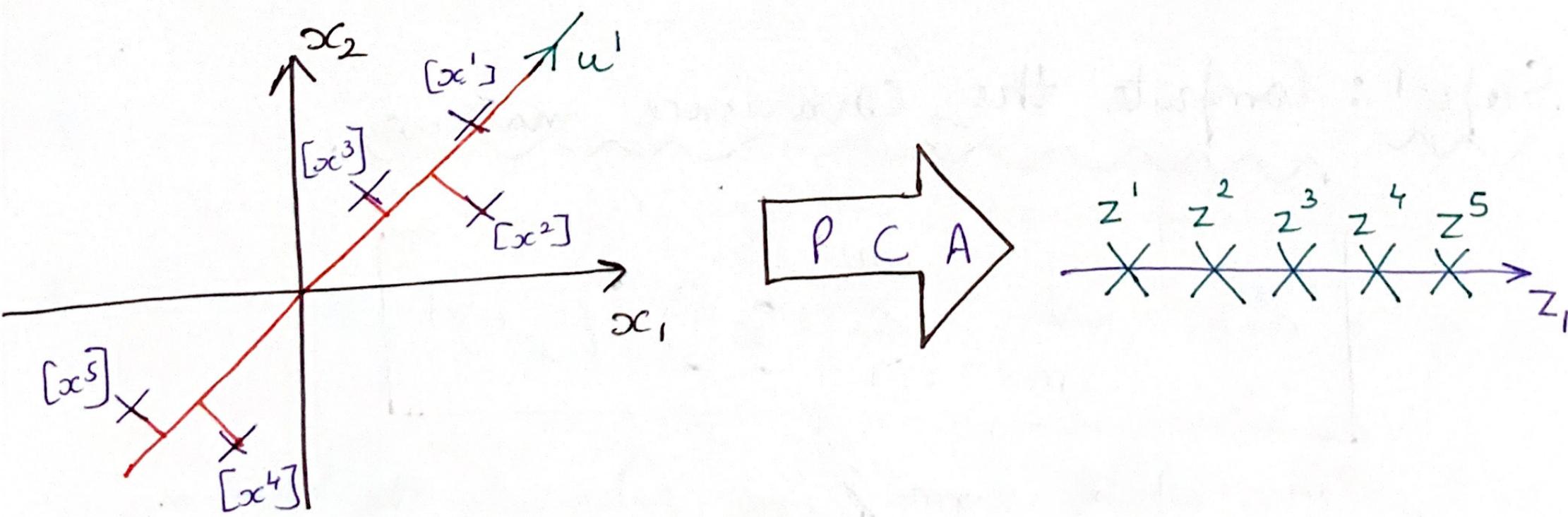
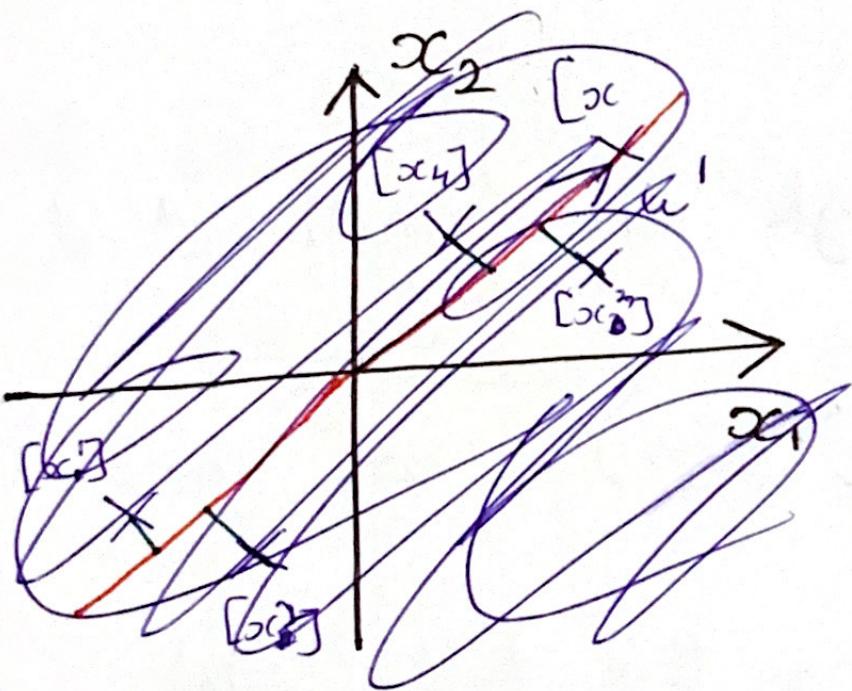


→ With preprocessing done, PCA finds the lower dimensional sub-space which minimizes the projection error.

→ In summary, we (for $2D \rightarrow 1D$, say) are doing the following:



$$x^i \in \mathbb{R}^2 \xrightarrow{\text{PCA}} z^i \in \mathbb{R}^1$$

→ So, we need to compute 2 things:

- 1) the \tilde{u}^i vectors: which "represent" the lower dimensional curve/surface/hyperplane.
- 2) the \tilde{x}^i vectors: which are the new, low dimensional mappings of our features/datapoints.

⇒ Algorithm description: Though the mathematical proof is a bit complicated (so), we can still apply PCA by doing the foll. steps:

Step-1: Compute the "covariance matrix"

$$\Sigma = \frac{1}{m} \sum_{i=1}^m (\tilde{x}^i) (\tilde{x}^i)^T$$

→ This capital sigma sign represents the covariance matrix
→ $\tilde{x}^i \rightarrow (n \times 1)$ matrix

∴ $\Sigma \rightarrow (n \times n)$ matrix

→ Step-2: Compute "eigenvectors" of matrix Σ

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

- * $\text{svd}()$ → singular value decomposition
- * $\text{eig}()$ → also gives eigenvector
- * ~~so~~ svd is more numerically stable than eig.
- * From $[U, S, V] = \text{svd}(\text{sigma})$, we get 3 matrices U, S & V .

Turns out the columns
of U are the u vectors we want!

- * U is also an $[m \times n]$ matrix
- * So, to reduce a system from n -D to k -D,
we just take the first k vectors from U .
(columns)

$$U = [u^1 \ u^2 \ \dots \ u^n] \in \mathbb{R}^{n \times n} \rightarrow U_{\text{reduce}} = [u^1 \ u^2 \ \dots \ u^k] \in \mathbb{R}^{n \times k}$$

→ Step-3 : Take the first k columns of the

U matrix & compute Z

* After computing U reduce, we can calculate the Z vector in the foll. way :

$$Z = (U_{\text{reduce}})^T \cdot x$$



$$Z^i = (U_{\text{reduce}})^T \cdot x$$

$$* Z = [u^1 \ u^2 \ \dots \ u^k]^T \cdot x$$

$$= \begin{bmatrix} (u^1)^T \\ \vdots \\ (u^k)^T \end{bmatrix} \cdot x$$

$\underbrace{\quad}_{k \times m} \qquad \underbrace{m \times 1}$

$$\therefore Z \in \mathbb{R}^{k \times 1}$$

⇒ Generalized MATLAB implementation

→ A tiny detail : Generally, we'll have m training examples having n features. We can represent them using X :

$$X = \begin{bmatrix} - (x^1)^T \\ - (x^2)^T \\ \vdots \\ - (x^m)^T \end{bmatrix}_{m \times n}$$

For computing sigma, we can use ~~log($x^T x$)~~
 $\frac{1}{m}(x^T x)$.

→ Implementation:

$$\text{sigma} = (1/m) * X^T * X \quad \% \text{ Covariance matrix}$$

$$[U, S, V] = \text{svd}(\text{sigma}) \quad \% \text{ Projected directions}$$

$$U_{\text{reduce}} = U(:, 1:k) \quad \% \text{ Take first } k \text{ directions}$$

$$Z = X * U_{\text{reduce}} \quad \% \text{ Projected data points}$$

→ In summary, for a single training example

$$x = \begin{bmatrix} x_1 \\ \vdots \\ x_m \end{bmatrix} \in \mathbb{R}^{m \times 1}$$

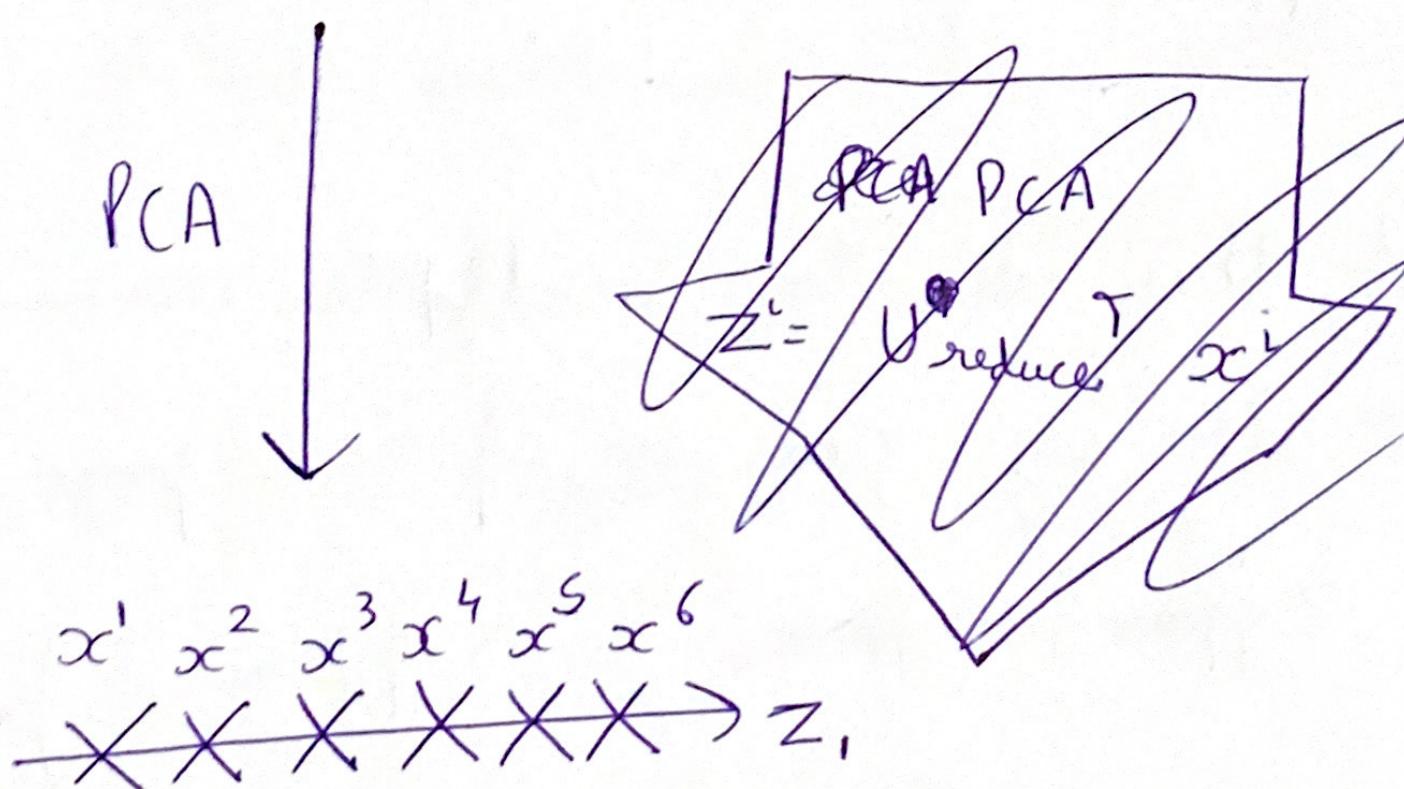
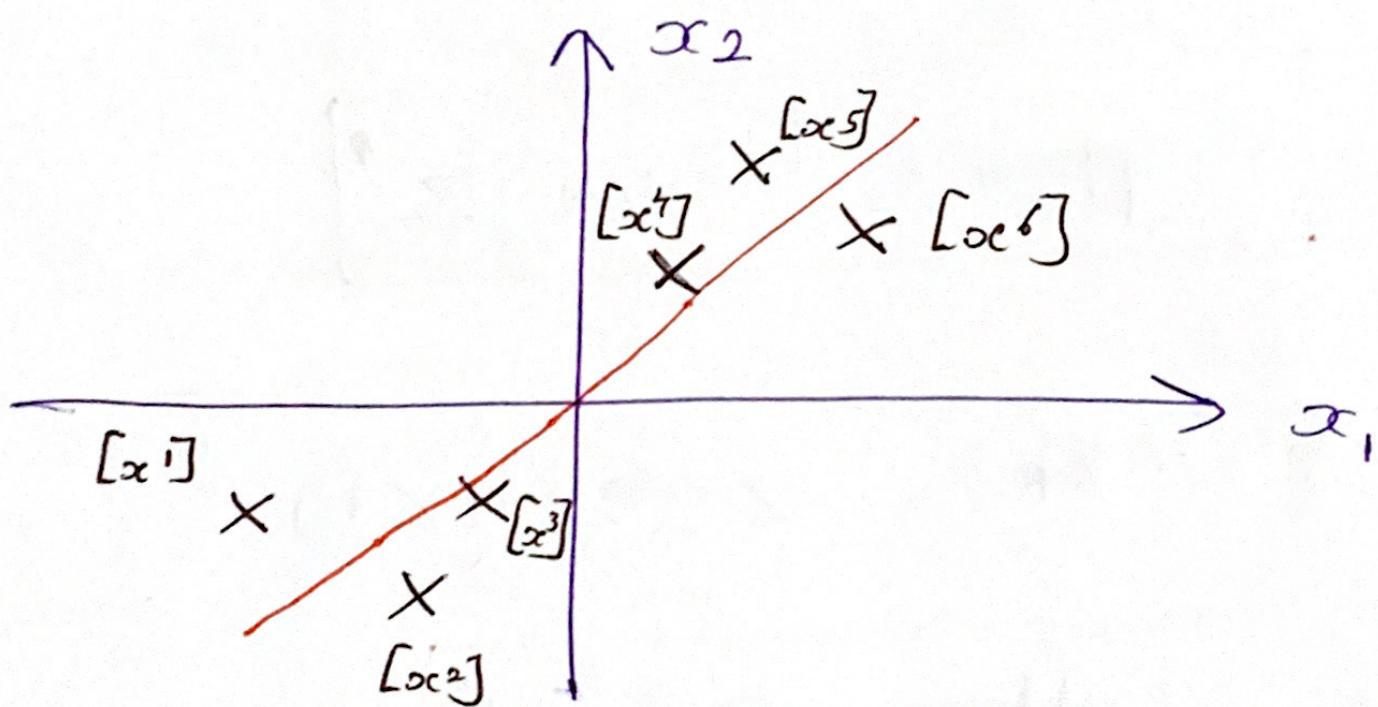
$$z = \begin{bmatrix} 1 & u^1 & u^2 & \dots & u^k \end{bmatrix}^T x = \begin{bmatrix} - (u^1)^T \\ - (u^2)^T \\ \vdots \\ - (u^k)^T \end{bmatrix} x$$

$$\delta z_j = (u^j)^T \cdot x$$

\Rightarrow Reconstruction from compressed representation

\rightarrow Is there a way, to decompress the compressed data, obtained from PCA back to the original?

\rightarrow eg:



→ For a single training example,

$$Z^* = U_{\text{reduce}}^T x^*$$

; $x^* \rightarrow (m \times 1)$, $U_{\text{reduce}} \rightarrow (k \times m)$

To go in the opt. dir^m, we can do

~~x^*~~

$$x_{\text{approx}} = U_{\text{reduce}} Z^*$$

→ $U_{\text{reduce}} \rightarrow (m \times k)$, $Z \rightarrow (k \times 1)$

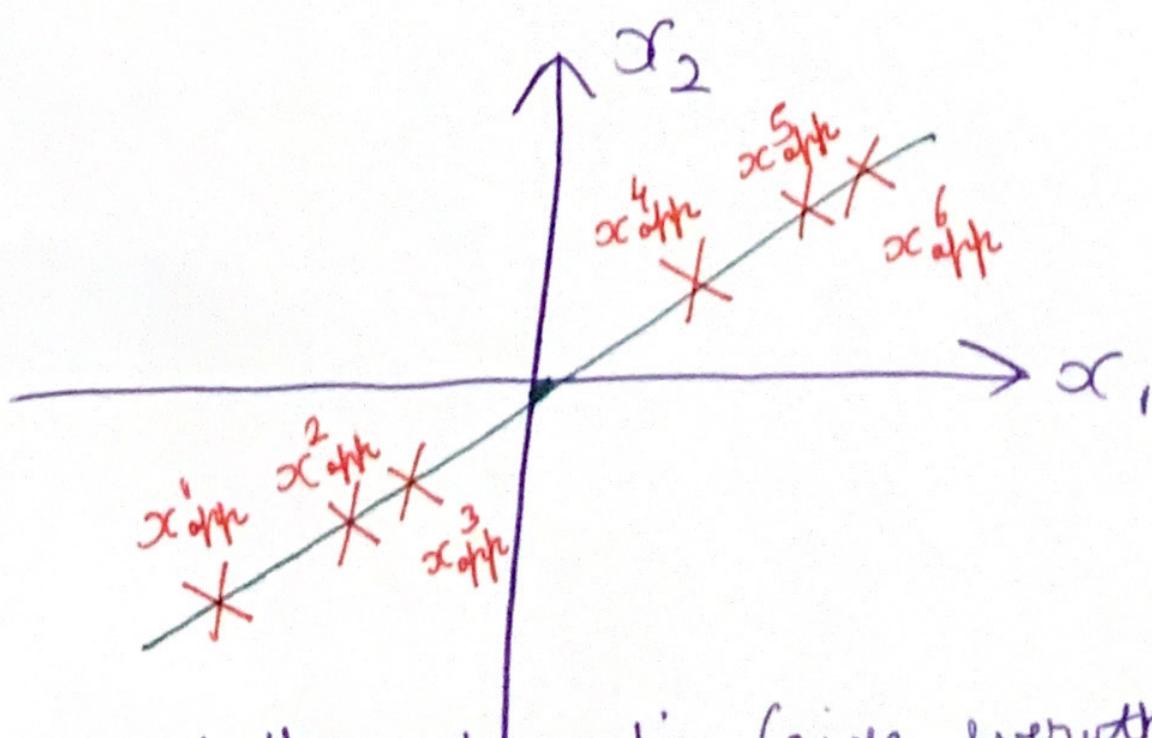
∴

$$x_{\text{approx}} \rightarrow (m \times 1)$$

→ We can extend this idea to multiple training examples! For the eg. given on prev. page, we

can calculate x_{approx}^1 , x_{approx}^2 , ..., x_{approx}^6 & plot

it:



We lose some of the information (given everything is projected + " + it) but it's now projected onto the 2-D space

→ U matrix has a special property that it is a unitary matrix

$$\Rightarrow \boxed{U^{-1} = U^*}$$
 where * means "conjugate transpose"

→ Since we're dealing with real numbers, this is equivalent to:

$$\boxed{U^{-1} = U^T}$$

So, we could use the inverse as well, but it would be a waste of energy & compute cycles.

\Rightarrow choosing the no. of principal components (K)

\rightarrow To choose a suitable val. of k , think about how PCA works:

* It tries to minimize average squared

projection error:

$$\frac{1}{m} \sum_{i=1}^m \|x_i - x_{\text{approx}}^i\|^2$$

* We can also calculate

the total variation in
the data as:

$$\frac{1}{m} \sum_{i=1}^m \|x_i\|^2$$

This basically calculates how far our
datapoints are from the origin.

\rightarrow When choosing k , it's typically to use a ratio
between average squared projection error ~~and~~
total variation with data:

$$\frac{\frac{1}{m} \sum_{i=1}^m \|x_i - x_{\text{approx}}^i\|^2}{\frac{1}{m} \sum_{i=1}^m \|x_i\|^2}$$

→ Algorithm for choosing k

Try PCA with $k = 1, \cancel{2}, \cancel{3}, \cancel{4}, \cancel{5}$

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

Compute the ratio & check if it is < 0.01

* If ratio > 0.01 , then increment k and try again.

→ This procedure is actually horribly inefficient!

We will call SVD as:

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

This returns a matrix S which is a $(m \times m)$ matrix & structured in foll. way:

$$S = \begin{bmatrix} S_{11} & 0 & \cdots & 0 \\ 0 & S_{22} & \cdots & 0 \\ \vdots & \cdots & \ddots & \vdots \\ 0 & \cdots & \cdots & S_{mm} \end{bmatrix}_{m \times m} \rightarrow \text{Diagonal matrix}$$

→ We typically want this ratio to be < 0.01
or $< 1\%$ (when expressed in percentages).

* In other words, the squared projection error divided by the total variation should be less than one percent,
so that "99 % of the variance is retained".

* Simple Block

Block

* ~~Simple Block~~ ^{said} : If you project your points on a surface $\in \mathbb{R}^k$ (s.t. if the ratio < 0.01 comes out), you will reconstruct the points with 99% accuracy. This is what is meant by "99 % variance is retained".

→ So, now we have to ensure that:

$$1 - \left(\frac{\sum_{i=1}^k s_{ii}}{\sum_{i=1}^m s_{ii}} \right) < 0.01$$

→ Now, this is extremely efficient!! Why?

2 reasons:

1) Now, we don't need to evaluate SVD

over S over for diff. values of k.

We'll compute S once ~~and keep it as~~

S ~~will~~ will ~~keep~~ using the rows
from S corresponding to the value of k in
the loop.

2) Since, we'll only proceed ~~on~~ to calculate

U, Z, Sx after the ratio is < 0.01 ,

we don't have to calculate these ~~as well~~

~~as~~ for every value of k.

→ Updated algo :

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

Pick smallest value of k for which :

$$\frac{\sum_{i=1}^k S_{ii}}{\sum_{i=1}^m S_{ii}} \geq 0.99$$

(99 % of variance retained)

Compute ~~K~~ PCA for the selected val. of k

- ⇒ Advice for applying PCA
- ⇒ Supervised learning speedup
- The most common use of PCA is to speed up supervised learning.
- Given a training set with a large no. of features (eg: $x^1, \dots, x^m \in \mathbb{R}^{1000}$) we can use PCA to reduce the no. of features in each example of the training set (eg: $z^1, \dots, z^m \in \mathbb{R}^{100}$)
- Applying PCA to map $x^i \rightarrow z^i$ should only be done on the training set:
- * Parameters The matrix U is "learned" through PCA.
 - * So, map z^i to your cross validation & test sets ~~@~~ ONLY after it (the matrix U) is defined on the training set.

→ Realistically, it is ~~computing power limit~~ to reduce data dimensionality by a factor of 5-10 without a major hit to the learning (or classification) algorithm.

7 Applications of PCA

→ Compression:

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

→ Visualization: As we can plot only in 2D/3D, we can use $k=2$ or $k=3$ to visualize large dimensional datasets. PCA is often used in this scenario.

* From ex 7 ~~PSDF~~: The PCA projection can be thought as a rotation that selects the view that maximizes the spread of the data, which often corresponds to the

⇒ MISUSE of PCA : To prevent overfitting



- A general misconception is that reducing the dimension of dataset by applying PCA can prevent ~~overfitting~~ overfitting as "fewer features, less likely to overfitting".
- This MIGHT work but it isn't a good way to address overfitting. Use regularization instead. It'll work way better!

→ PCA doesn't consider the values of the outcome's



* It is very possible that PCA throws away some useful data which might have led to better results.

* Regularization does a much better job as it keeps in mind, the effect of outcome values (y) while preventing overfitting.

⇒ PCA is sometimes used where it shouldn't be

Before implementing

→ Don't assume that you absolutely need to do PCA !

* Try out your full machine learning algorithm/pipeline/workflow without PCA first. Then use PCA if you find that you need it.

* Of course, if you have a dataset requiring more disk space than available, then you have to do it !