduce the number of features to $k < n$
→ this makes it less likely that overfitting occurs.

⇒ BAD use!

This might work ok, but isn't a good way to address overfitting. Use regularization instead.

$$\Rightarrow \min_{\theta} \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2$$

⇒ PCA throws away ⇒ therefore might lose valuable information

PCA is sometime used where it shouldn't be

Design of ML system:

- Get training set: $(x^{(1)}, y^{(1)})$
 - Run PCA to reduce $x^{(i)}$ in dimension to get $z^{(i)}$
 - Train logistic regression on $(\underset{x^{(1)}}{z^{(1)}}, y^{(1)}), (\underset{x^{(2)}}{z^{(2)}}, y^{(2)}), \dots, (\underset{x^{(m)}}{z^{(m)}}, y^{(m)})$
 - Test on Test set: Map $x_{\text{test}}^{(i)}$ to $z_{\text{test}}^{(i)}$. Run $h_{\theta}(z)$ on $(\underset{z_{\text{test}}^{(1)}}{z_{\text{test}}^{(1)}}, \underset{y_{\text{test}}^{(1)}}{y_{\text{test}}^{(1)}}), \dots, (\underset{z_{\text{test}}^{(m)}}{z_{\text{test}}^{(m)}}, \underset{y_{\text{test}}^{(m)}}{y_{\text{test}}^{(m)}})$
- } Project plan

⇒ How about doing without PCA?

Before implementing PCA, first try running whatever you want to do by using the original/raw data $x^{(i)}$. Only if that doesn't do what you expect, then implement PCA and consider using $z^{(i)}$