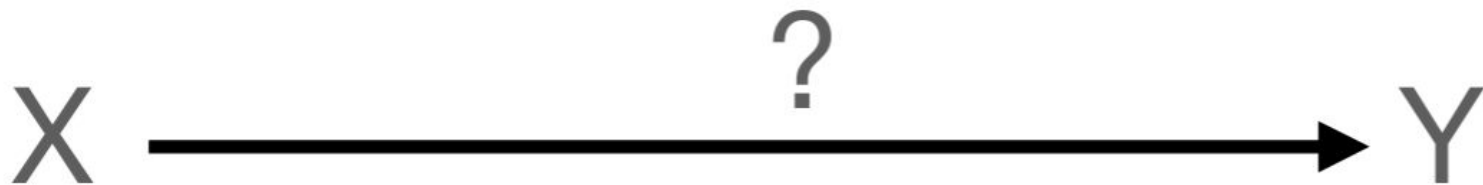


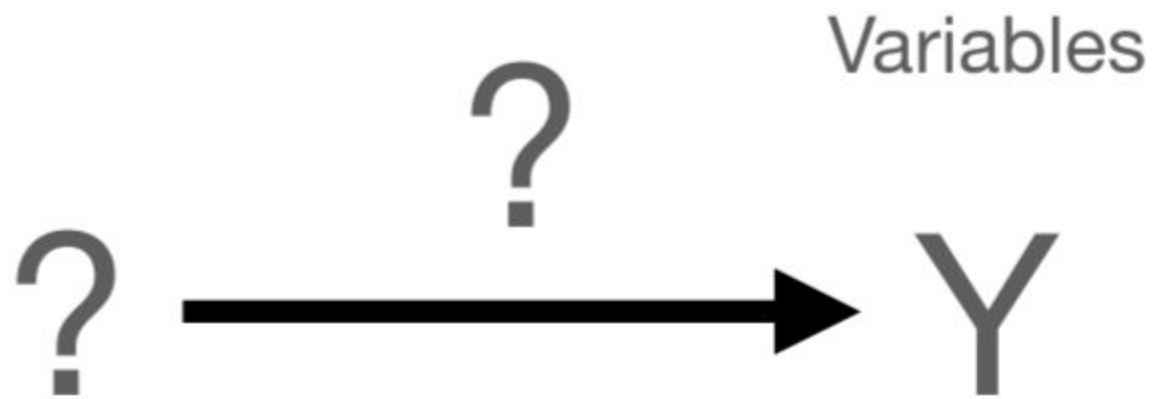
Unsupervised Learning



Predictor variables

Response variables



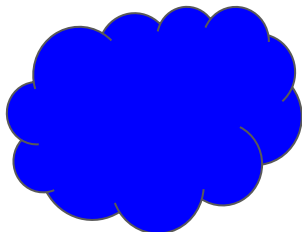
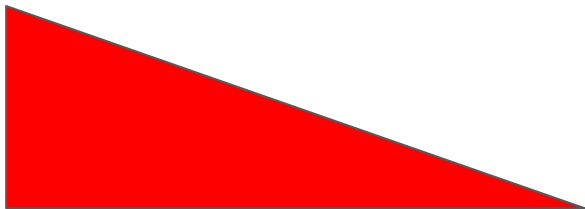
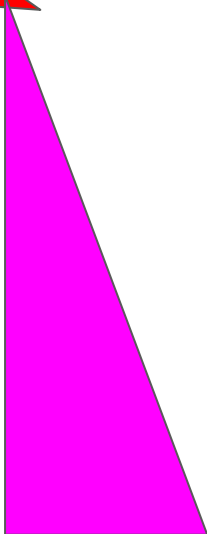
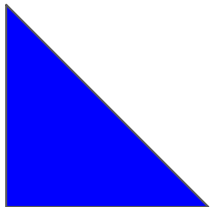
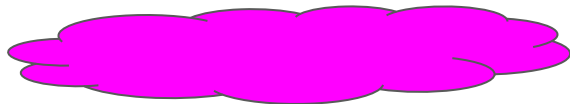
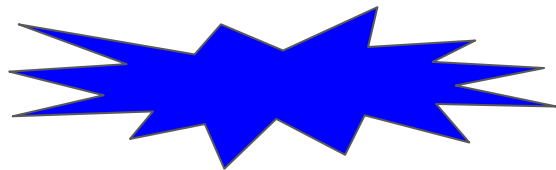
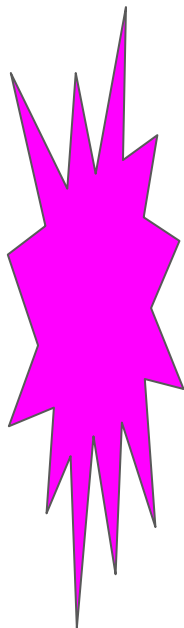
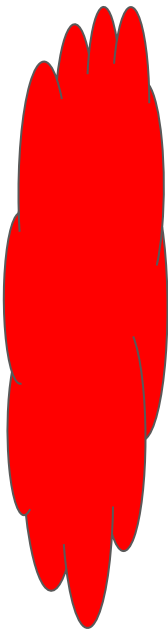
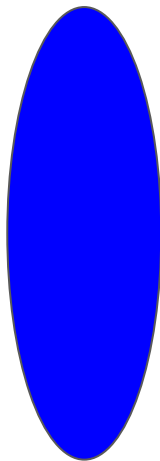
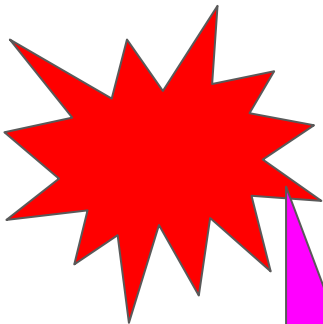
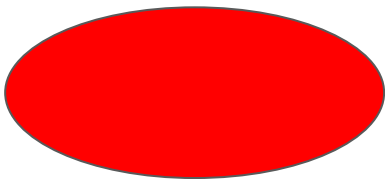
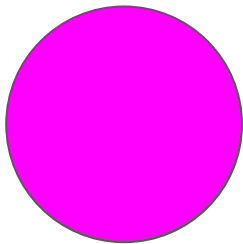


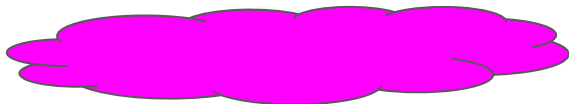
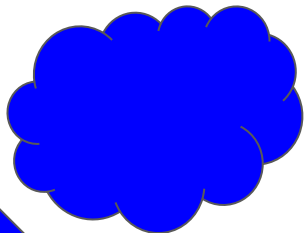
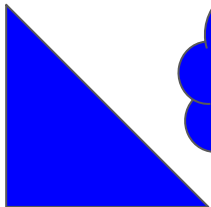
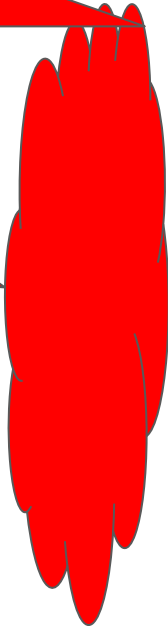
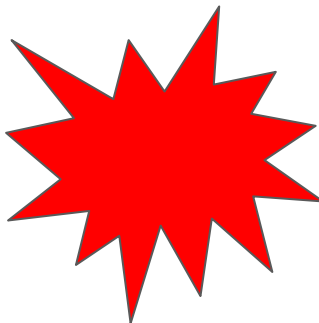
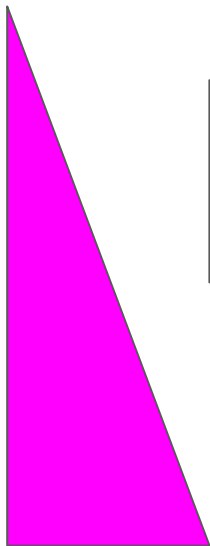
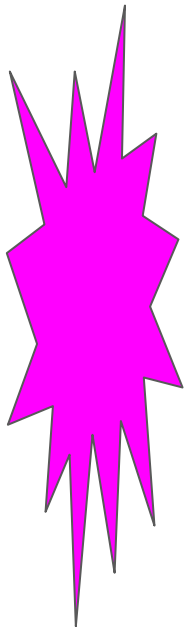
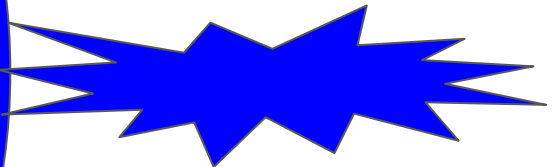
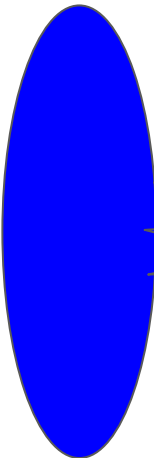
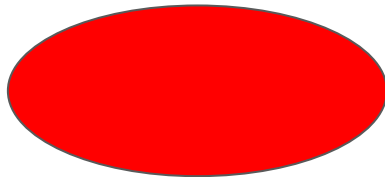
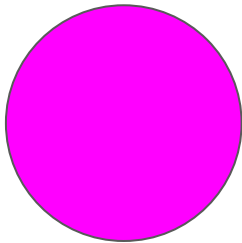
DATA: X

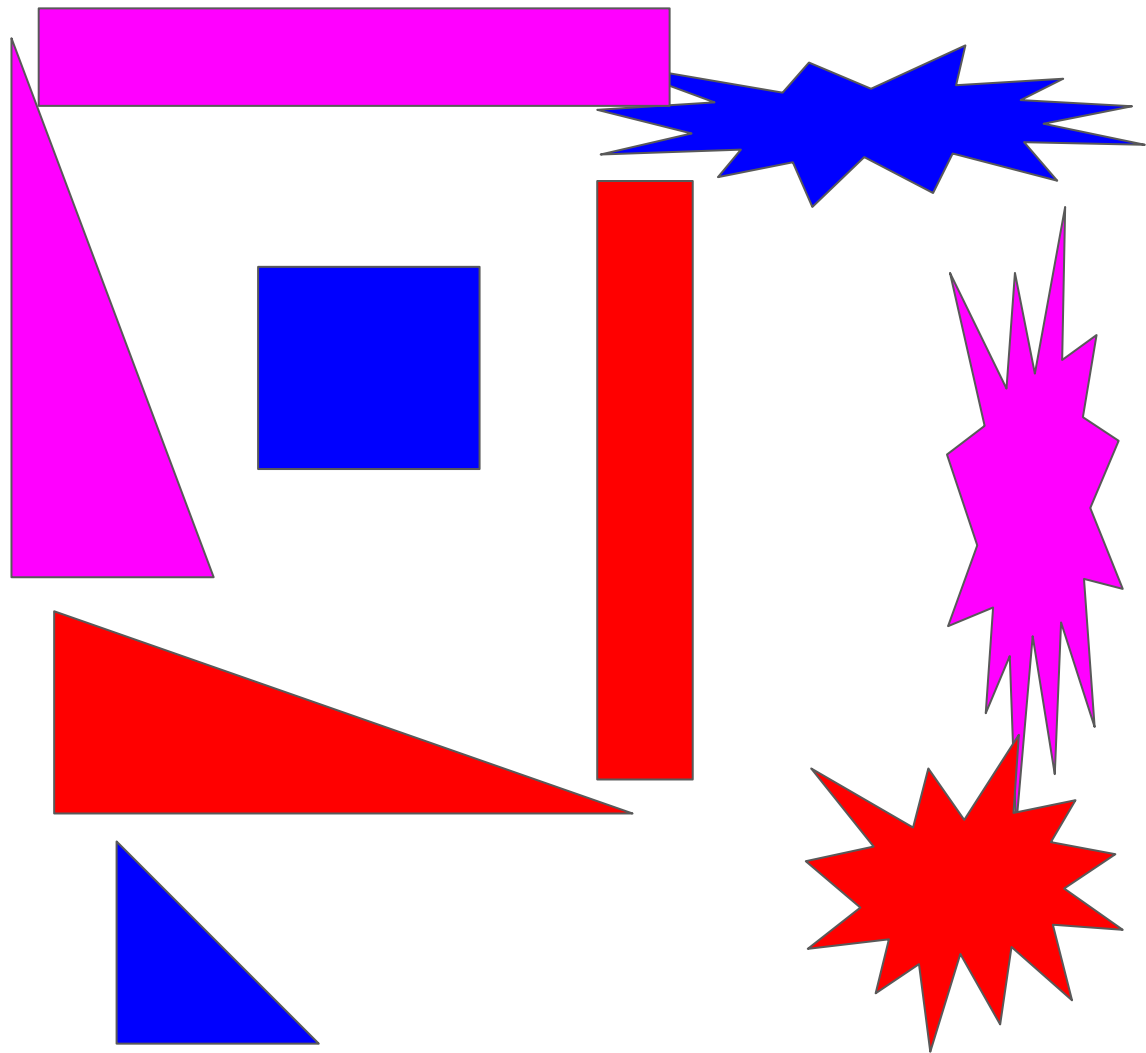
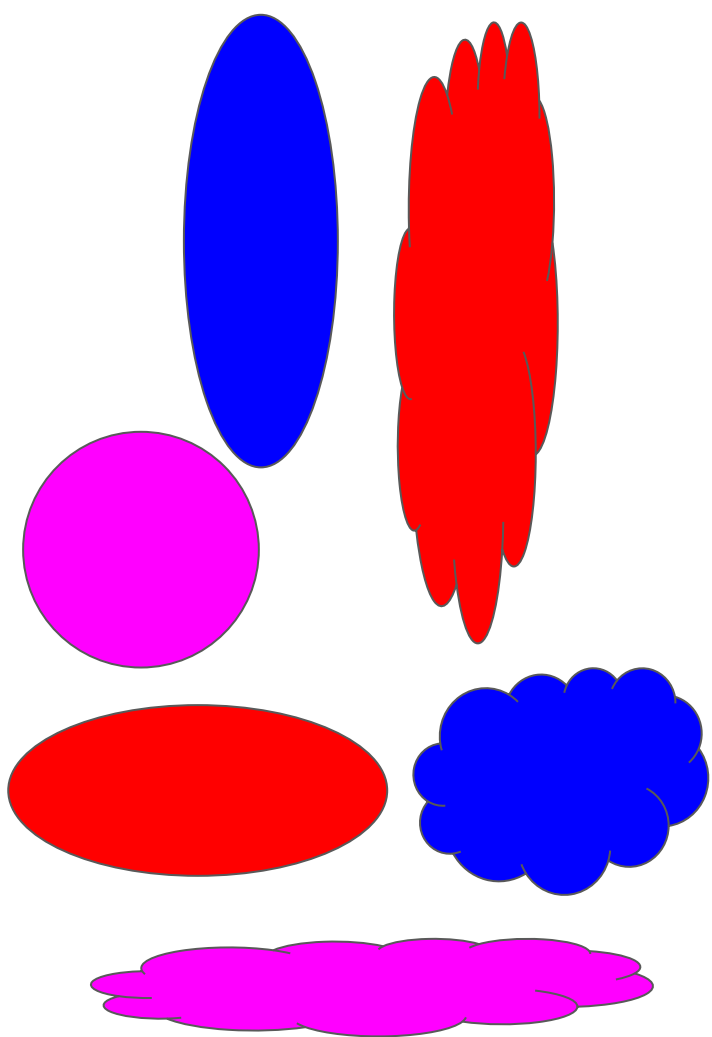
Using
“properties”
of X

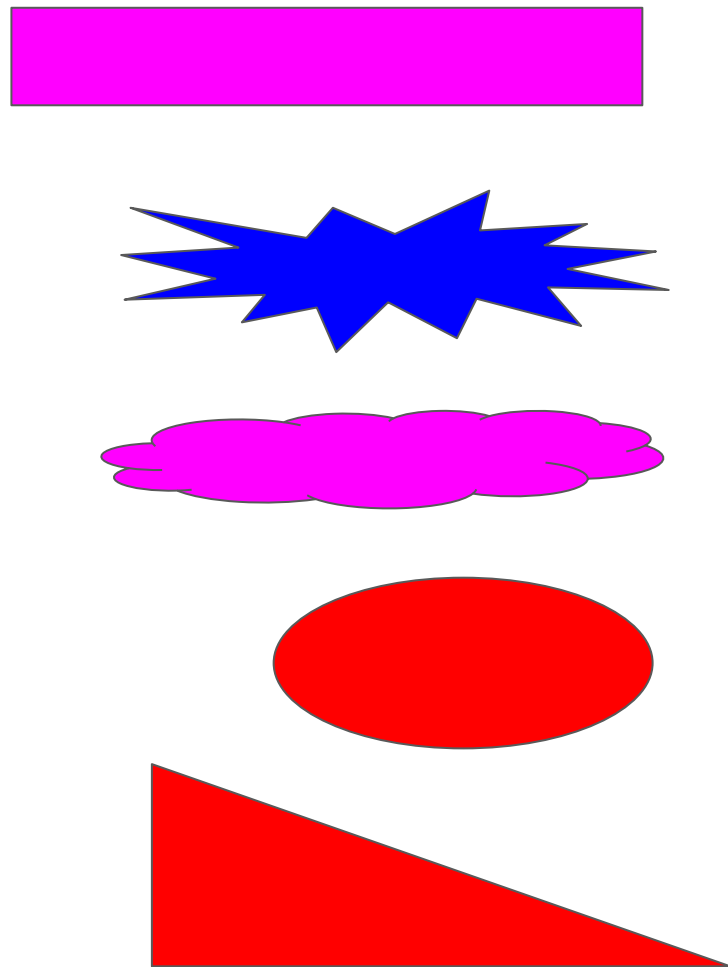
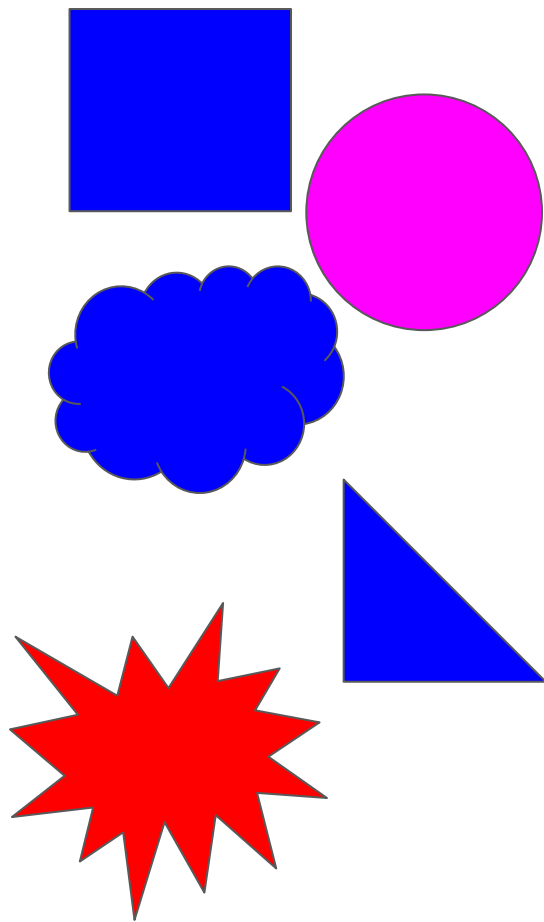
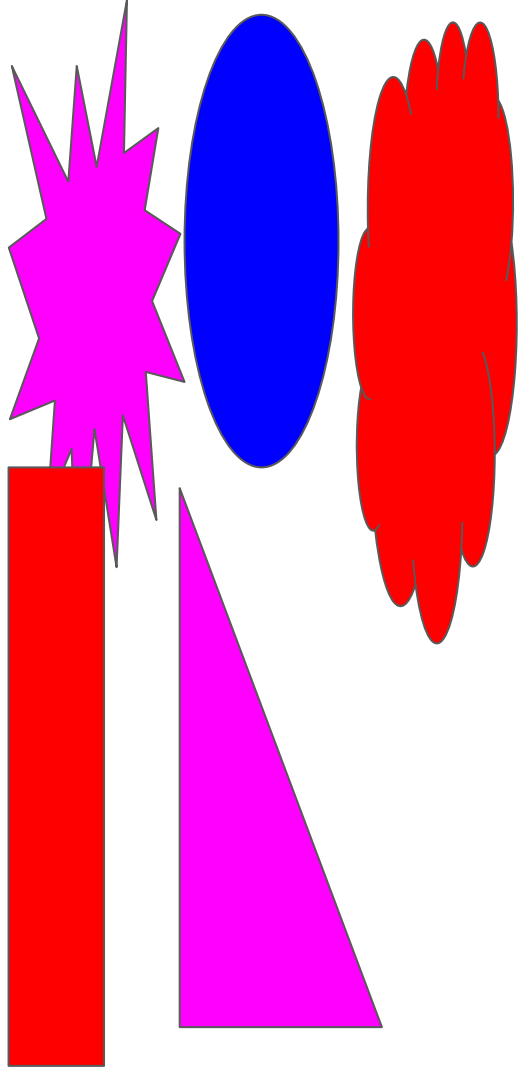
MODEL: $\beta_0, \beta_1, \beta_2, \dots, \beta_n$

Let's start with a dataset









What should these properties be?

- Well defined
 - No ambiguity
- Computable
 - Always has an answer

Data Centric and multidimensional

Variables:	var_1	var_2	var_3	...	var_{n-1}	var_n
Samples:						
sample ₁						
sample ₂						
sample ₃						
...						
sample _{m-1}						
sample _m						

Almost always very multidimensional

$$N_{\text{variables}} \gg M_{\text{samples}}$$

- Peptides in CE-MS
- Genotypes after sequencing
- Species in a given area
- Transcripts in a tissue
 - Transcripts in a cell

Techniques in unsupervised learning

Hodgepodge of nonsense!!!!

Principal Component Analysis

Nearest Neighbours

- K nn
- Fixed-radius nn
- Approximate nn

Cluster Analysis

- Connectivity based CA
- Centroid based
- Grid based

Multidimensional Scaling

- Principal coordinate analysis
- Metric
- Non metric
- Generalized multidimensional scaling

Neural Network Autoencoders

- Classic
- Variational
- Denoising

Mixture Models

- Gaussian
- Categorical
- Multivariate

Neighbour Embedding

- TSNE
- UMAP

What we will do today

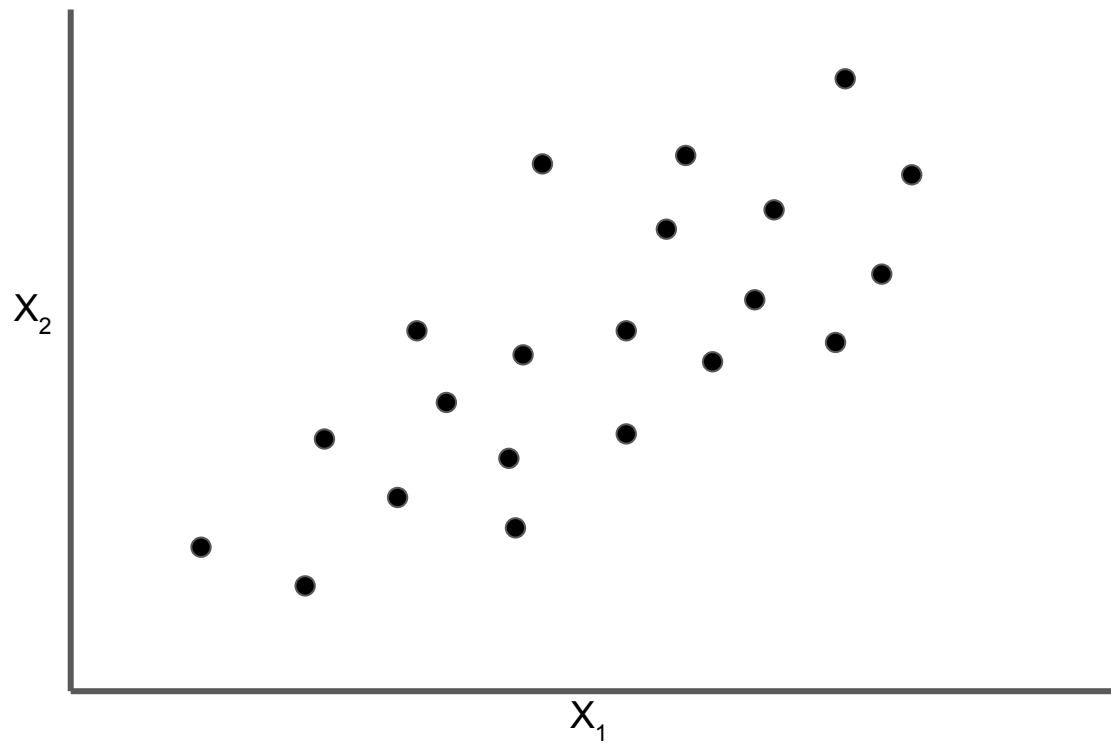
One technique - PCA

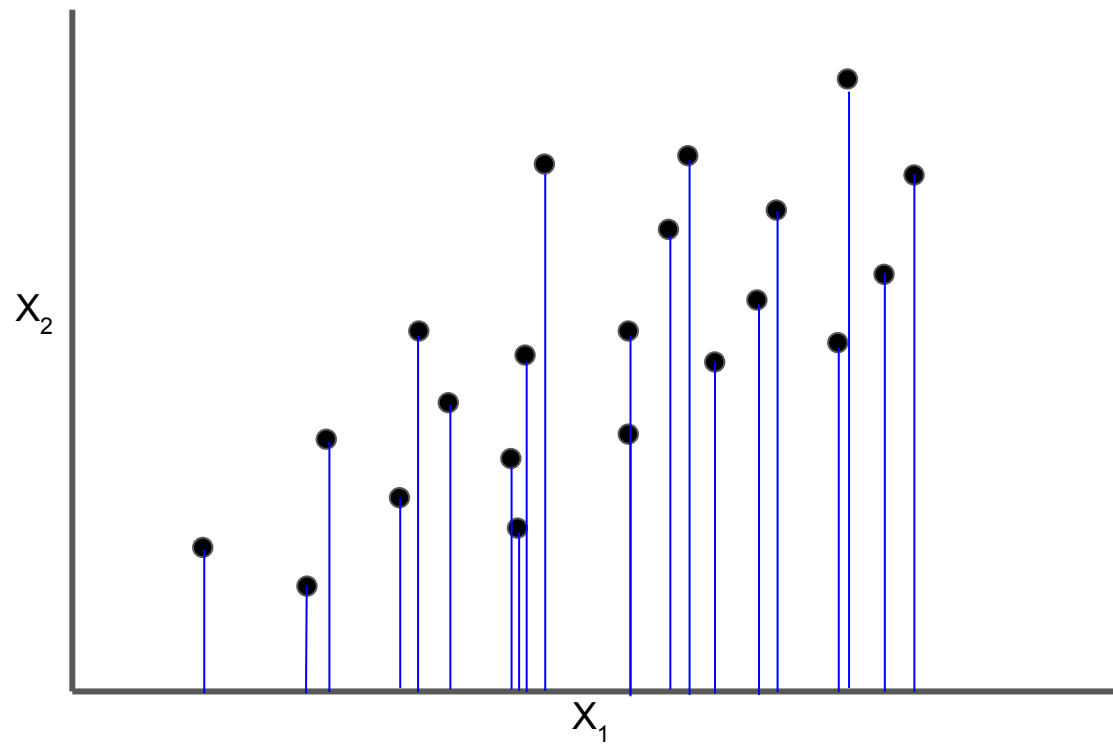
- Cover it broadly, without brushing off the important details.
- Set the groundwork for other techniques.
- Not learn how to “compute” a PCA, but understand its principles.

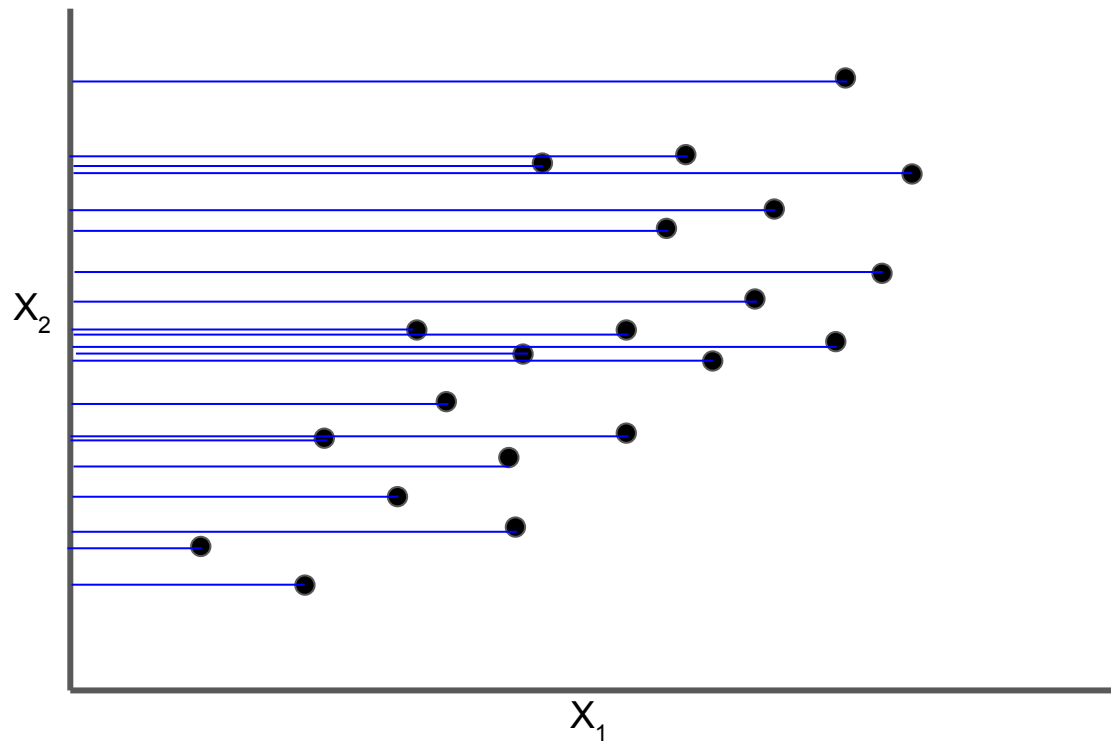
Data X

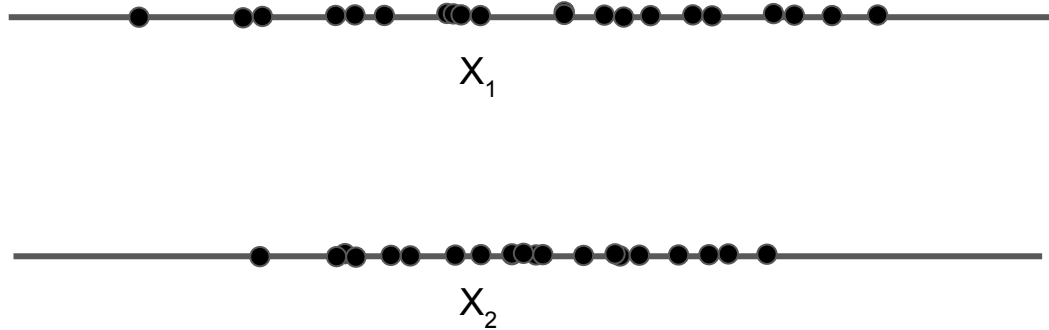
	\mathbf{x}_1	\mathbf{x}_2
sample ₁		
sample ₂		
sample ₃		
...		
sample _{m-1}		
sample _m		

We have some data $X \in \mathbb{R}^2$





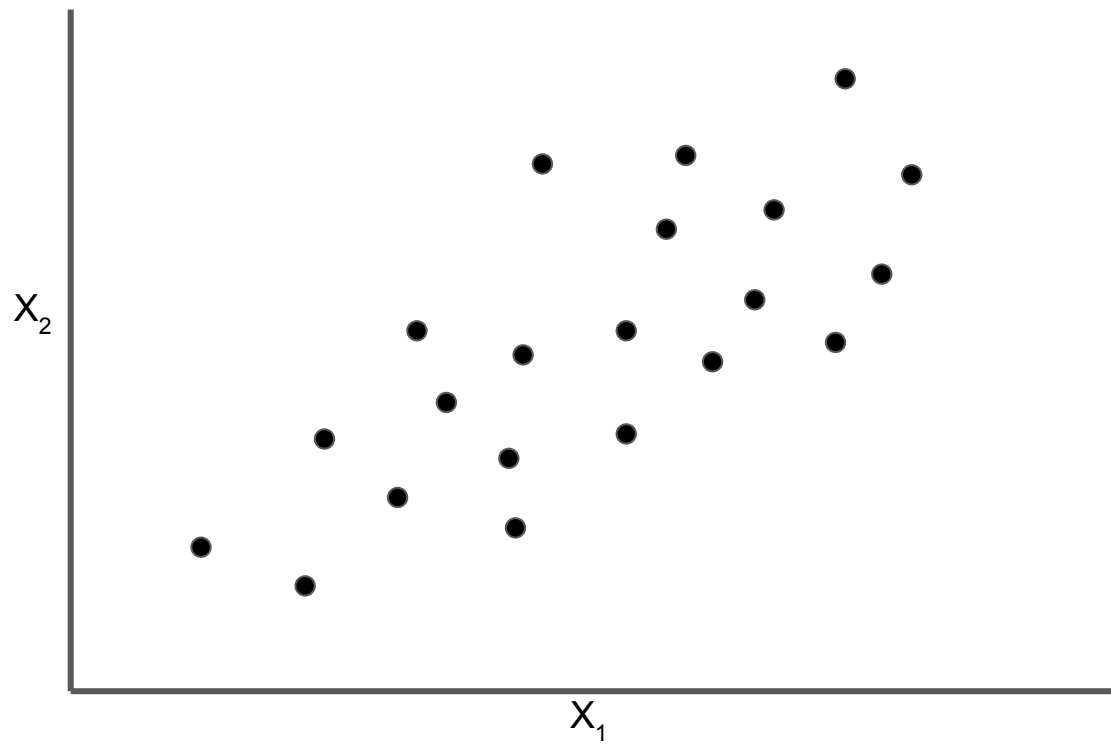




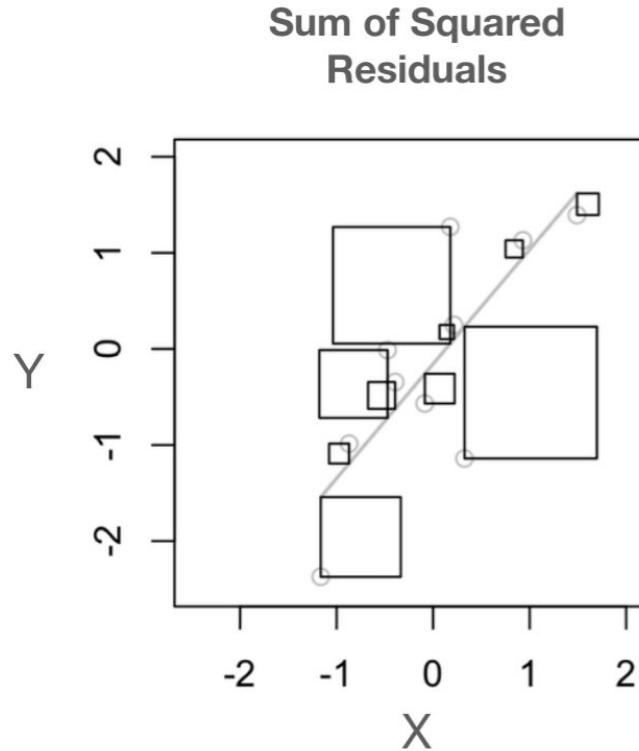
$$s^2 = \frac{\sum (x_i - \bar{x})^2}{n - 1}$$

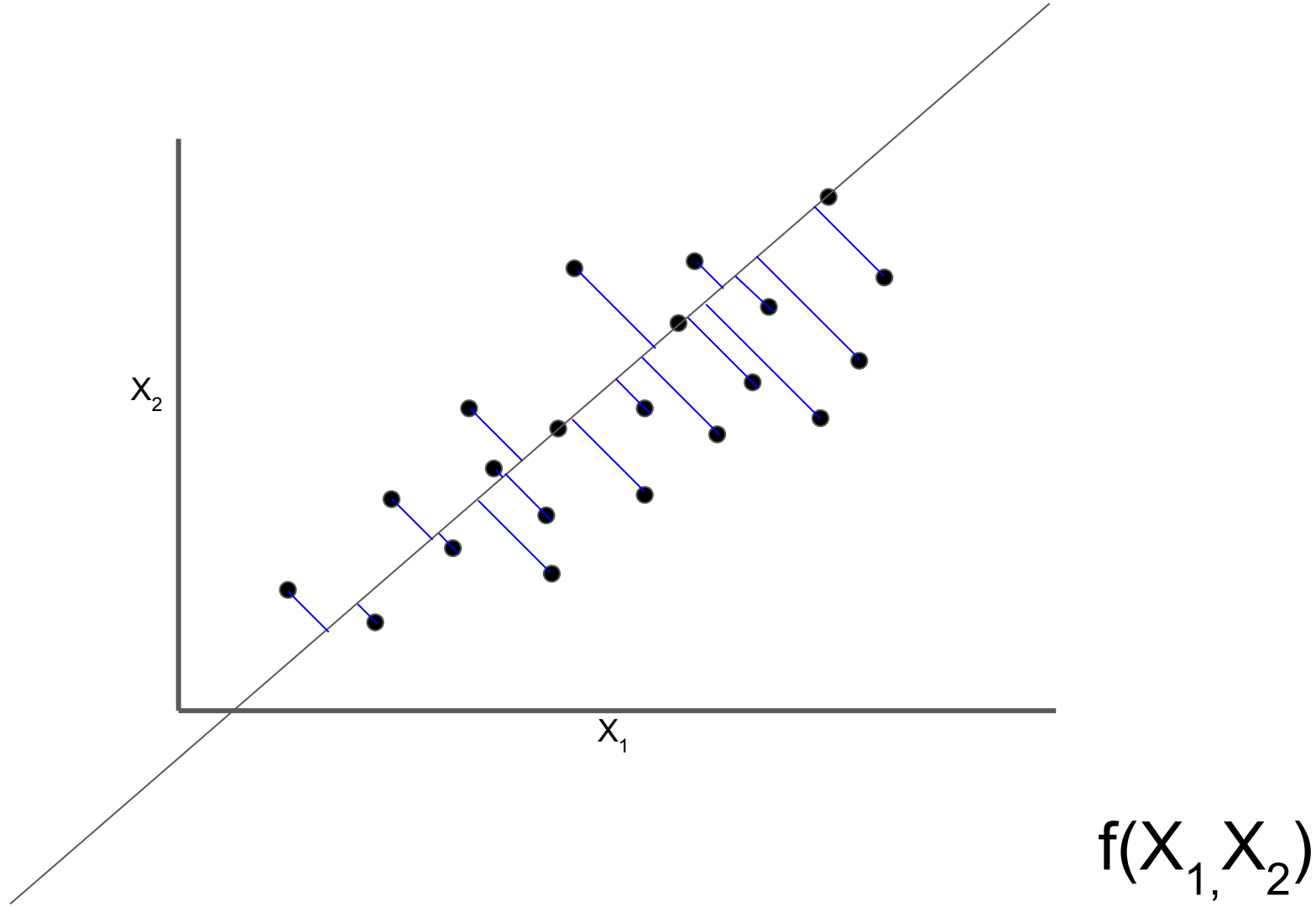
var(x1) and var(x2)

Total variance in the sample:
 $\text{var}(x_1) + \text{var}(x_2)$



Correlated variables

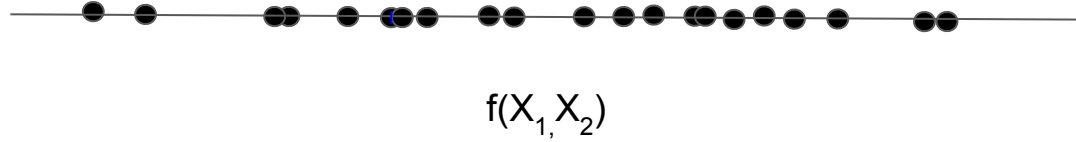
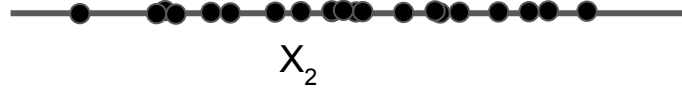
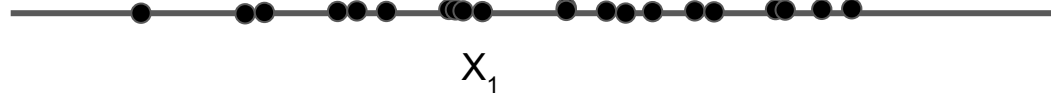




$f(X_1, X_2)$ must be a linear function

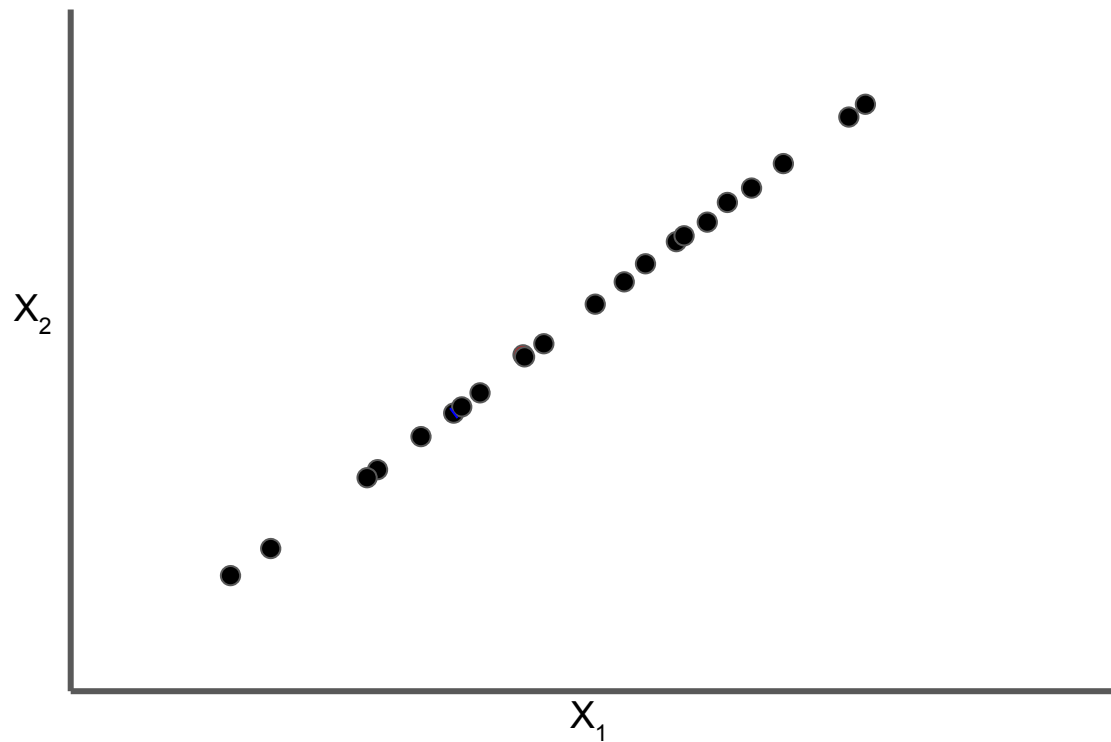
Addition

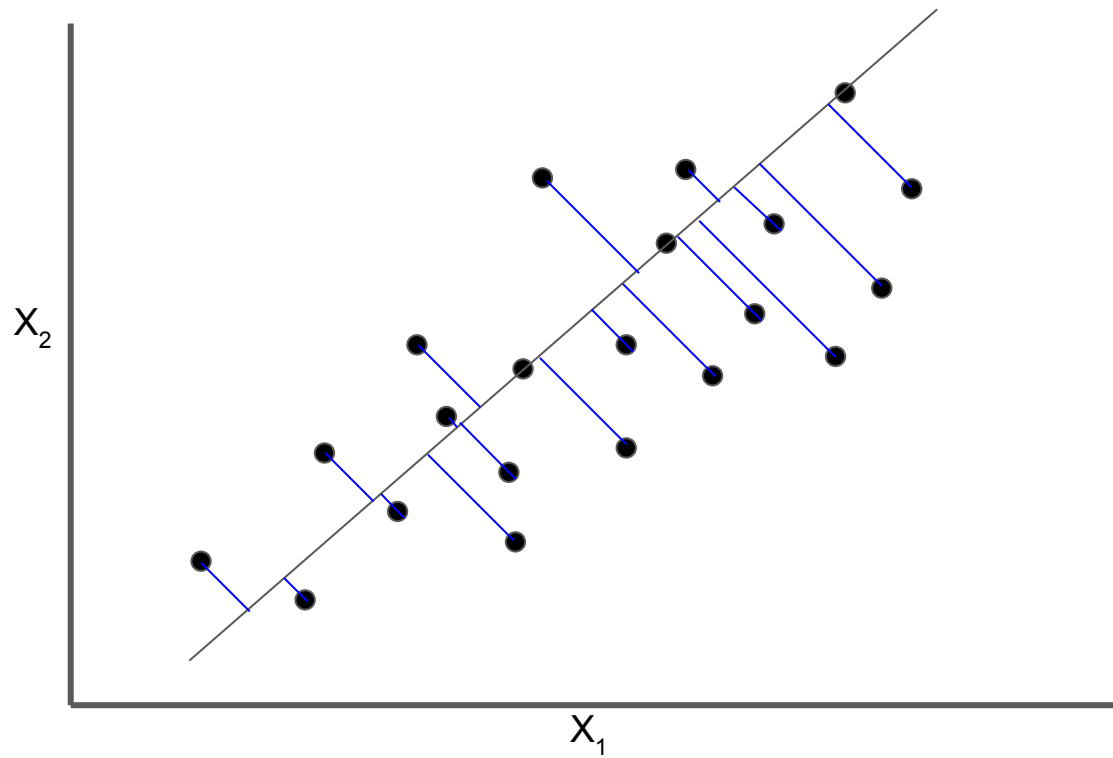
Multiplication



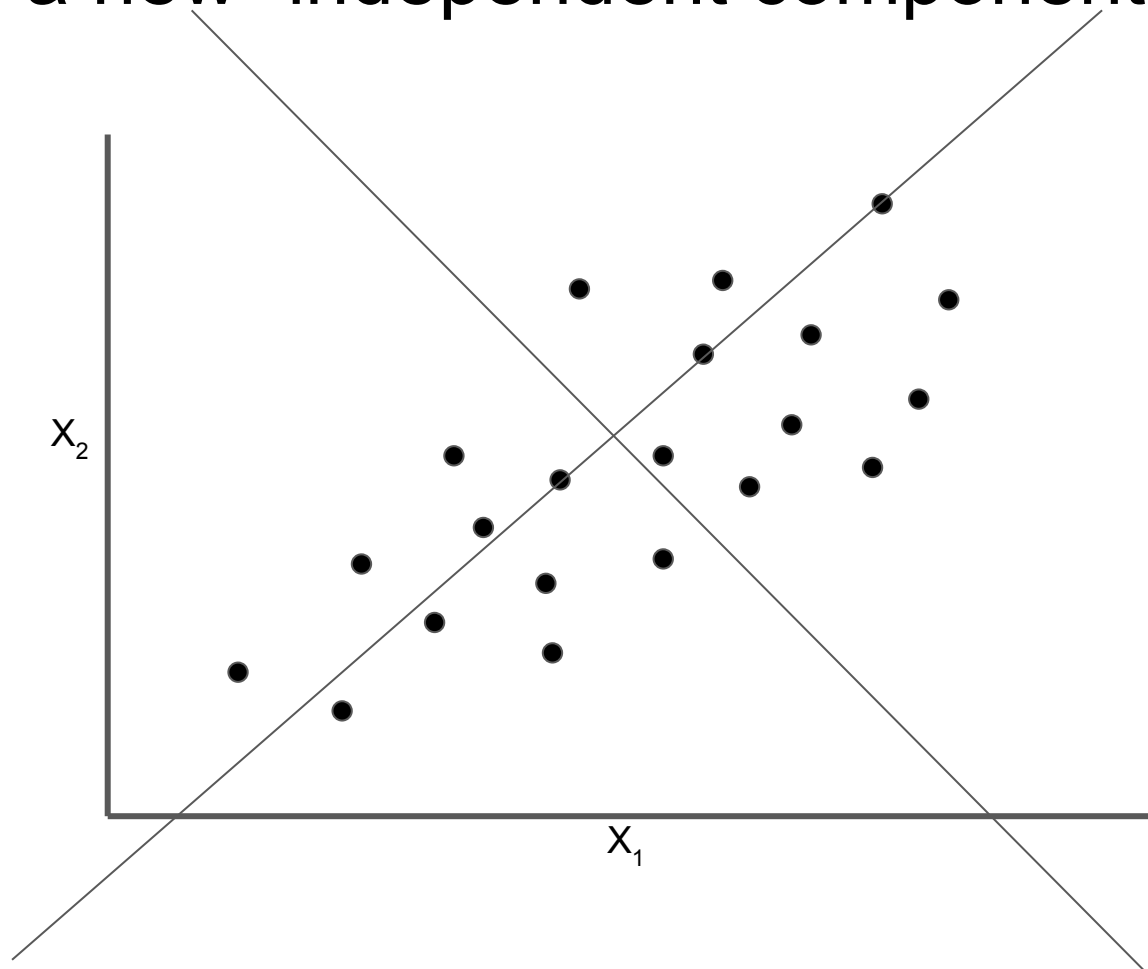
$$\text{var}(f(X_1, X_2)) > \text{var}(x_1) \text{ or } \text{var}(x_2)$$

We have missing information

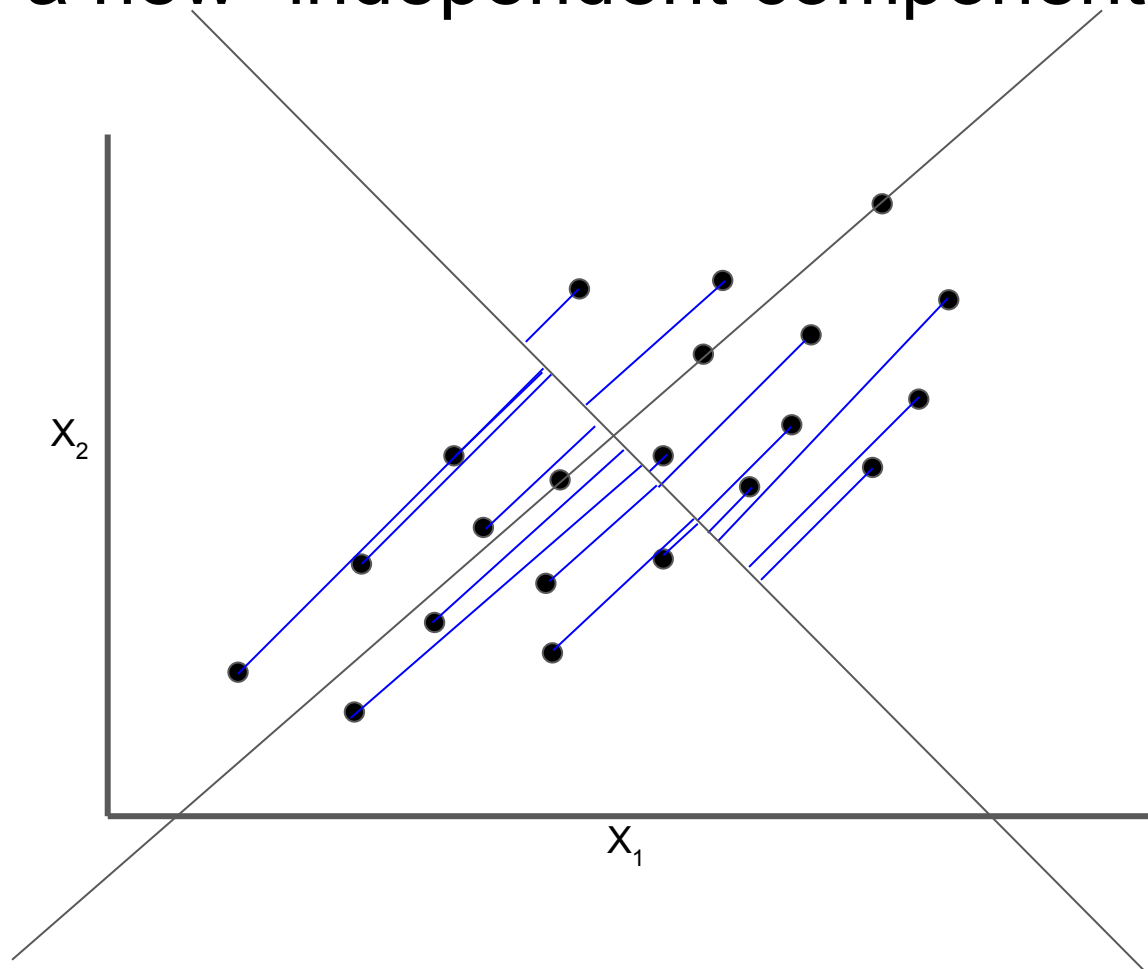




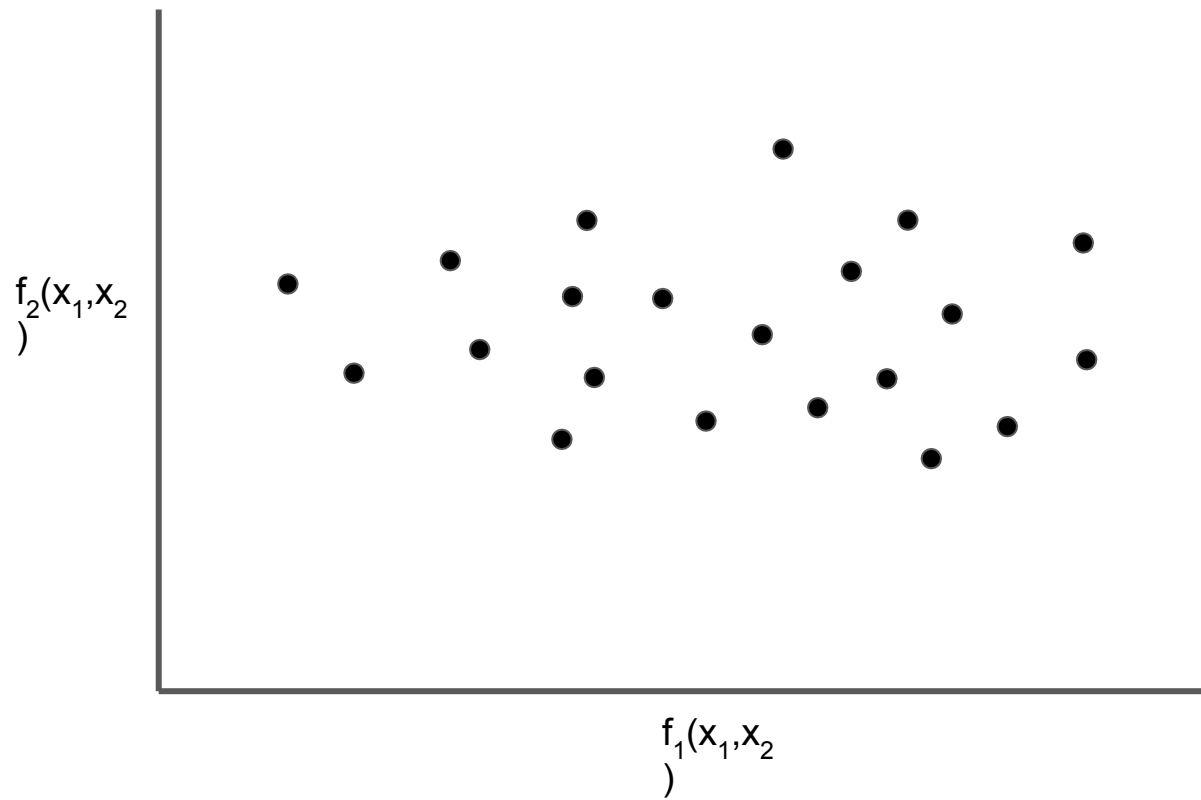
We need a new “independent component”

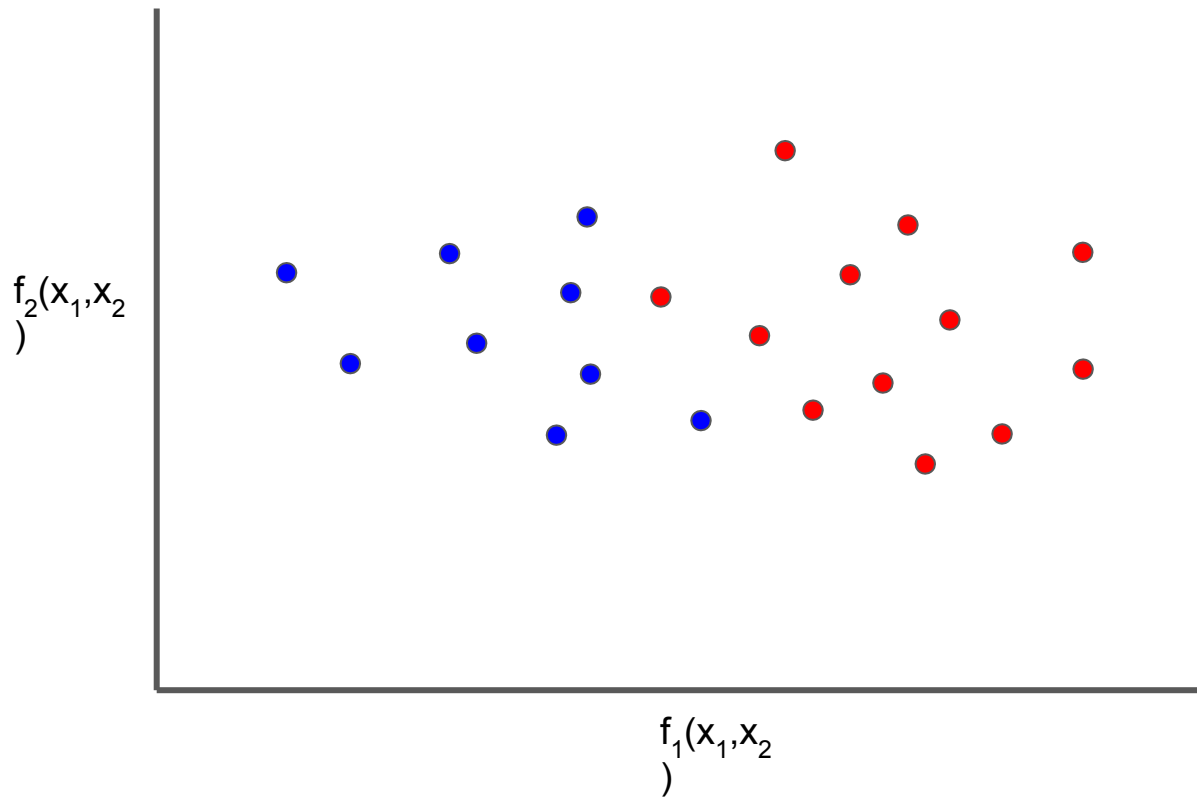


We need a new “independent component”



$$\text{var}(f_1(X_1, X_2)) + \text{var}(f_2(X_1, X_2)) == \text{var}(x1) + \text{var}(x2)$$



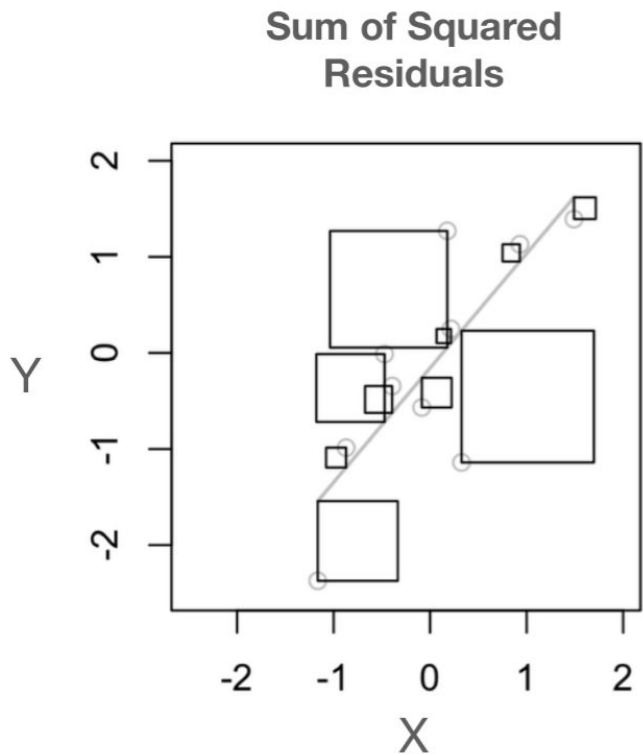


$$\text{var}(f_1(X_1, X_2)) + \text{var}(f_2(X_1, X_2)) == \text{var}(x_1) + \text{var}(x_2)$$

Two dimensions are cool, but we have bigger problems than that.

Let's define the “constraints”

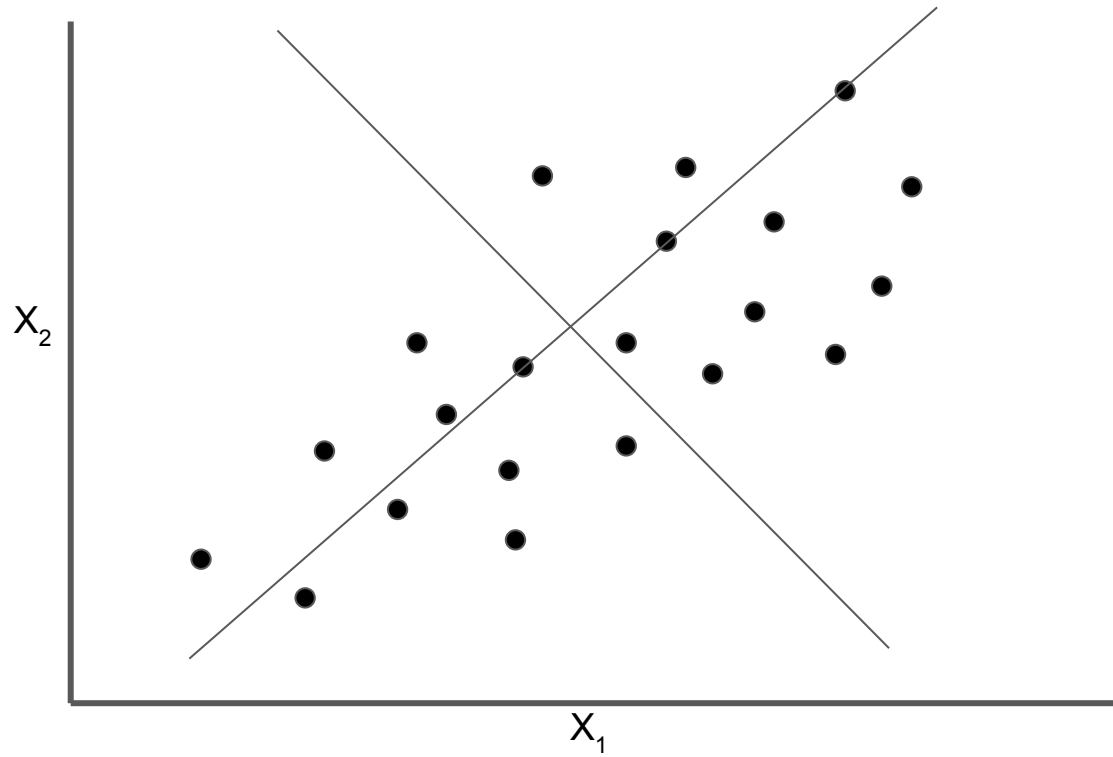
Constrains



$$SS_{res} = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

but our estimate \hat{y}_i is simply
a linear function of x_i :

$$\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$$



Constrains

Maximum variance

Linear independence

How to compute?

Nothing in linear algebra makes sense except in the light of eigen-analysis.

Matrices as linear transformation

The covariance matrix

	var_1	var_2	var_3	...	var_{n-1}	var_n
var_1	$S^2(\text{var}_1)$					
var_2	$\text{COV}(\text{var}_1, \text{var}_2)$	$S^2(\text{var}_2)$				
var_3	$\text{COV}(\text{var}_1, \text{var}_3)$		$S^2(\text{var}_3)$			
...		
var_{n-1}	$\text{COV}(\text{var}_1, \text{var}_{n-1})$				$S^2(\text{var}_{n-1})$	
var_n	$\text{COV}(\text{var}_1, \text{var}_n)$					$S^2(\text{var}_n)$

