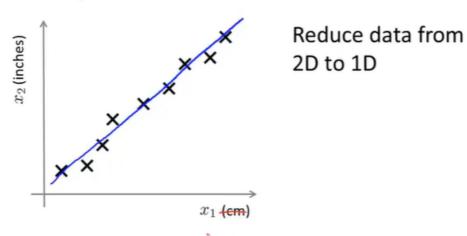
## 降维

### 应用1.数据压缩

比如一个特征是长度(cm),另一个是长度(inches),这样特征就高度冗余,两个特征可以降维为一个。

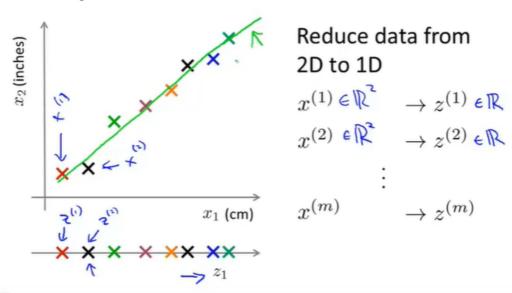
或者一个特征表示飞行员的技能,另一个是他享受飞行的程度,这两个特征也是高度相关的,这种情况也需要降维。

### **Data Compression**

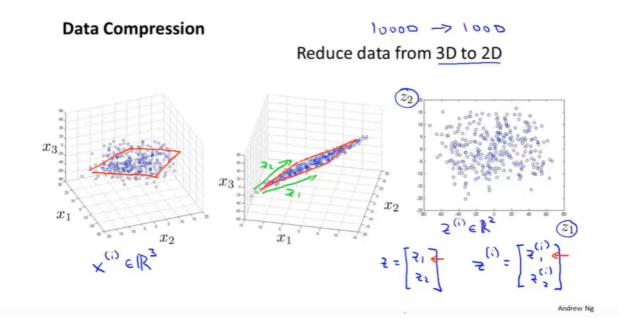


我们将坐标投影到z1这条线上,本来用二维向量表示的一个数据点,现在用一个实数就可以表示,可以将内存需求减半,另外这也可以让学习算法运行地更快。

### **Data Compression**



三维降维到两维:将三维空间内的数据点投影到一个二维平面上。这样只需要z1,z2两个点就可以确认一个点的位置。



# 应用2.可视化数据

比如有50个指标数据来衡量一个国家,我们不想用50个数字去代表一个国家。

Data Visualization							
						XL	
	XI	X2	V		Xs	Mean	
		Per capita	X3	Хų	Poverty	household	
	GDP	GDP	Human	7-4	Index	income	
	(trillions of	(thousands	Develop-	Life	(Gini as	(thousands	
Country	US\$)	of intl. \$)	ment Index	expectancy	percentage)	of US\$)	
Canada	1.577	39.17	0.908	80.7	32.6	67.293	
China	5.878	7.54	0.687	73	46.9	10.22	
India	1.632	3.41	0.547	64.7	36.8	0.735	
Russia	1.48	19.84	0.755	65.5	39.9	0.72	
Singapore	0.223	56.69	0.866	80	42.5	67.1	
USA	14.527	46.86	0.91	78.3	40.8	84.3	

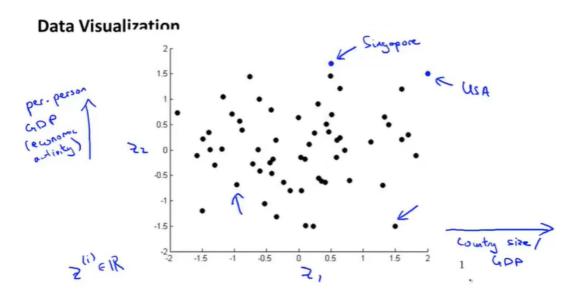
[resources from en.wikipedia.org]

Andrew Ng

我们可以尝试用一个二维的向量z来表示一个国家。相当于将50D降维到2D。可以在一个二维图中表示每个国家。

Data Visualization	2 (1) EIR2		
Country	$z_1$	$z_2$	£ 6/10
Canada	1.6	1.2	
China	1.7	0.3	
India	1.6	0.2	
Russia	1.4	0.5	
Singapore	0.5	1.7	
USA	2	1.5	

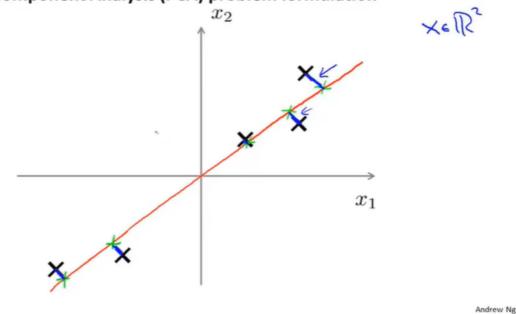
50个特征可以用国家面积和GDP来表示,并很好的理解这俩这两个数据。



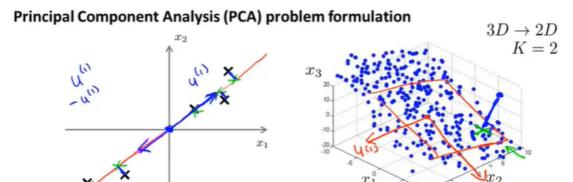
# PCA(主成分分析方法)降维 (1)

在二维降维到一维过程中,需要找到一条直线,将各个点投影到这条直线上,投影的蓝色线段长度称为投影误差。PCA需要找到一条投影误差最小的线。

### Principal Component Analysis (PCA) problem formulation



PCA就是找出一条直线,一个平面,或者其他维空间,来对数据样本点进行投影。

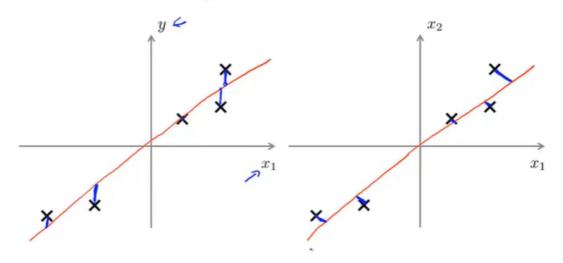


Reduce from 2-dimension to 1-dimension: Find a direction (a vector  $\underline{u}^{(1)} \in \mathbb{R}^n$ ) onto which to project the data so as to minimize the projection error. Reduce from n-dimension to k-dimension: Find k vectors  $\underline{u}^{(1)},\underline{u}^{(2)},\ldots,\underline{u}^{(k)}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

### PCA不是线性回归

左图是线性回归,右图是PCA,线性回归目标是最小化蓝色线段是y值的差值, PCA最小化的是正交距离。

### PCA is not linear regression



线性回归试图预测变量y,PCA中我们同等对待一系列数据集中的x

# PCA(主成分分析方法)降维 (2)

### 数据预处理:

均值标准化: 用每个值减去均值得到的数来代替每个数, 这样使得每个数据样本

的均值为0

特征缩放: (数据样本特征值-平均值) / 标准偏差

#### Data preprocessing

Training set:  $x^{(1)}, x^{(2)}, \dots, x^{(m)} \leftarrow$ 

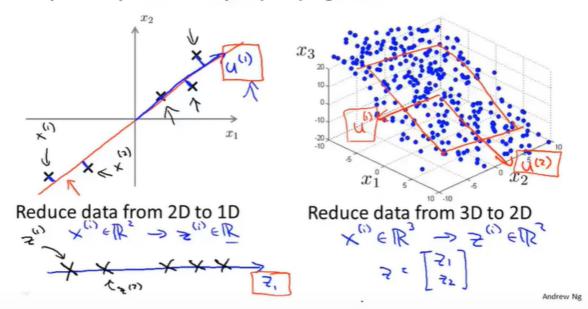
Preprocessing (feature scaling/mean normalization):

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

 $\mu_j = \frac{1}{m} \sum_{i=1}^m x_j^{(i)}$  Replace each  $x_j^{(i)}$  with  $x_j - \mu_j$ . If different features on different scales (e.g.,  $x_1 =$  size of house,  $x_2 =$  number of bedrooms), scale features to have comparable range of values.

### 降维:

#### Principal Component Analysis (PCA) algorithm



### PCA分析:

协方差矩阵: sigma符号, 用来表示一个矩阵(看起来很像求和符号)

svd 奇异值分解

#### Principal Component Analysis (PCA) algorithm

Reduce data from n-dimensions to k-dimensions

Compute "covariance matrix":

Compute "eigenvectors" of matrix 
$$\Sigma$$
:

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

$$= \sum_{i=1}^{n} \sum_{k=1}^{n} (x^{(i)})(x^{(i)})^{T}$$

$$= \sum_{i=1}^{n} \sum_{k=1}^{n} (x^{(i)})(x^{(i)})^{T}$$

$$= \sum_{k=1}^{n} \sum_{k=1}^{n} (x^{(i)})(x^{(i)})^{T$$

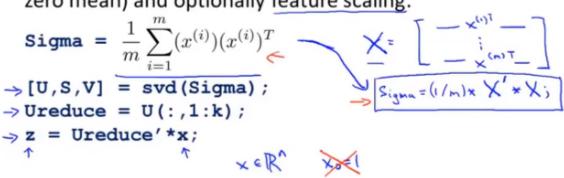
取U矩阵(n\*n)的前k列

U矩阵(n\*n)的前k列的转置,再乘以n\*1维的x(i)向量,即可得到k维的向量z Principal Component Analysis (PCA) algorithm

From [U,S,V] = svd(Sigma), we get:

### Principal Component Analysis (PCA) algorithm summary

After mean normalization (ensure every feature has zero mean) and optionally feature scaling:



PCA最终得到了最小平方差投影误差