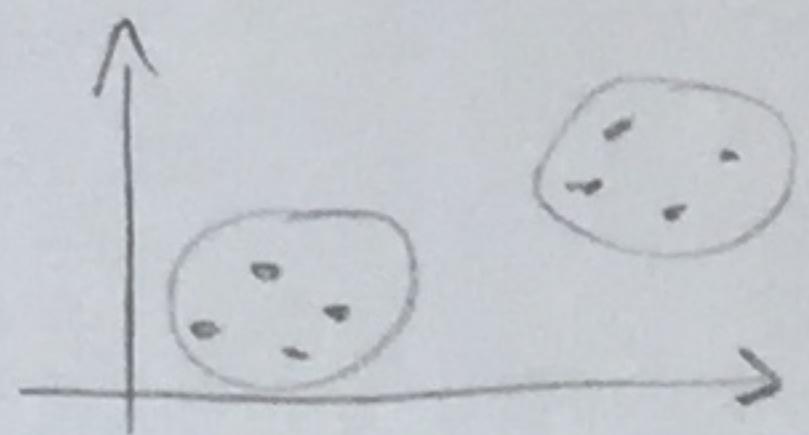


Unsupervised Learning: Clustering

\Rightarrow Unlabeled data (no y vector) \rightarrow given data, find structure in data!

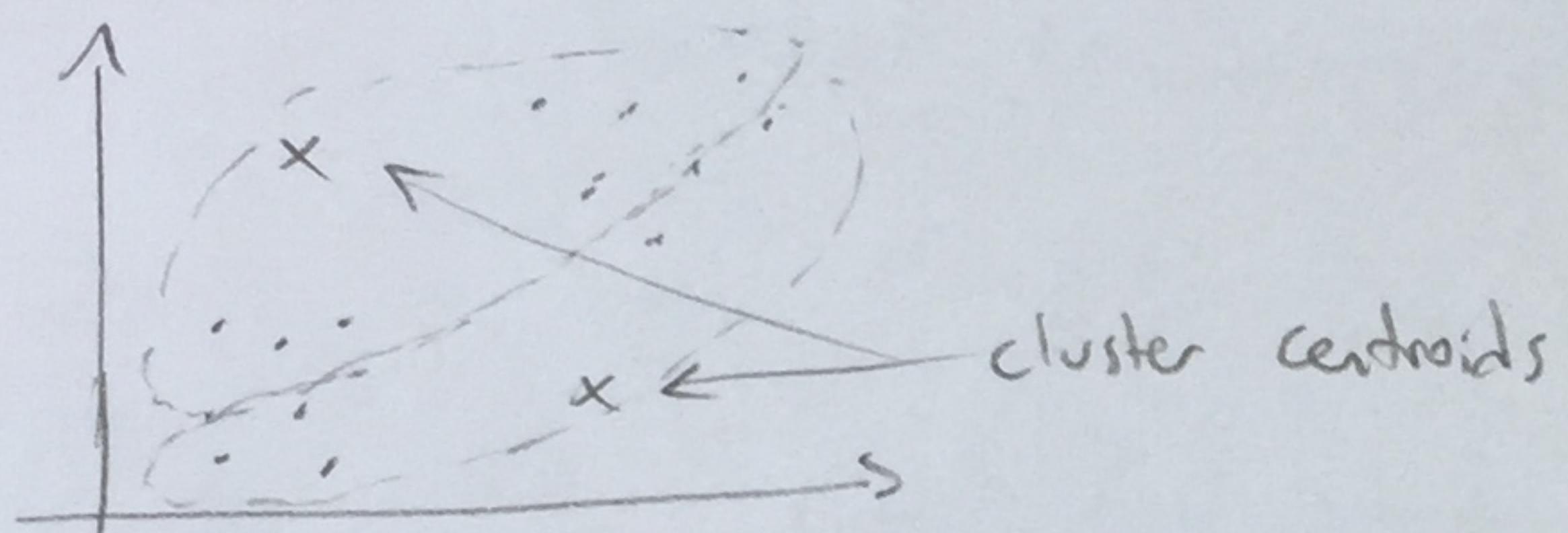
e.g. find clusters:



\Rightarrow Clustering \rightarrow social networks / graphs, market segmentation, resource allocation, galaxies, etc.

\Rightarrow K-Means Algorithm: groups data into coherent subsets (clusters)

e.g.



\Rightarrow assign all pts to one centroid or another (cluster assignment)

\Rightarrow then move cluster centroids to mean location of all pts in their group

\hookrightarrow eventually, cluster centroids converge to single location.

\Rightarrow parameters in k-means: $\cdot K$ (# clusters)

\cdot training set $(x^{(1)}, \dots, x^{(m)})$

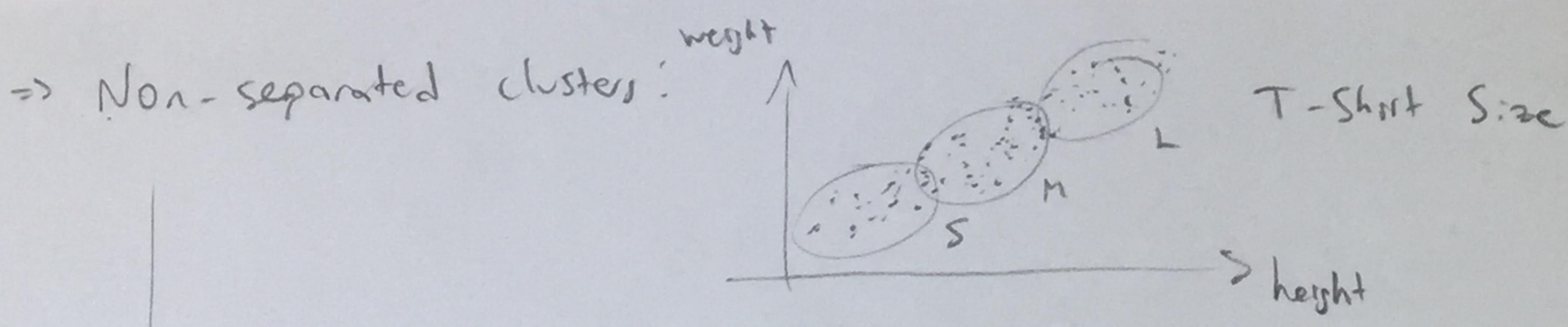
\hookrightarrow note no $x_0 \rightarrow x^{(i)} \in \mathbb{R}^n$

Algorithm: Randomly initialize K cluster centroids $\mu_1, \mu_2, \dots, \mu_K \in \mathbb{R}^n$

repeat:

CLUSTER ASSIGNMENT $\left\{ \begin{array}{l} \text{for } i = 1 \text{ to } m \text{ (all pts.)} \\ c^{(i)} = \text{index (from 1 to } K \text{) of centroid closest to } x^{(i)} \quad (\min \|x^{(i)} - \mu_k\|^2) \end{array} \right.$

MOVE CLUSTER CENTROIDS $\left\{ \begin{array}{l} \text{for } k = 1 \text{ to } K \\ \mu_k = \text{average of pts. assigned to cluster } k \quad (\mu_k \in \mathbb{R}^n) \end{array} \right.$



↳ useful for market segmentation!

Optimization Objective of K-Means:

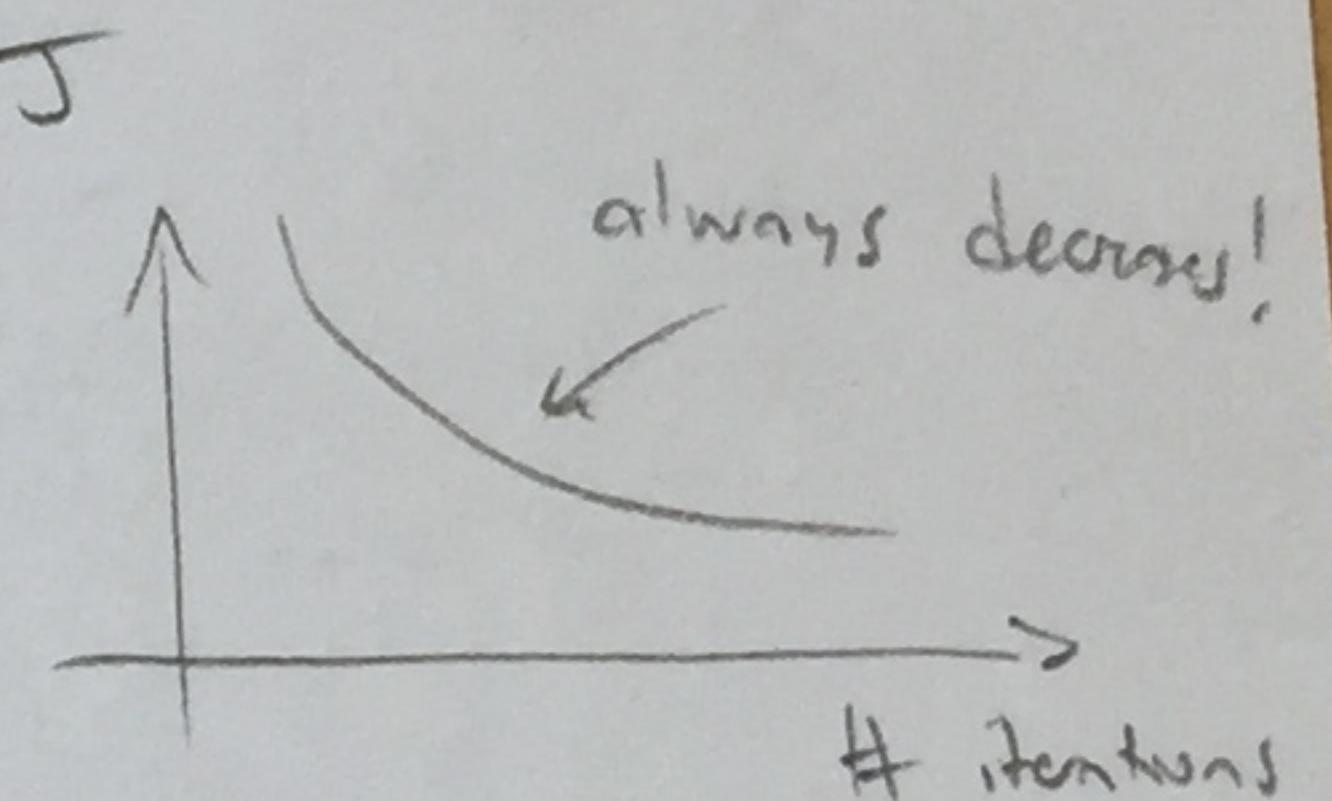
⇒ keeping track of $c^{(i)}$'s (cluster index to which given $x^{(i)}$ is assigned) and μ_k (loc. of cluster centroid) → combined, $\mu_{c^{(i)}}$ is centroid to which $x^{(i)}$ assigned

$$J(c^{(1)} \dots c^{(m)}, \mu_1 \dots \mu_K) = \frac{1}{m} \sum_{i=1}^m \|x^{(i)} - \mu_{c^{(i)}}\|^2 \rightarrow \text{average squared distance b/w } x^{(i)} \text{ & centroid it is assigned to}$$

$$\min_{c^{(1)} \dots c^{(m)}, \mu_1 \dots \mu_K} J(c^{(1)} \dots c^{(m)}, \mu_1 \dots \mu_K) \quad (\text{Distortion cost function})$$

↳ cluster assignment step minimizes $J_{c^{(1)} \dots c^{(m)}}$ holding μ 's fixed

move centroid step minimizes J w.r.t. $\mu_1 \dots \mu_K$



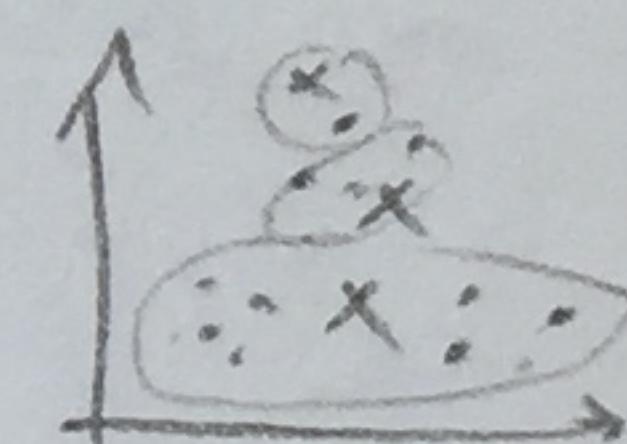
Random Initialization: best method:

⇒ $K < m$. Randomly pick K examples → set $\mu_1 \dots \mu_K =$ these examples

⇒ depending on initialization, can hit local optimum →

↳ try multiple initial guesses to get global optimum

for $i=1 \rightarrow 100$: run K-means, get $J^{(i)}$ → pick lowest $J^{(i)}$
(w/ unique initialization)

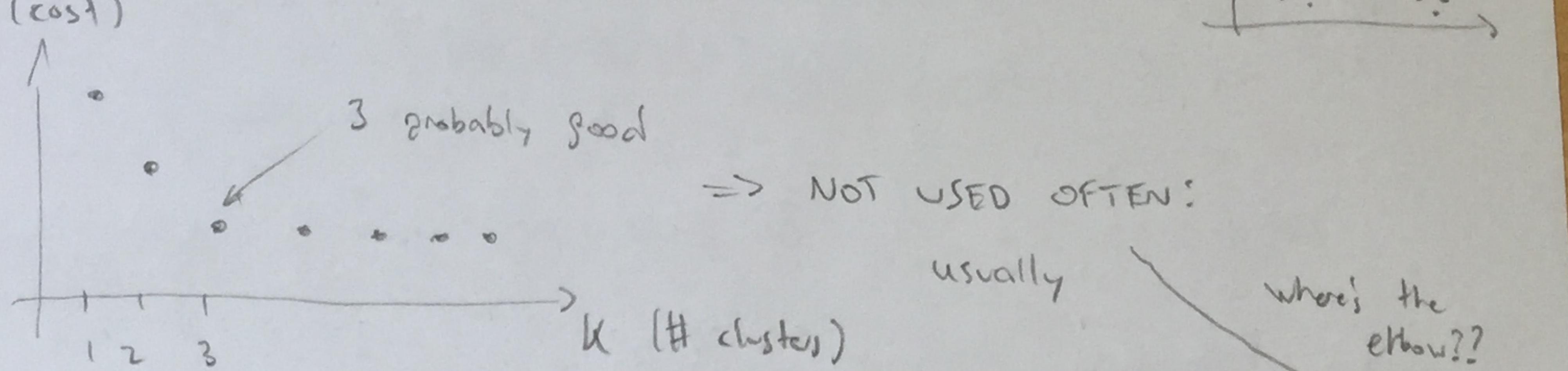


\Rightarrow Choosing K , # clusters:

\hookrightarrow most commonly \rightarrow visually / by hand

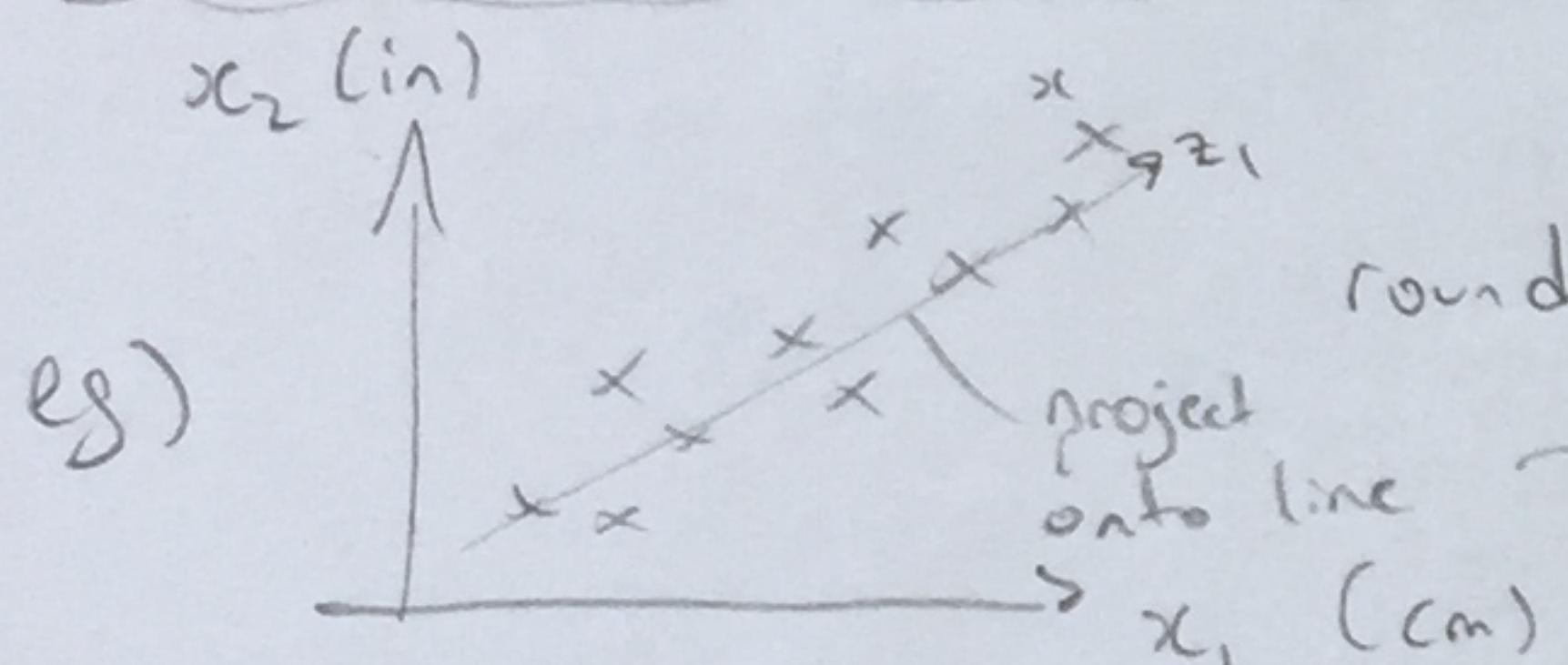
\hookrightarrow sometimes hard / ambiguous (part of unsupervised learning)

\Rightarrow elbow method:



\Rightarrow can choose # clusters based off results in some downstream purpose
eg. # t-shirt sizes

Dimensionality Reduction: compression



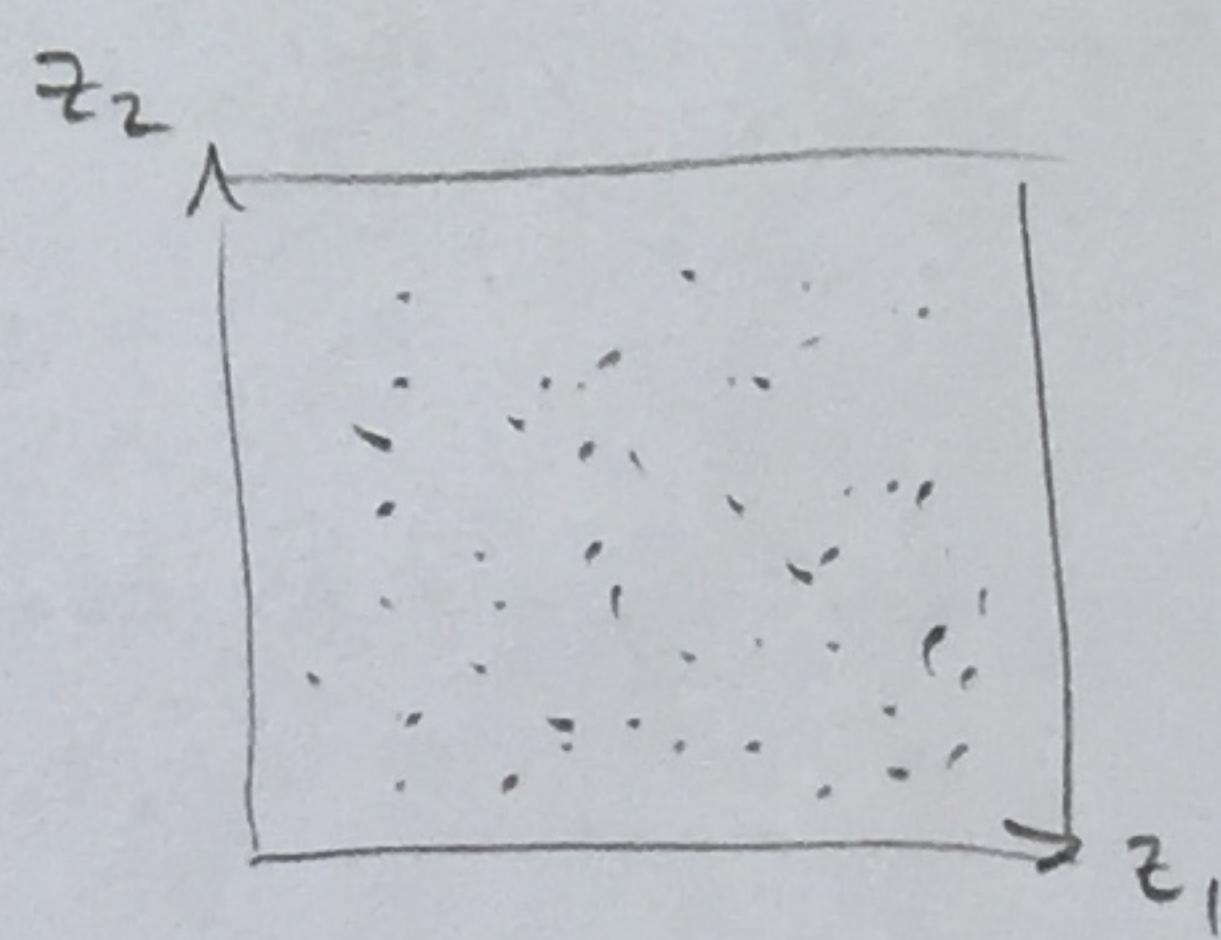
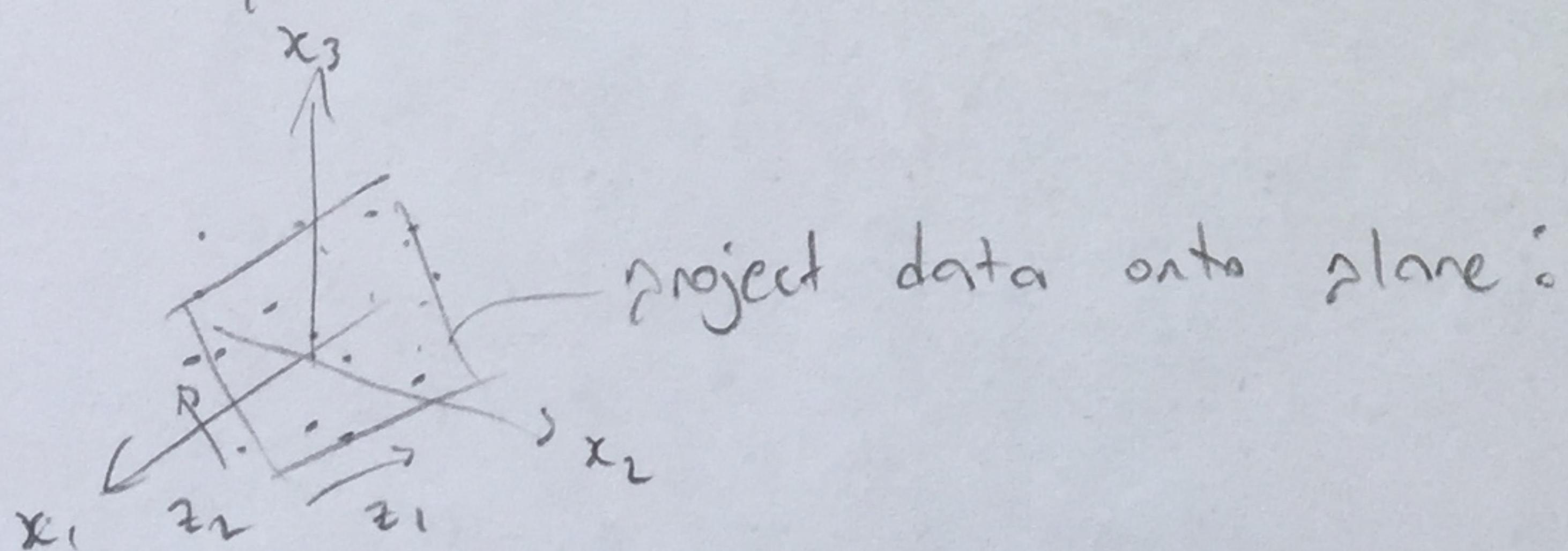
\Rightarrow redundant! Only need 1D, not 2D

$$x \in \mathbb{R}^2$$
$$z \in \mathbb{R}^1$$

\hookrightarrow approximation of original training set!

\Rightarrow reduces disc space req'd & speeds up learning algorithms

\Rightarrow Compression \rightarrow 3D to 2D:



\hookrightarrow data should already more or less lie on plane (strong correlation)

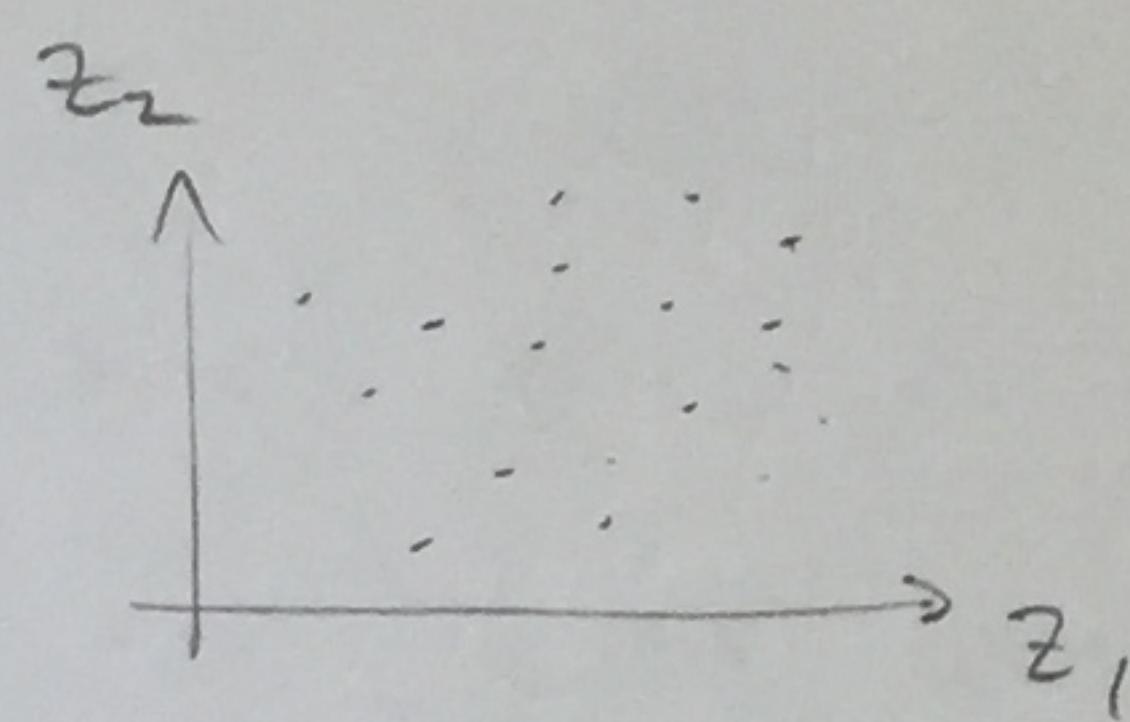
Higher Dimensional Data Viz:

Say $x \in \mathbb{R}^{50}$. Reduce $x \rightarrow z \in \mathbb{R}^2$

e.g. $z_1 \propto$ country size $\quad x_1 = \text{pop.}$

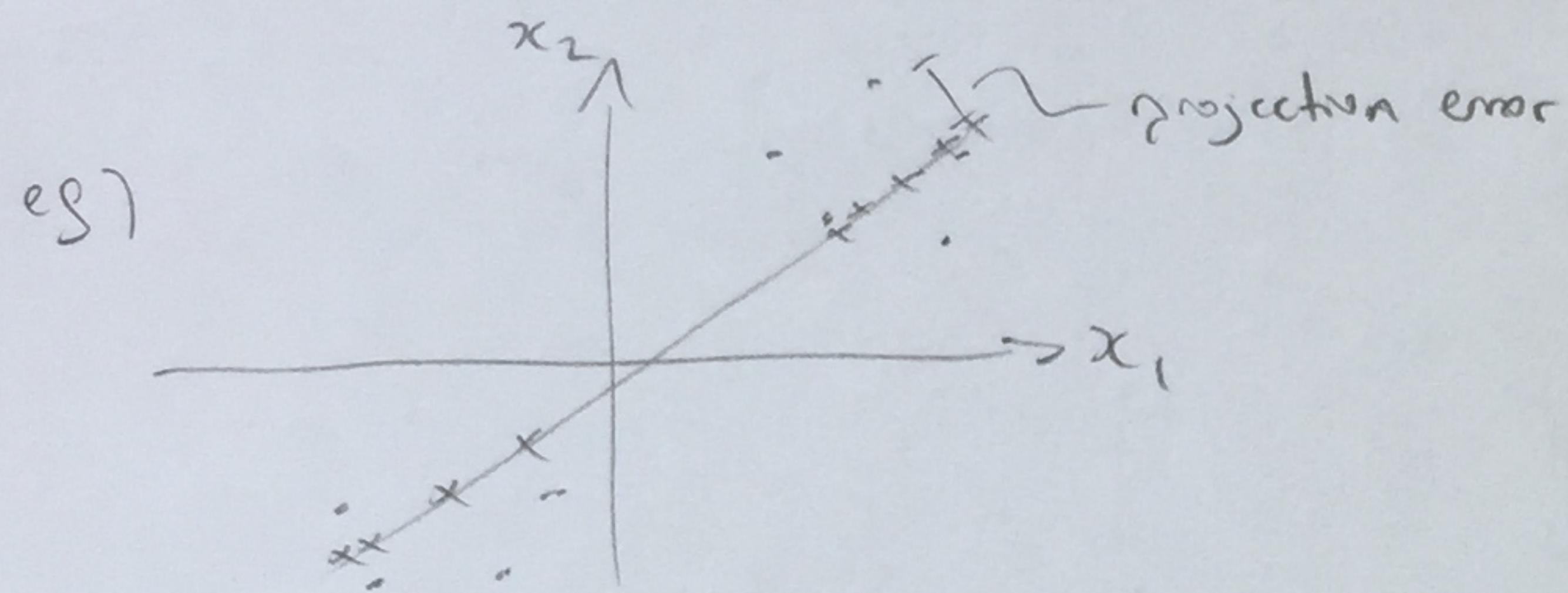
$z_2 \propto$ per person GDP $\quad x_2 = \text{GDP}$

$x_3 = \text{land area etc.}$



\Rightarrow reduce to 2 or 3D for plotting \rightarrow PCA helps do this

Principal Component Analysis: PCA



1) Mean normalization & feature scaling x_1, x_2

2) Minimize sum-of-squares of projection error

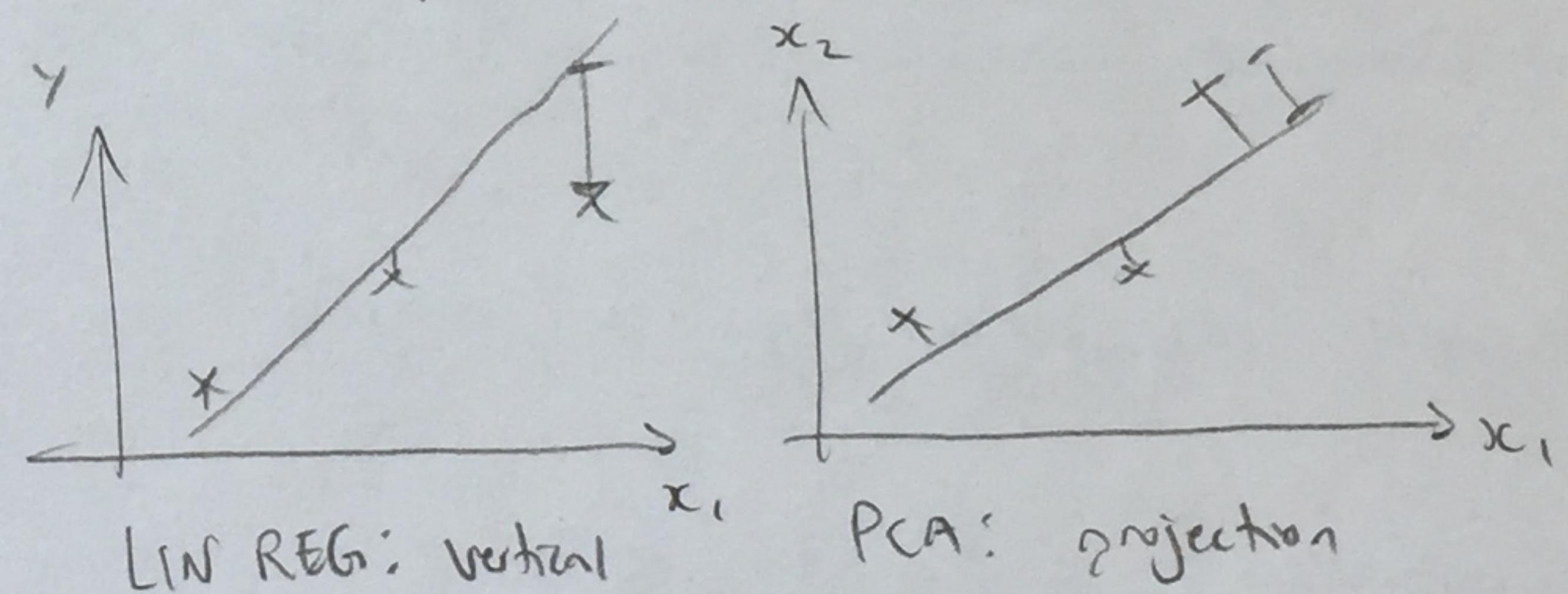
\hookrightarrow find $\underline{u}_{\text{best}}$ to project data

\Rightarrow reducing from n -dim to k -dim: Find k -vectors $\underline{u}^{(1)}, \underline{u}^{(2)}, \dots, \underline{u}^{(k)}$

onto which to project data (e.g. 3D \rightarrow 2D, find $\underline{u}^{(1)}$ & $\underline{u}^{(2)}$ to define plane)

\Rightarrow PCA is NOT linear regression!

\hookrightarrow no special feature "y" being predicted



\Rightarrow Algorithm:

1) Pre-processing: mean norm/feature scaling: $\mu_j = \frac{1}{m} \sum_{i=1}^m x_j^{(i)}$

\hookrightarrow replace $x_j^{(i)}$ w/ $\frac{x_j^{(i)} - \mu_j}{s_j}$

s_j \leftarrow standard dev of feature j

2) Reducing n -dim to k -dim:

\hookrightarrow compute covariance matrix: $\Sigma = \frac{1}{m} \sum_{i=1}^n (x^{(i)}) (x^{(i)})^T$

\hookrightarrow compute eigenvectors of $\Sigma \rightarrow [U, S, V] = \text{svd}(\Sigma)$,

\hookrightarrow singular value decomposition

$\Rightarrow \underline{U}$ matrix is $n \times n$, columns will be $\underline{u}^{(1)}, \underline{u}^{(2)}, \underline{u}^{(3)}$ etc.

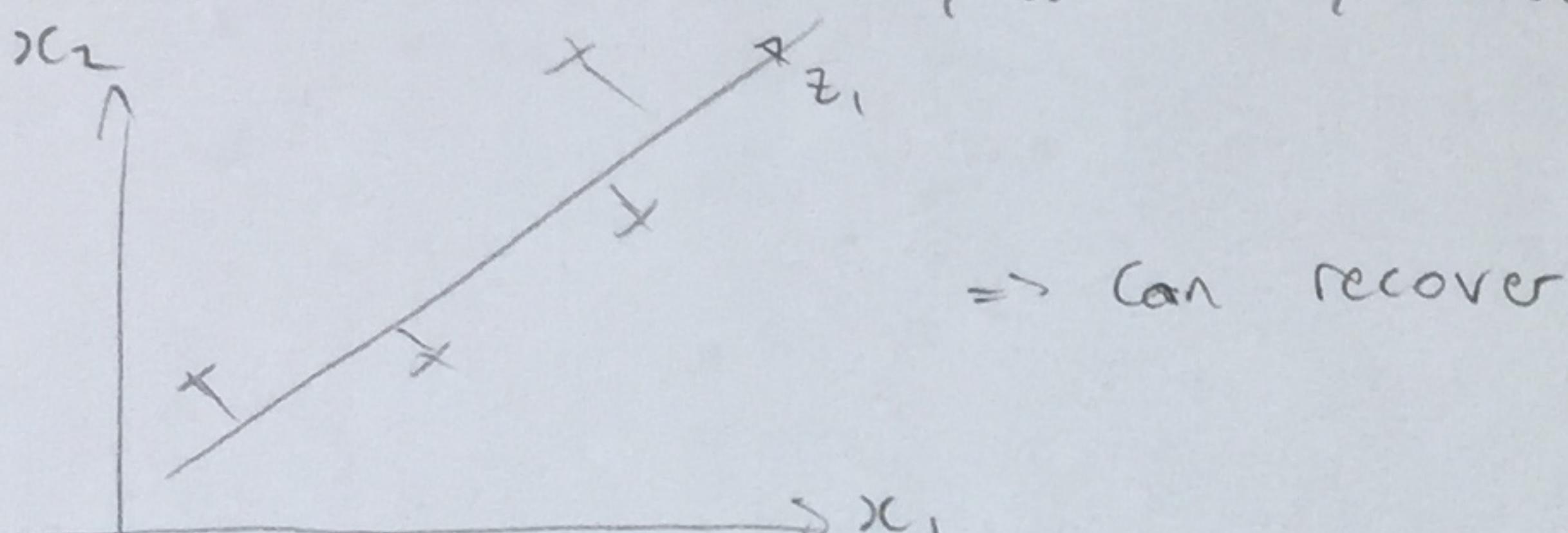
\hookrightarrow take first k $\underline{u}^{(i)}$ vectors!

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

$$\hookrightarrow \underline{U}_{\text{reduced}} = [n \times k] \rightarrow z = \underline{U}_{\text{red}}^T x \in [k \times 1]$$

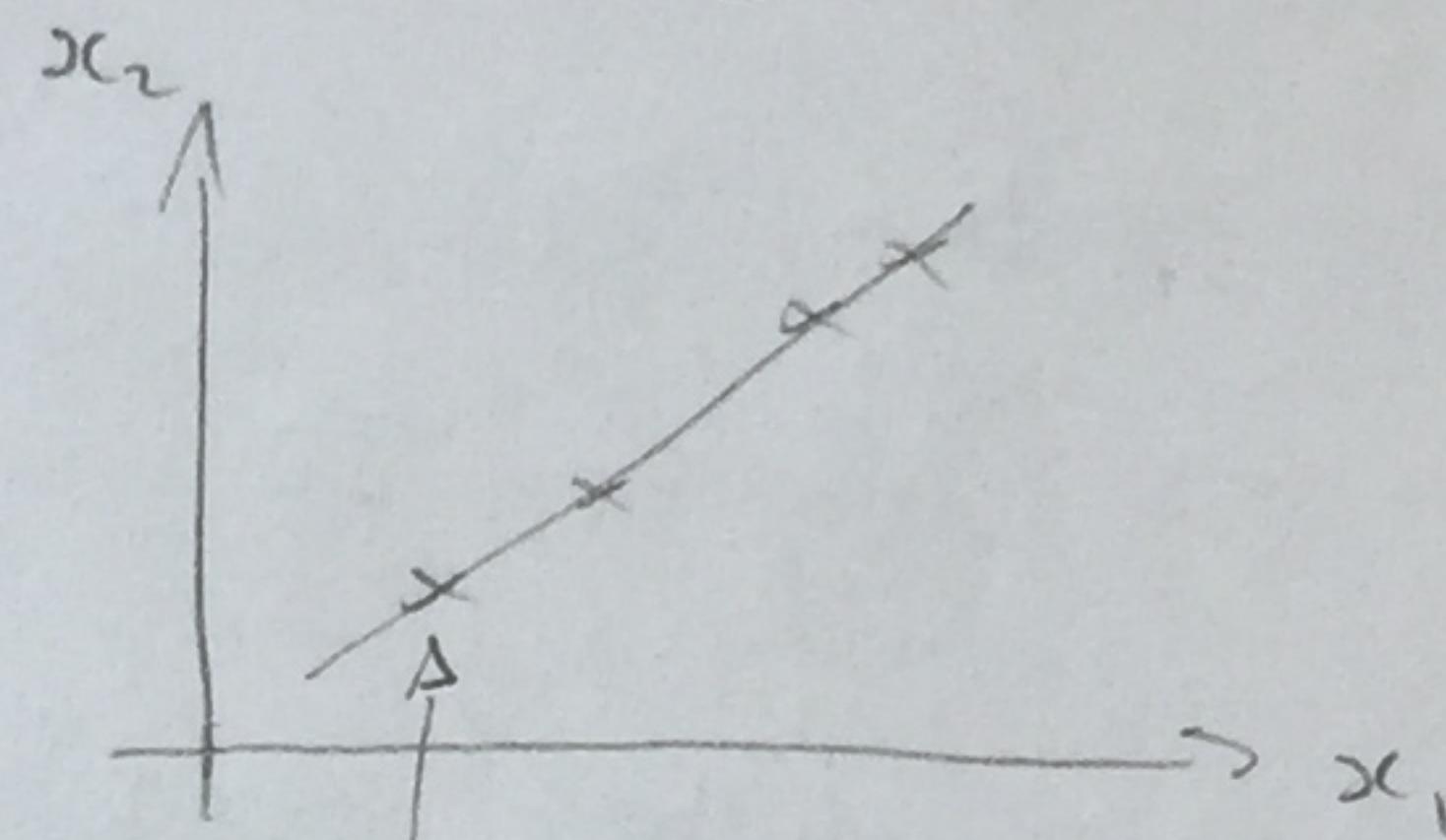
$[k \times n] \quad [n \times 1]$

\Rightarrow Reconstruction from Compressed Representation $\rightarrow z \Rightarrow \underline{x}$!



$$x_{\text{approx}}^{(i)} = \underline{U}_{\text{red}} z^{(i)} \in [n \times 1]$$

$[n \times k] \quad [k \times 1]$



approximate reconstruction of data

\Rightarrow Choosing # Principal Components: Choosing k

\Rightarrow PCA minimizes $\frac{1}{m} \sum_{i=1}^m \|x^{(i)} - x_{\text{approx}}^{(i)}\|^2$ Define variation $\frac{1}{m} \sum_{i=1}^m \|x^{(i)}\|^2$

\hookrightarrow choose smallest k such that $\frac{\text{avg. square proj. error}}{\text{variation}} \leq 0.01$ (1%)

\hookrightarrow "99% of variance is retained"

algorithm: PCA w/ $k=1 \rightarrow \underline{U}_{\text{red}}, \underline{x} \& \underline{z} \rightarrow$ check variance retained

\hookrightarrow increase k until variance constraint satisfied

$$\Rightarrow \text{helpful: for given } k, \frac{\text{avg sq. proj. error}}{\text{variation}} = 1 - \frac{\sum_{i=1}^k S_{ii}}{\sum_{i=1}^n S_{ii}} \stackrel{\leq \text{ diagonal}}{=} \begin{bmatrix} S_{11} & 0 & 0 \\ 0 & S_{22} & 0 \\ 0 & 0 & \ddots \end{bmatrix}$$

\hookrightarrow don't need to find eigenvalues every time!

$$\Rightarrow \text{pick smallest } k \text{ for } \frac{\sum_{i=1}^k S_{ii}}{\sum_{i=1}^n S_{ii}} \geq 0.99$$

Advice for PCA:

⇒ Speeding up algorithms: Say 10,000 features, $\underline{x} \in \mathbb{R}^{10,000}$

↳ extract inputs $\underline{x}^{(1)} \rightarrow \underline{x}^{(m)}$, apply PCA such that $\underline{z}^{(1)} \rightarrow \underline{z}^{(m)}$

↳ then pair $(\underline{z}^{(1)}, y^{(1)}) \dots (\underline{z}^{(m)}, y^{(m)})$ → get $h_{\theta}(z) = \frac{1}{1 + e^{-\theta^T z}}$

⇒ take $\underline{x} \rightarrow$ map to $\underline{z} \rightarrow h_{\theta}(z)$

* Note: map $\underline{x}^{(1)} \rightarrow \underline{z}^{(1)}$ only on training set! Apply mapping to $\underline{x}_{cv}, \underline{x}_{test}$

⇒ For compression → retain 99% variance

For visualization → 2D or 3D

⇒ Misuse of PCA → avoiding overfitting (reducing # features)

↳ use regularization instead! PCA loses data/info w/o looking at y vals.

⇒ Always implement model w/o PCA first (most info retained), then use PCA if memory or speed unsatisfactory