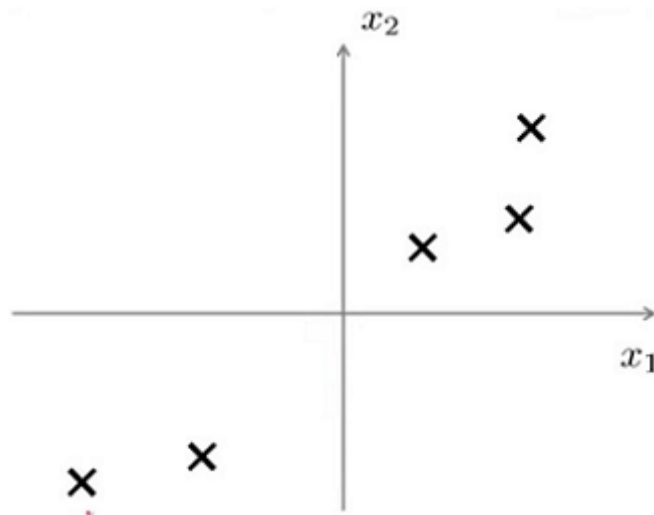


Principal Component Analysis (PCA): Problem Formulation

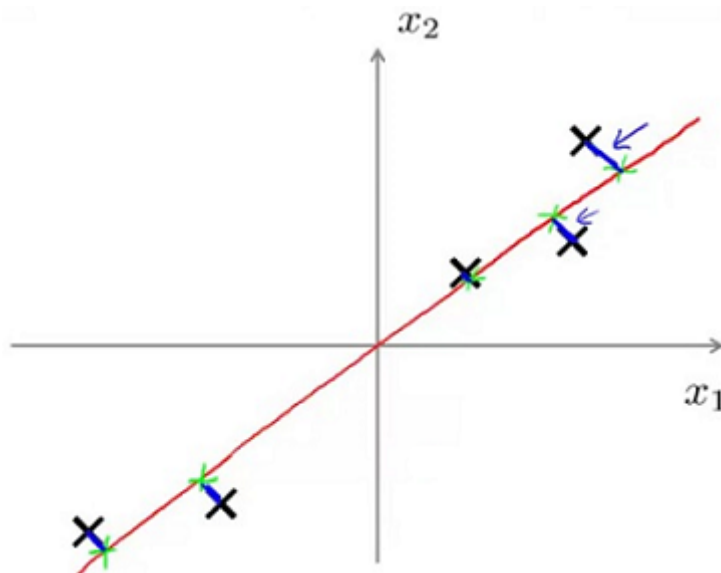
For the problem of dimensionality reduction, by far the most popular and the most commonly used algorithm is something called principal components analysis, or PCA. Here, we'll start talking about how we formulate precisely what we want PCA to do.

Let's say we've a data set with $x \in \mathbb{R}^2$ and let's say we want to reduce the dimension of the data from two-dimensional to one-dimensional, in other words we would like to find a line onto which to project the data.

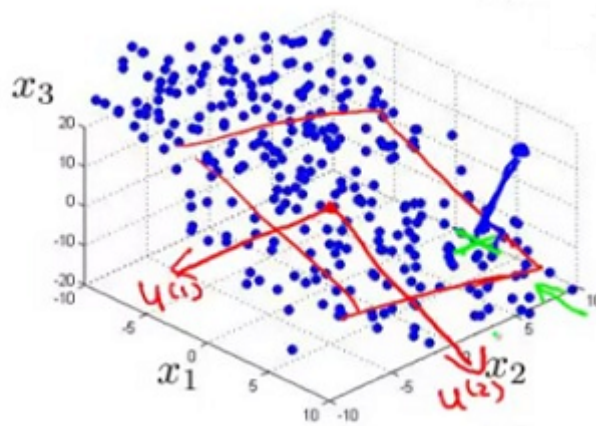
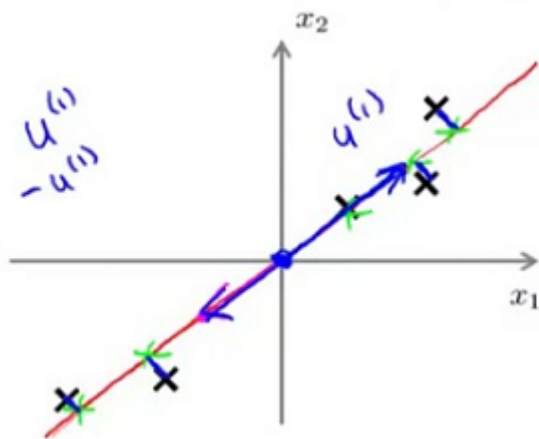


How do we determine this line?

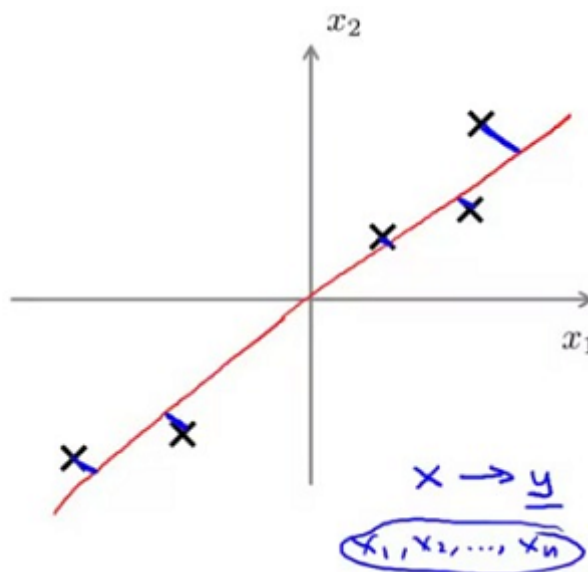
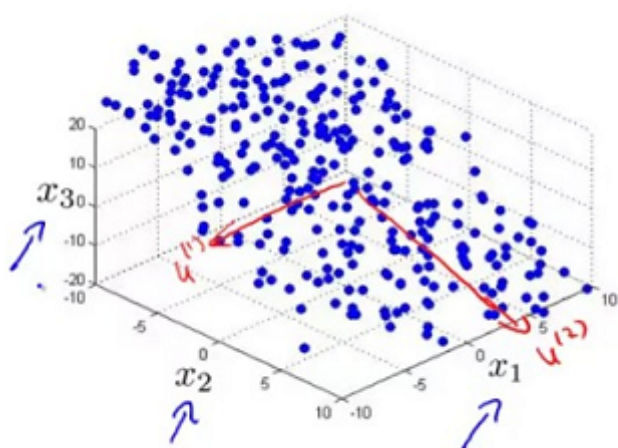
- What we find is that the distance between each point and the projected version is pretty small.
- PCA tries to find a lower dimensional surface so the sum of squares onto that surface is minimized
- The blue lines are sometimes called **the projection error**
 - PCA tries to find the surface (a straight line in this case) which has the minimum projection error



As an aside, before applying PCA, it's standard practice to first perform **mean normalization** and **feature scaling** on our data before PCA.



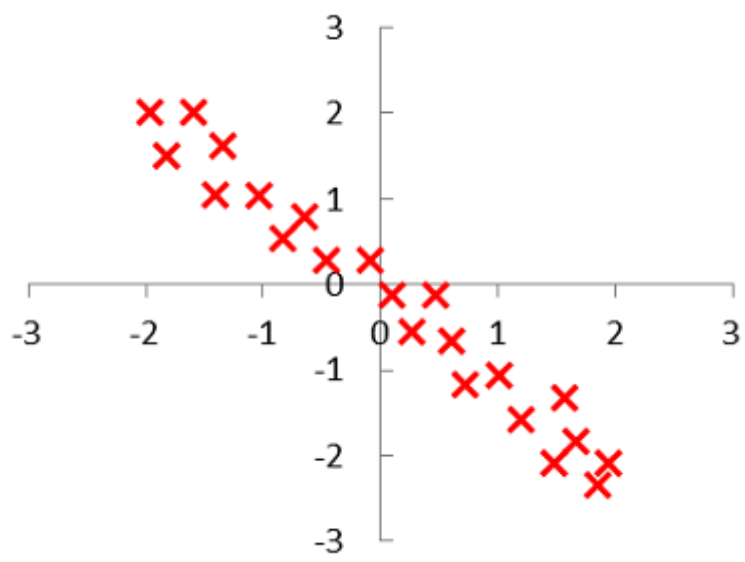
- A more formal description is:
 - **Goal of PCA** - Reduce from 2-dimension to 1-dimension: Find a direction (a vector $u^{(1)} \in \mathbb{R}^n$) onto which to project the data so as to minimize the projection error.
 - $u^{(1)}$ can be positive or negative ($-u^{(1)}$) which makes no difference
- In the more general case:
 - Reduce from n -dimensional to k -dimension: Find k vectors $u^{(1)}, u^{(2)}, \dots, u^{(k)}$ onto which to project the data so as to minimize the projection error.
 - We can define a point in a plane with k vectors
- If we have a 3D point cloud - ($3D \rightarrow 2D$):
 - Find pair of vectors which define a 2D plane (surface) onto which we're going to project our data



- How does PCA relate to linear regression?
 - PCA is not linear regression (Despite cosmetic similarities, very different)
 - For linear regression, fitting a straight line to minimize the straight line between a point and a squared line
 - For PCA minimizing the magnitude of the shortest orthogonal distance
 - More generally
 - With linear regression we're trying to predict "y"
 - With PCA there is no "y" - instead we have a list of features and all features are treated equally
 - If we have 3D dimensional data $3D \rightarrow 2D$
 - Have 3 features treated symmetrically

Video Question: Suppose you run PCA on the dataset below. Which of the following would be a reasonable vector $u^{(1)}$ onto which to project the data? (By convention, we choose $u^{(1)}$ so that

$$\|u^{(1)}\| = \sqrt{(u_1^{(1)})^2 + (u_2^{(1)})^2}, \text{ the length of the vector } u^{(1)}, \text{ equals 1.})$$



- $u^{(1)} = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$
- $u^{(1)} = \begin{bmatrix} 0 \\ 1 \\ 1 \end{bmatrix}$
- $u^{(1)} = \begin{bmatrix} 1/\sqrt{2} \\ 1/\sqrt{2} \end{bmatrix}$

$$u^{(1)} = \begin{bmatrix} -1/\sqrt{2} \\ 1/\sqrt{2} \end{bmatrix}$$