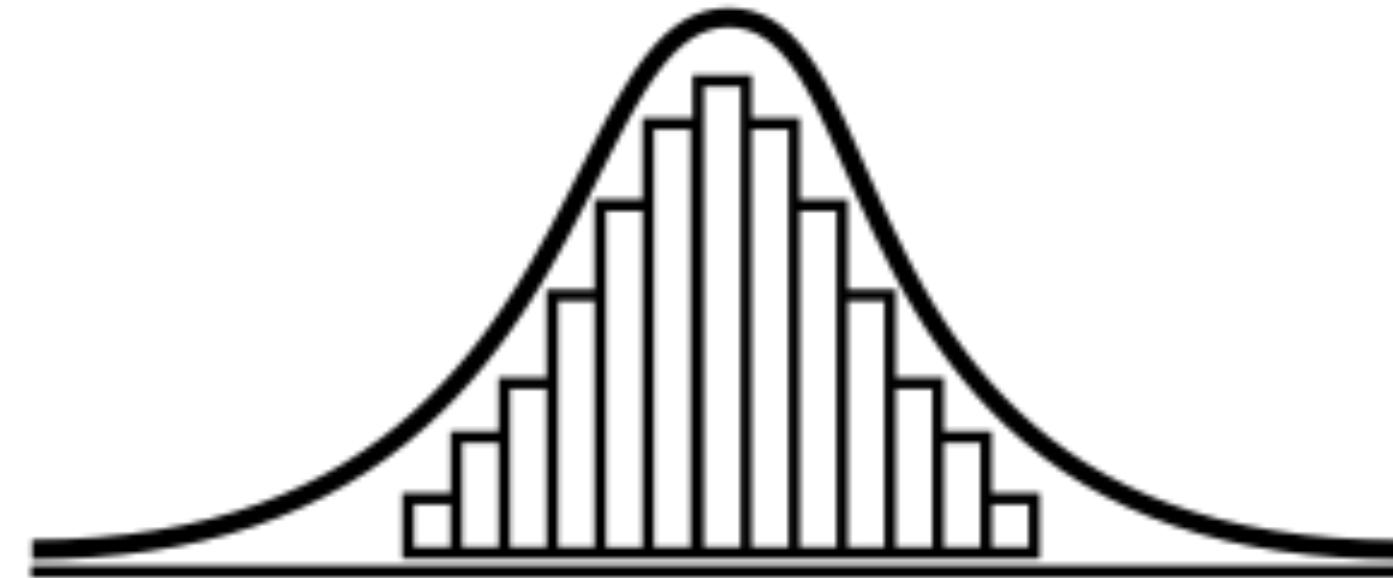


Gaussian Progress

(an utmost normal topic)



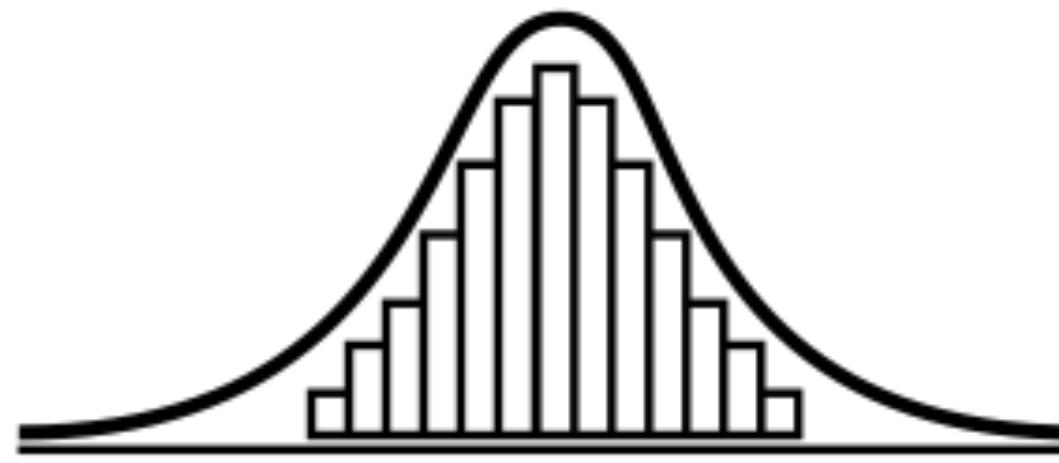
a homage to gaussian methods, tricks and progress



When you're a machine learning professional you might feel like you need to learn so many algorithms that it can be hard to keep up. It can be very demotivating.



This talk is not about downplaying this feeling but it is about demonstrating a lovely hack; understanding a mother algorithm.



It turns out that if you appreciate what the gaussian distribution can do then there are lot's of algorithms that are much easier to grasp. This talk is an attempt at explaining the power of the Gaussian[tm] by stepping up the ladder of complexity of algorithms.

Today

1. Literally Introduce Gauss
2. Explain *the gaussian* trick for
 - Classification, Clusters and Outliers
3. Explain normal neural benefits
 - Mixed Density Networks/Mixture Encoders
4. Live Code a Gaussian Process
 - A really meta trick

The problem

Teacher wants to read a newspaper so he comes up with a mathematical challenge for the kids to keep 'em busy.

So the teacher said ...

"Add all the numbers 1 till 100."

... and then continued reading his newspaper.

This is the task:

$$1 + 2 + \dots + 99 + 100 = ?$$

This is the task:

$$1 + 2 + \dots + 99 + 100 = ?$$

But here's a rewrite.

$$\begin{array}{r} 1 + 2 + \dots + 49 + 50 \\ 100 + 99 + \dots + 52 + 51 \\ \hline 101 + 101 + \dots + 101 + 101 = 101 * 50 = 5050 \end{array}$$

This is the task:

$$1 + 2 + \dots + 99 + 100 = ?$$

But this task is similar.

$$\begin{array}{r} 1 + 2 + \dots + 99 + 100 \\ 100 + 99 + \dots + 2 + 1 \\ \hline \end{array}$$

$$101 + 101 + \dots + 101 + 101 \rightarrow 101 * 100/2 = 5050$$

This is the task:

$$1 + 2 + \dots + (n-1) + n = ?$$

This is the task:

$$1 + 2 + \dots + (n-1) + n = ?$$

This is what math is!

$$\begin{array}{r} 1 \quad \quad + \quad \quad 2 \quad + \quad \dots \quad + \quad (n-1) \quad + \quad n \\ n \quad \quad + \quad (n-1) \quad + \quad \dots \quad + \quad \quad 2 \quad + \quad 1 \\ \hline \end{array} +$$
$$(n+1) + (n+1) + \dots + (n+1) + (n+1) \rightarrow (n+1) * n/2$$

Math is a compiler for numbers.

But math is best explained with stories.

The story I just told is more commonly told via this formula:

$$\sum_{i=1}^n i = \frac{n(n + 1)}{2}$$

I personally think the story of the boy who bested his teacher is a better way of explaining things. Especially since this story is something that *actually* happened.

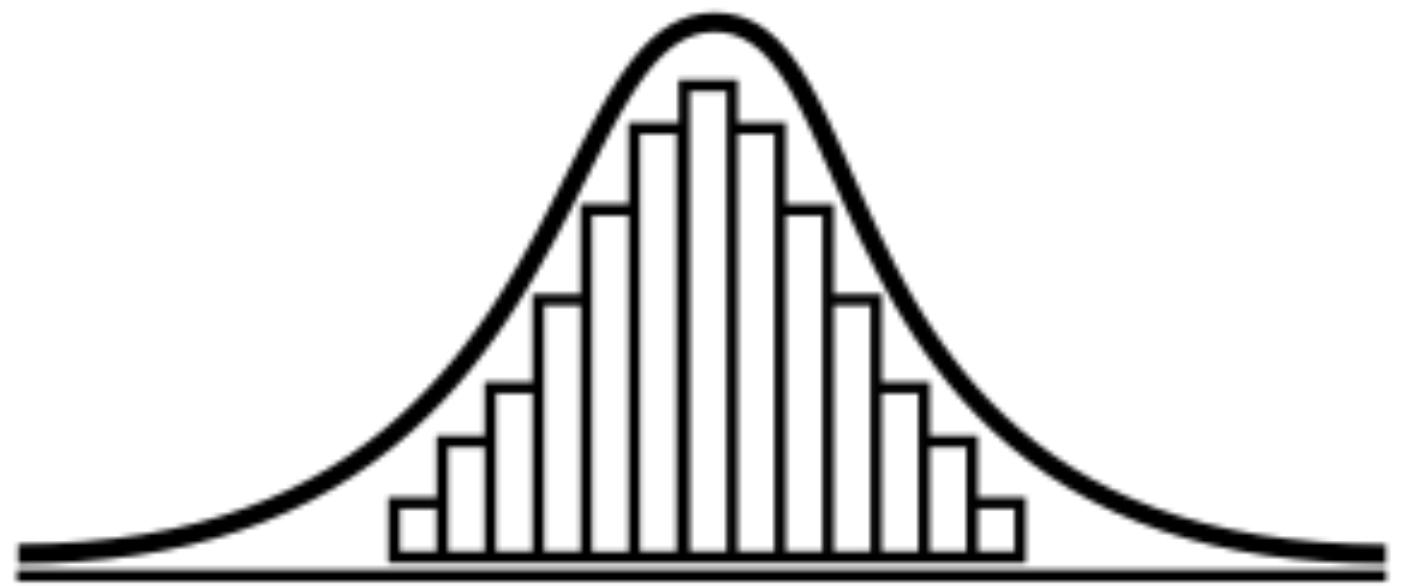
It supposedly actually happened to a kid named "Carl Friedrich Gauss".

And something very popular got named after him;

- the gaussian distribution
- the normal distribution
- the bell curve

These are all the same thing.



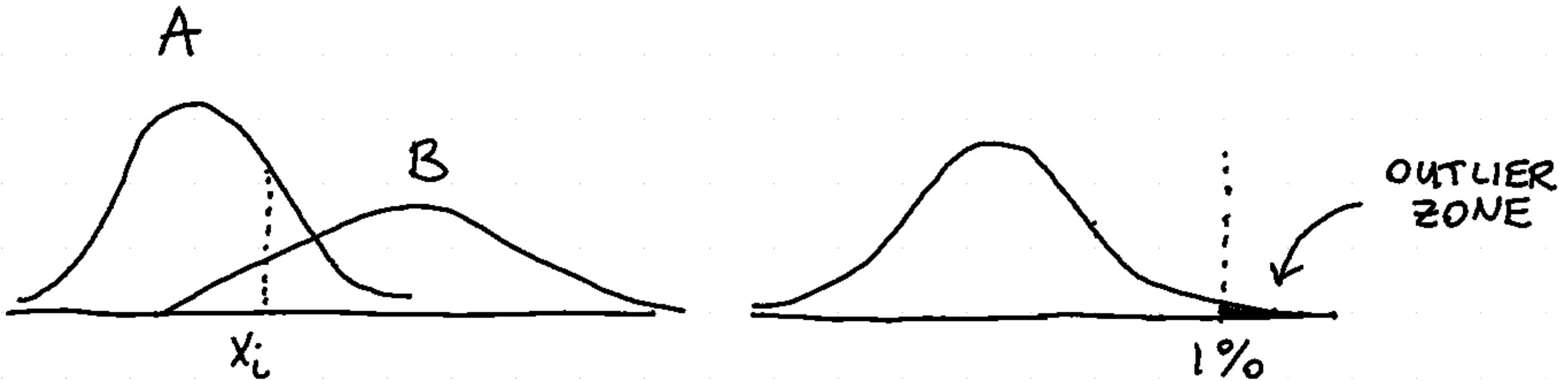


Gauss As A Building Block

This is the maths in one dimension.

$$f(x|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

But here's some usecases:



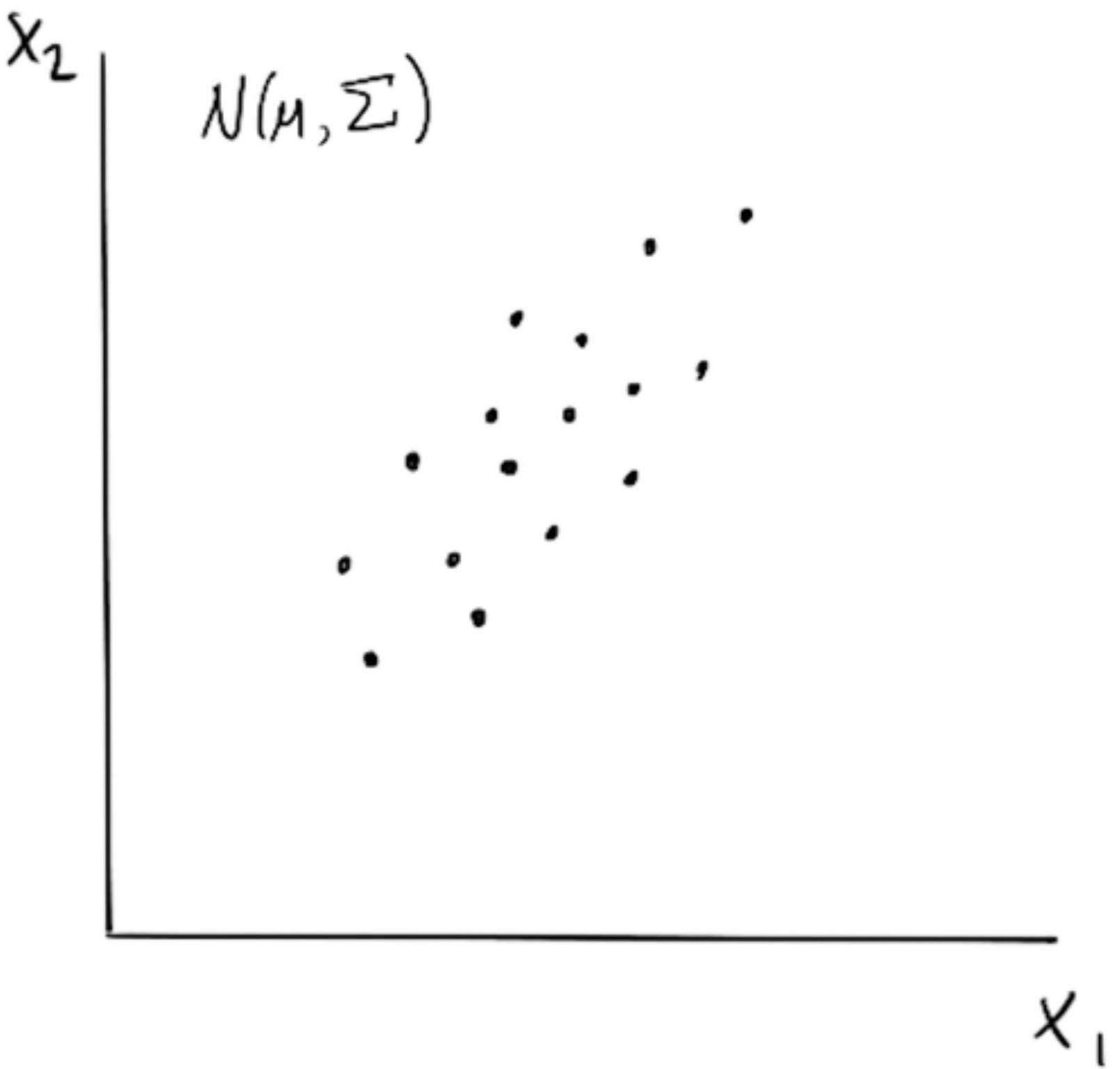
$$p(x_i = A) = \frac{f_A(x_i)}{f_A(x_i) + f_B(x_i)}$$

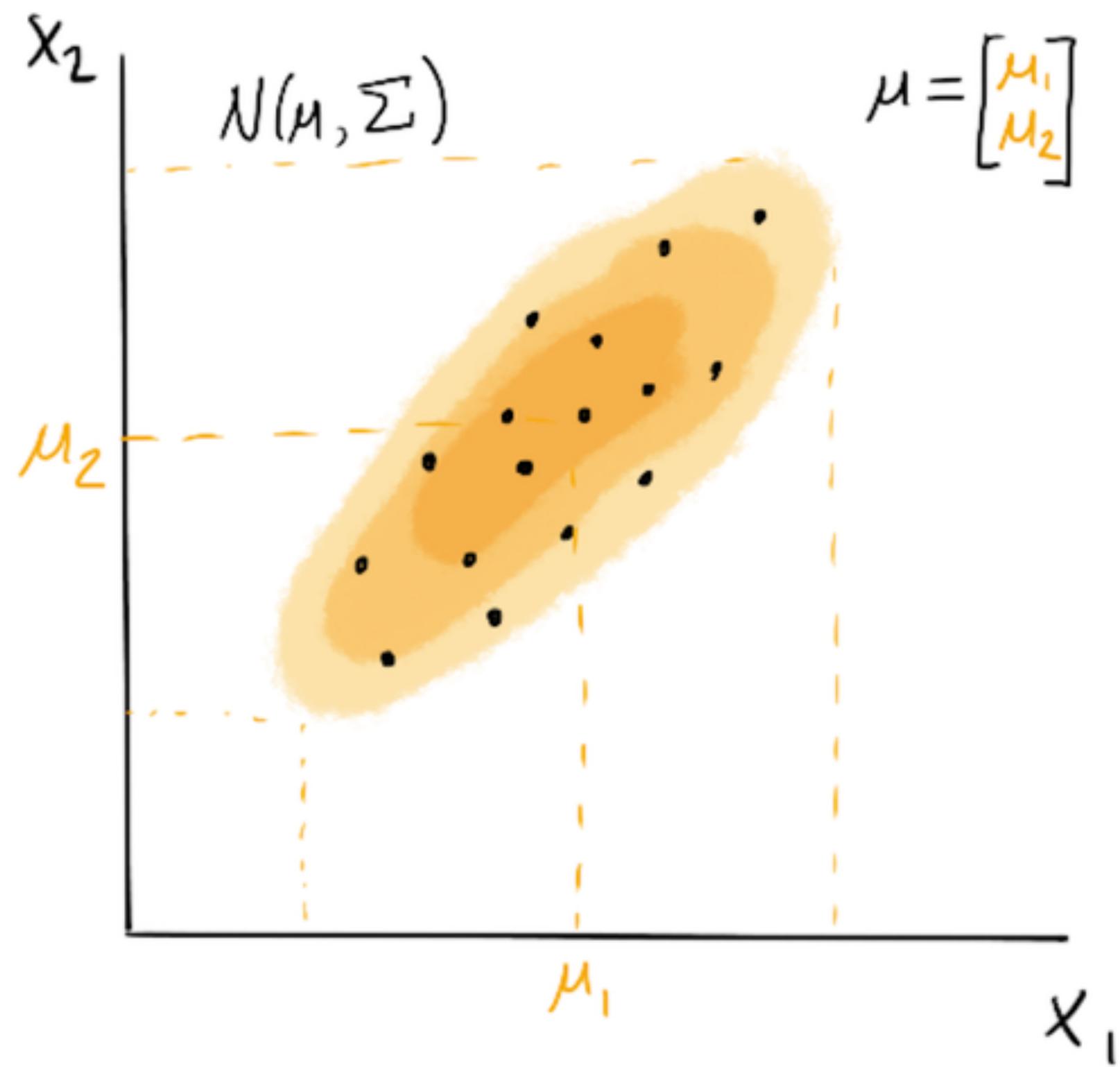
This is the maths in one dimension.

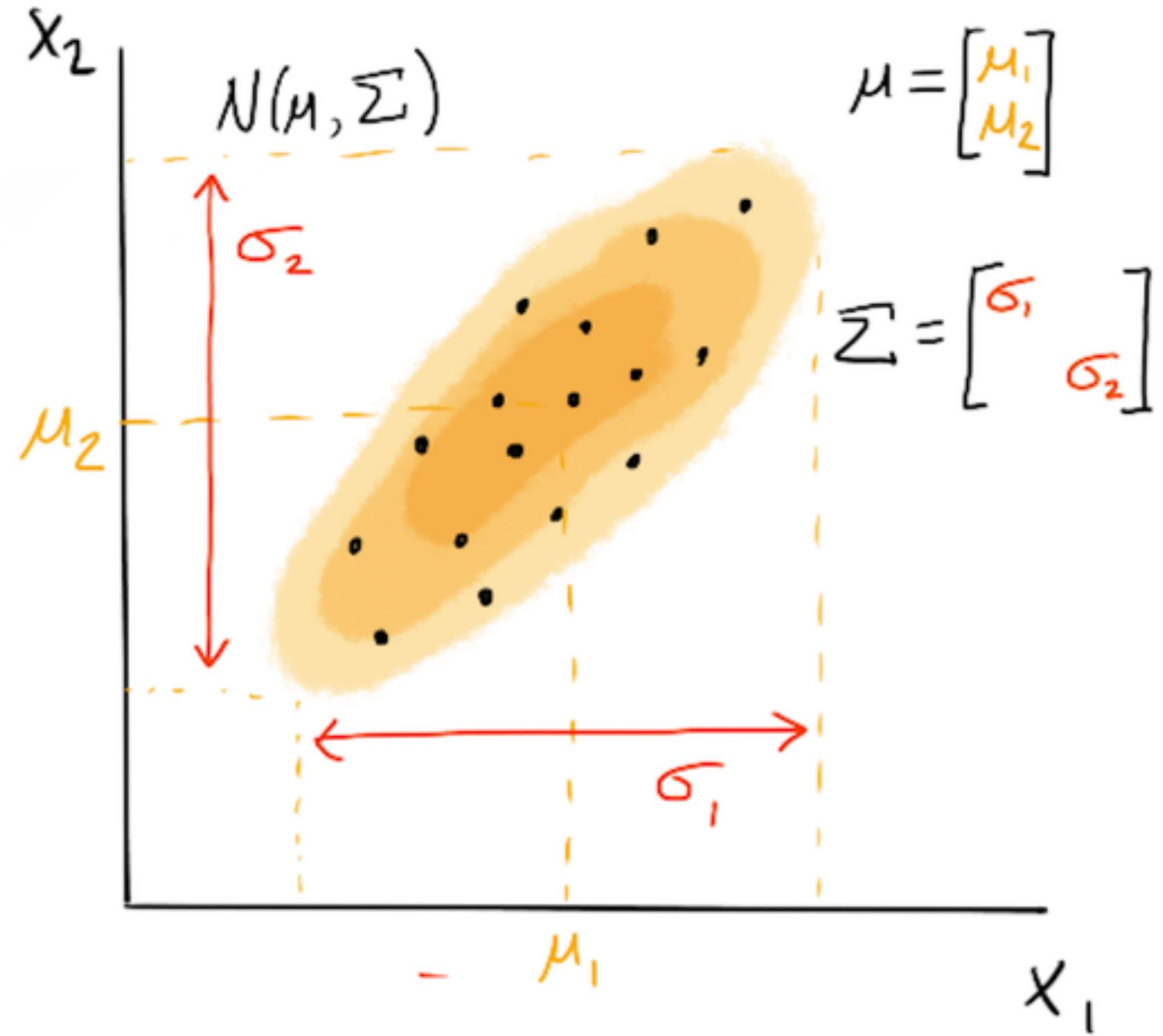
$$f(x|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

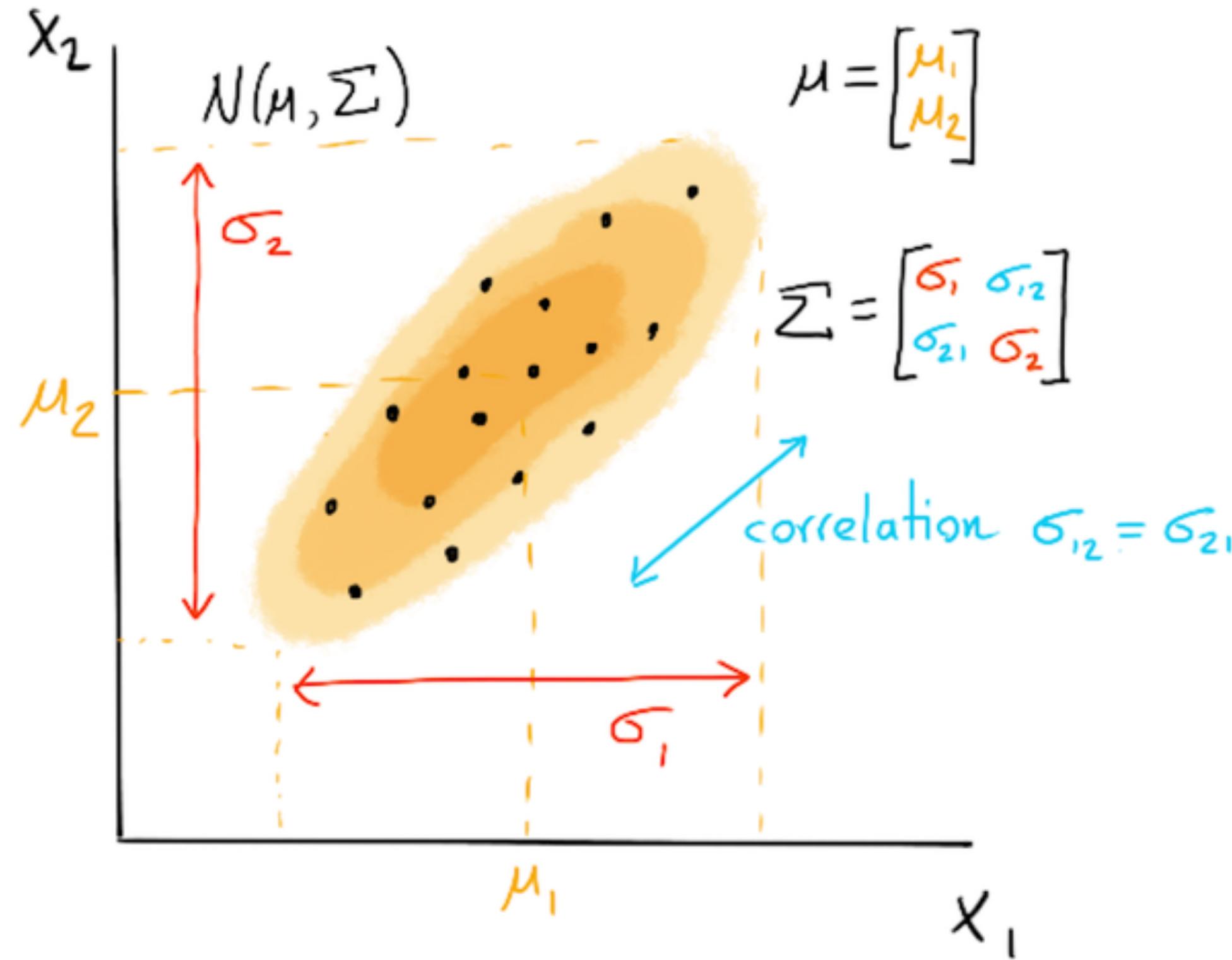
This is the formula for 2+ dimensions.

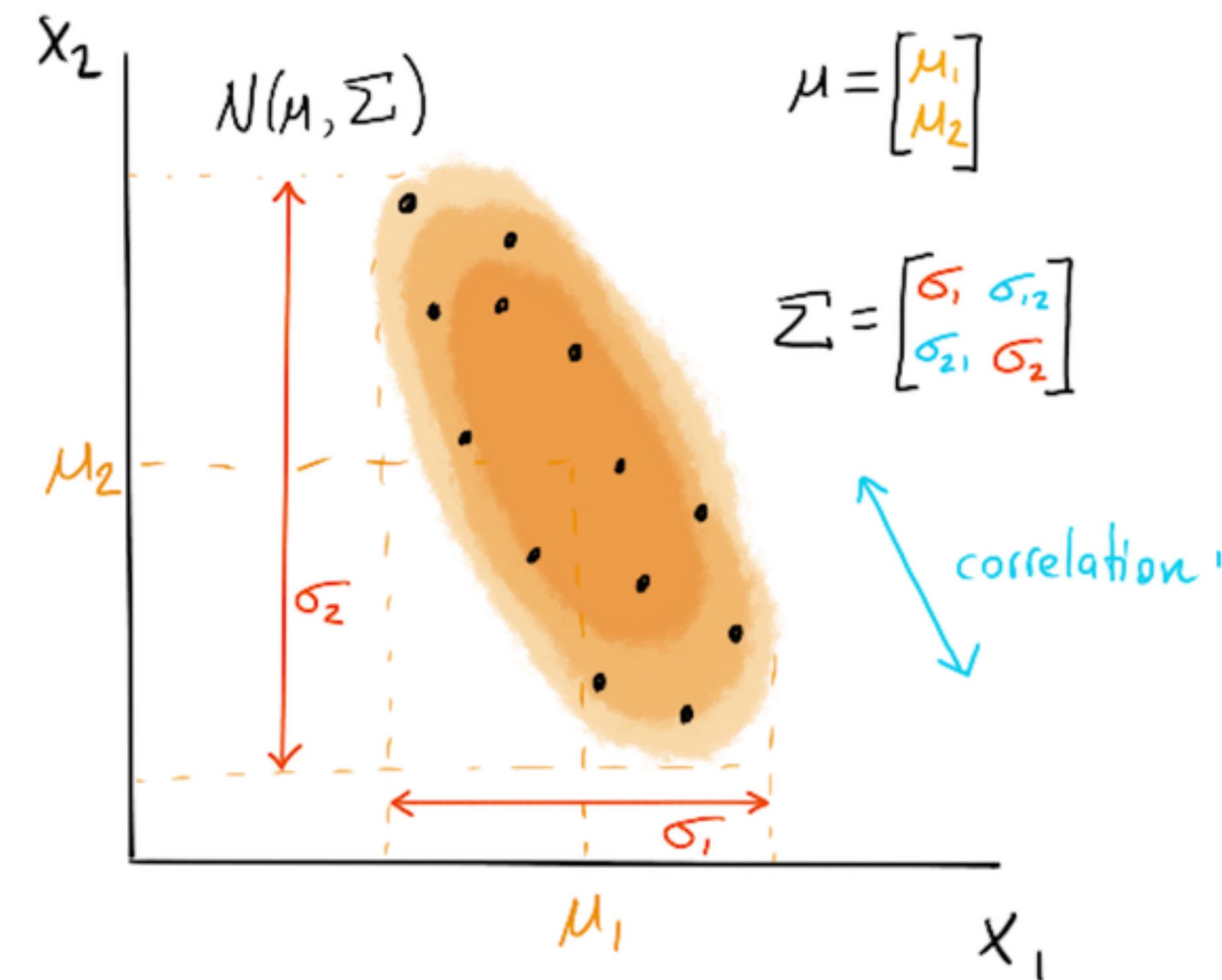
$$f_{\mathbf{x}}(x_1, \dots, x_k) = \frac{\exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})\right)}{\sqrt{(2\pi)^k |\boldsymbol{\Sigma}|}}$$

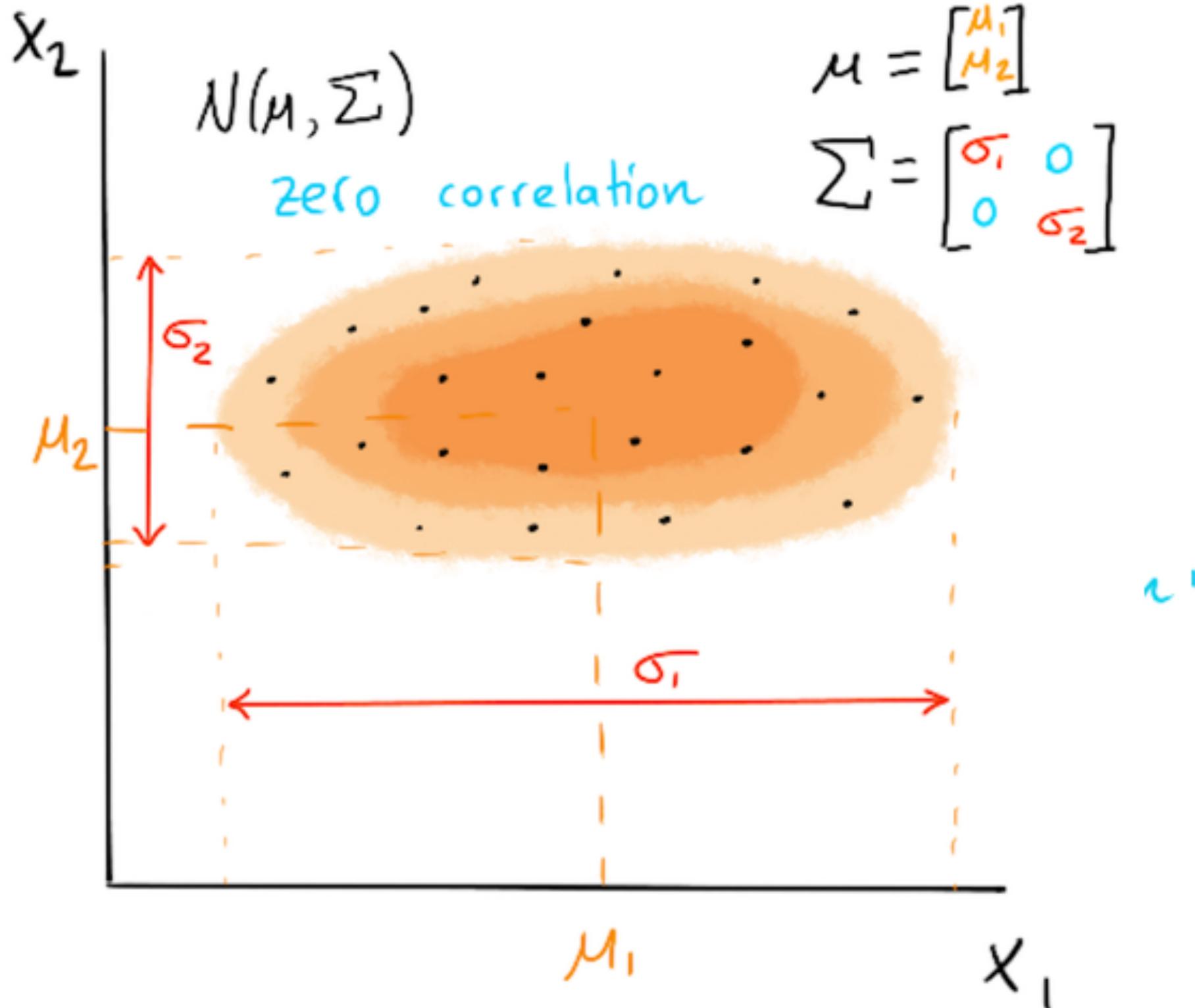






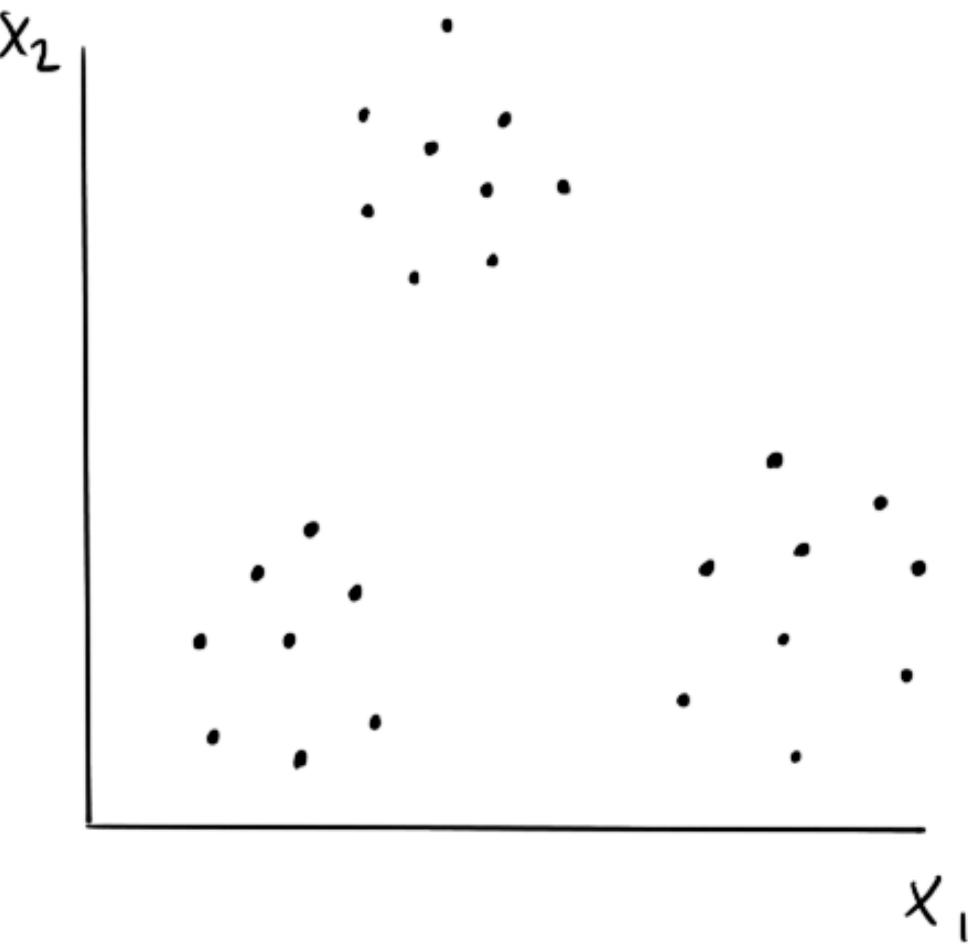


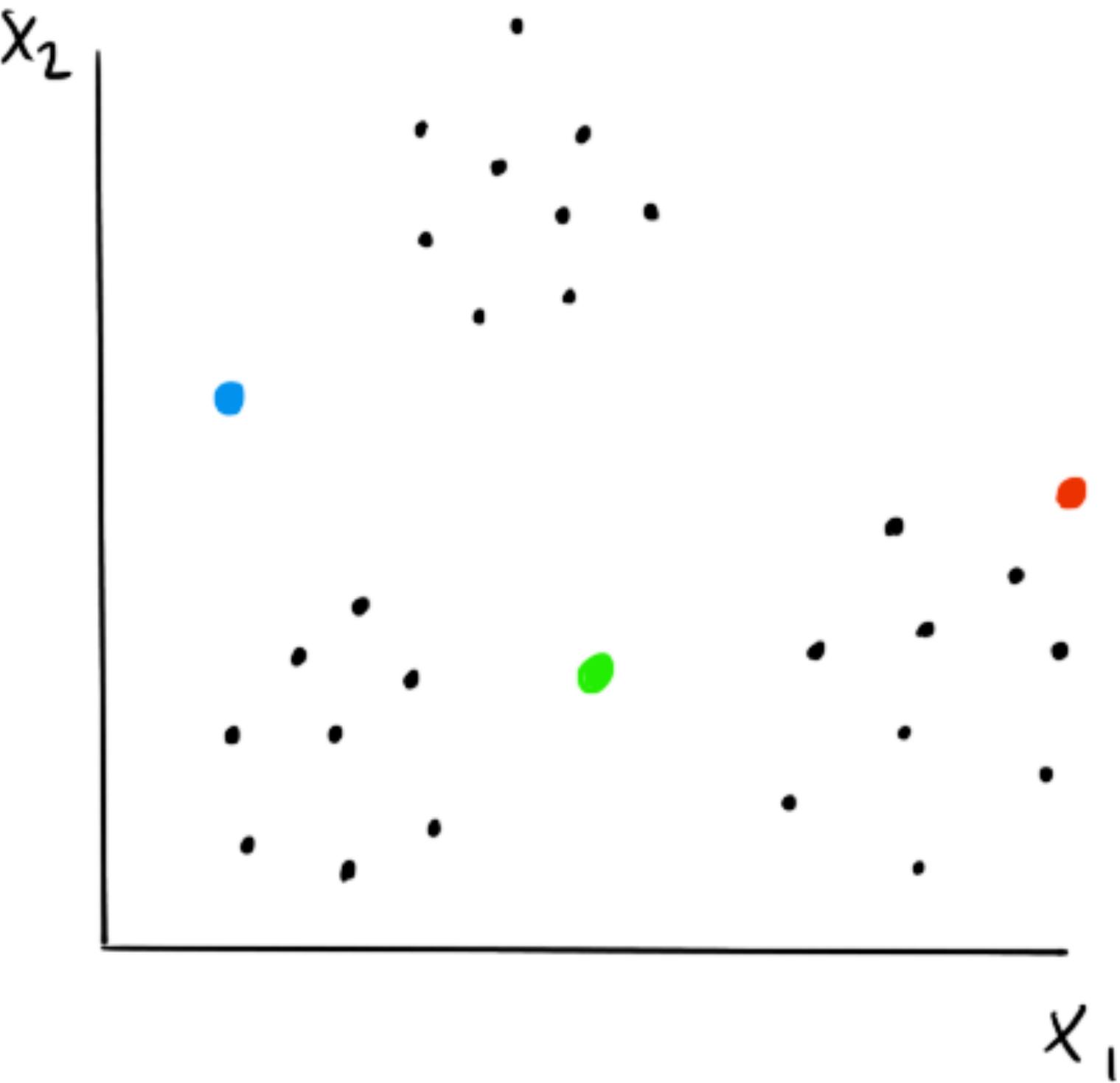


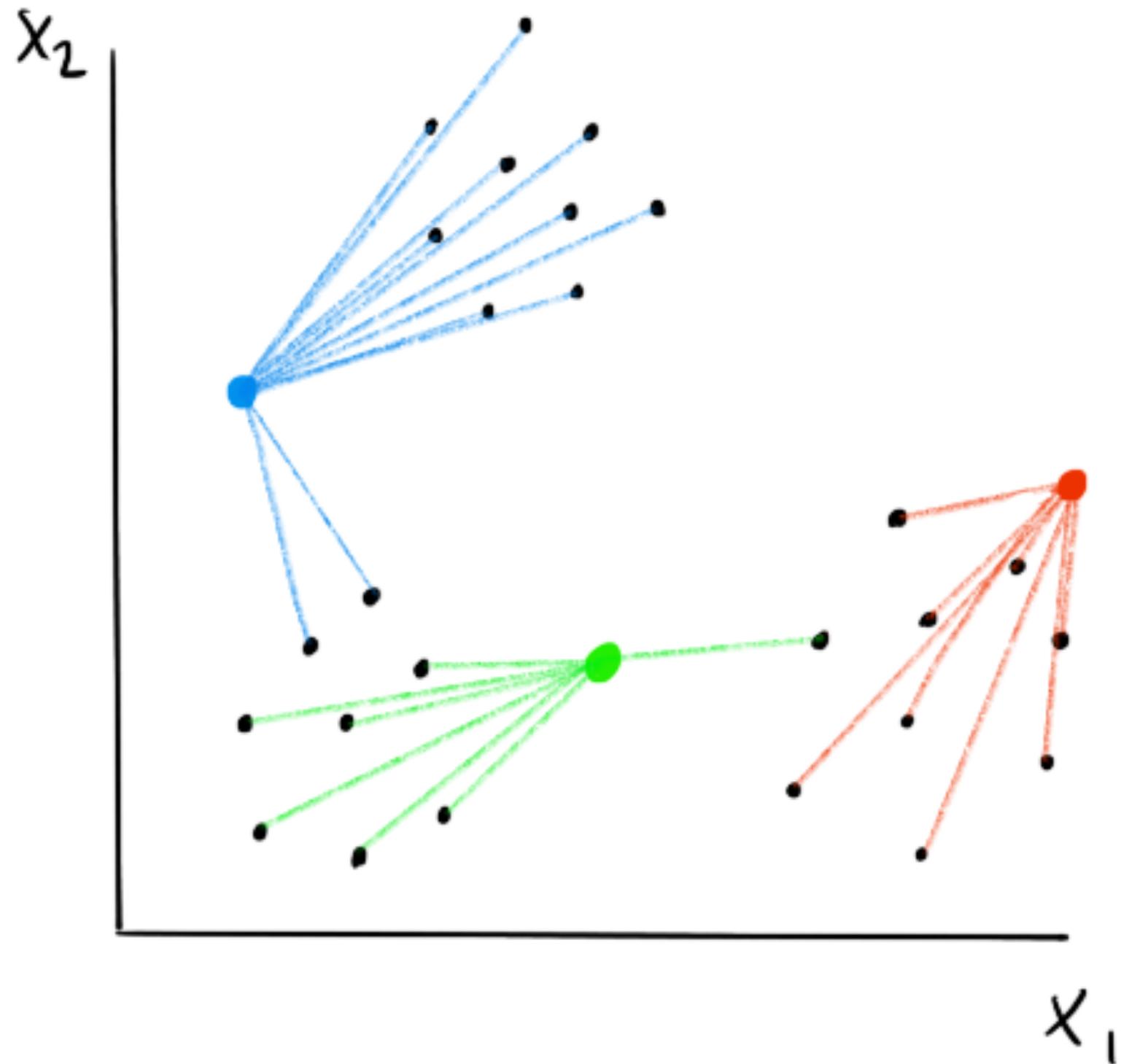


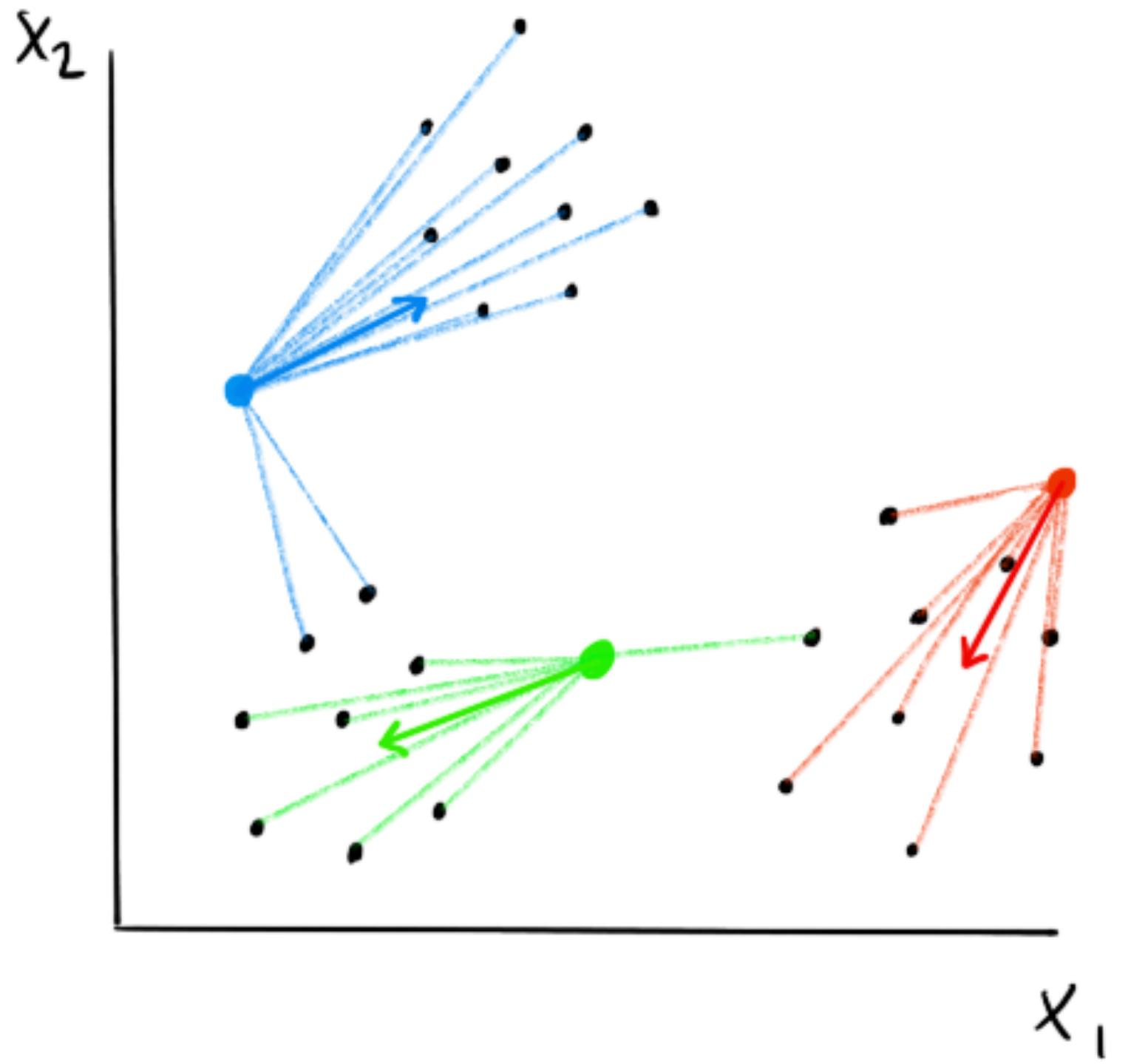
So what can you do with this? Plenty!

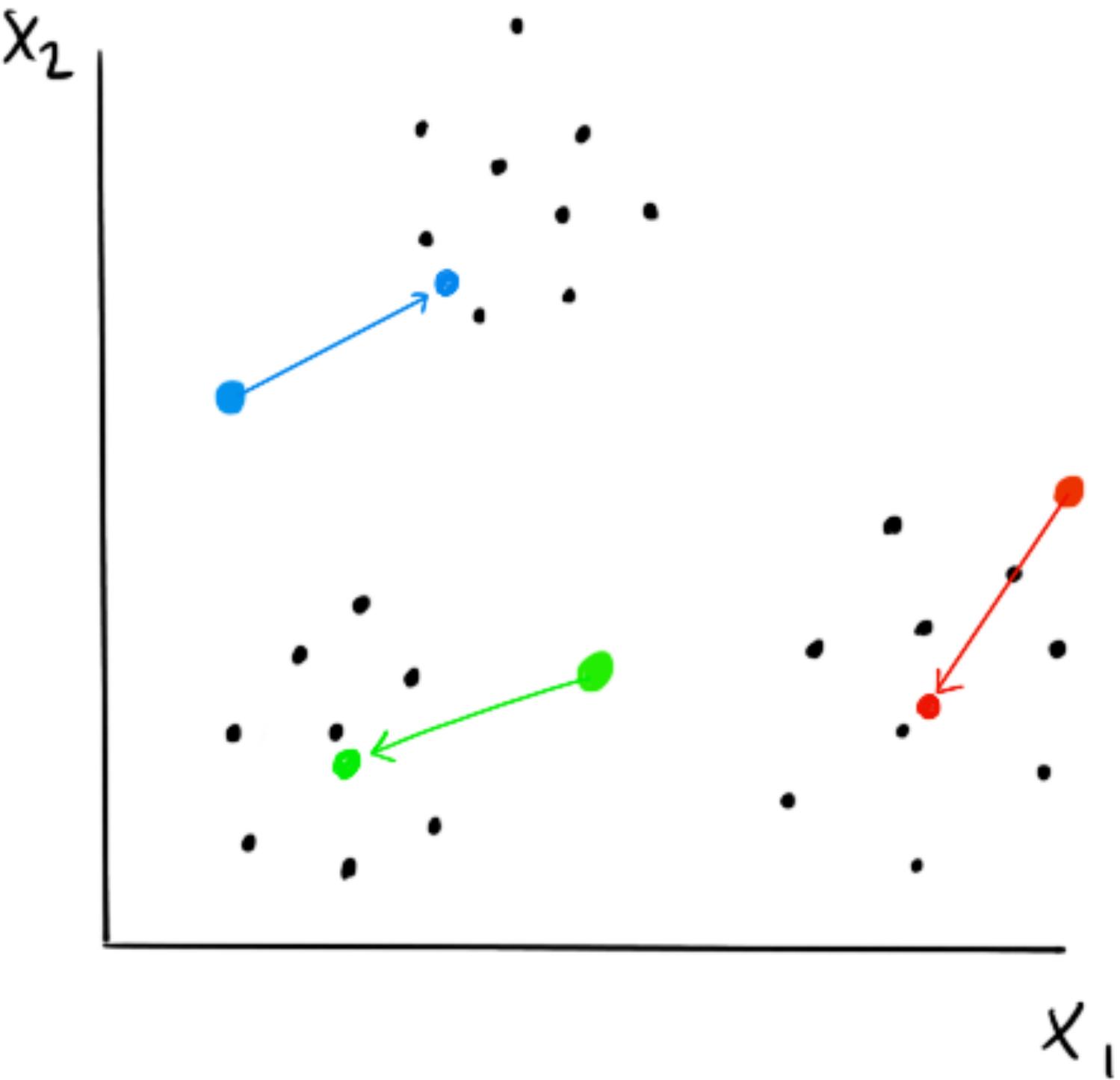
For starters; we can take k -means and make it better.

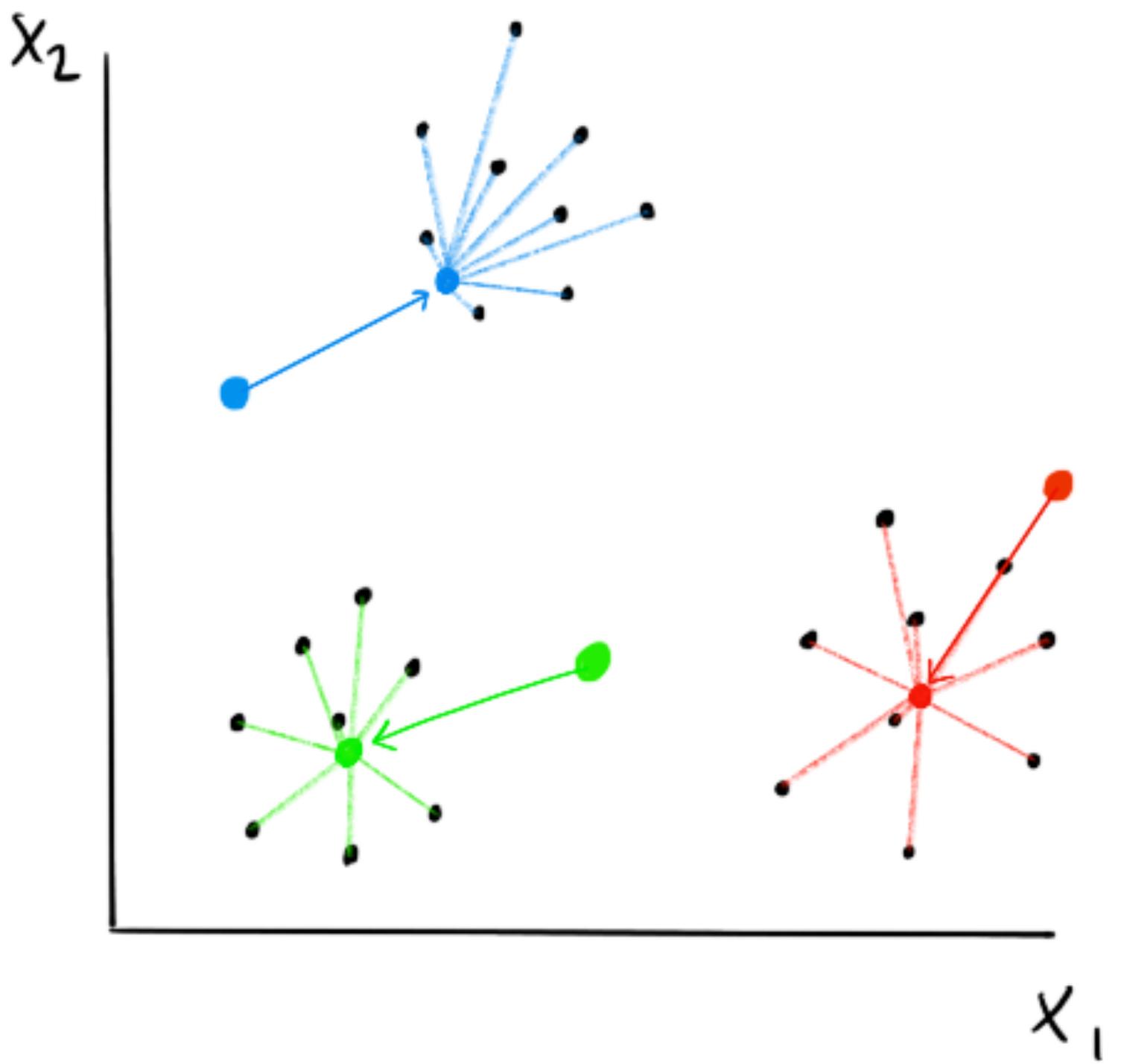


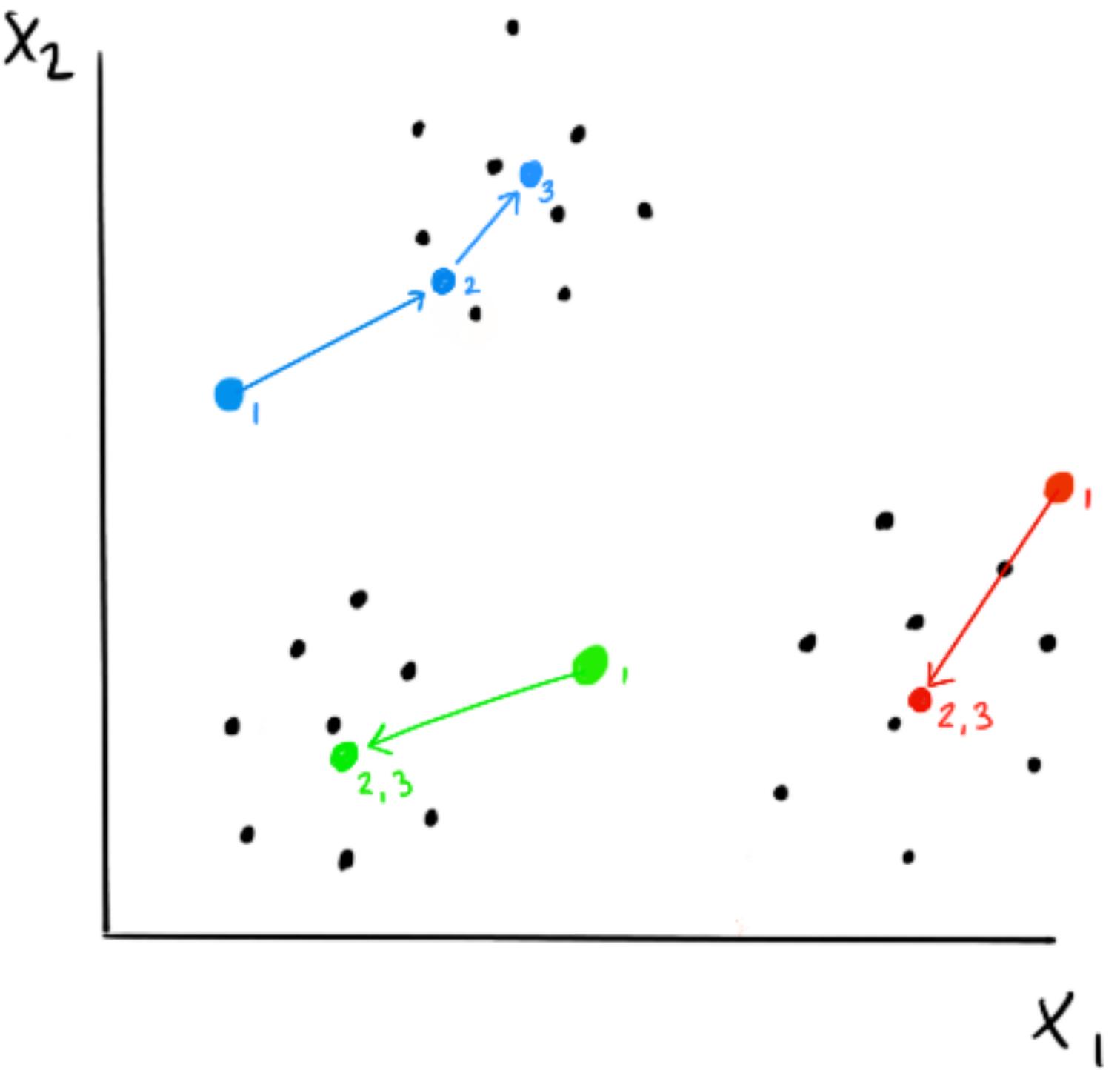










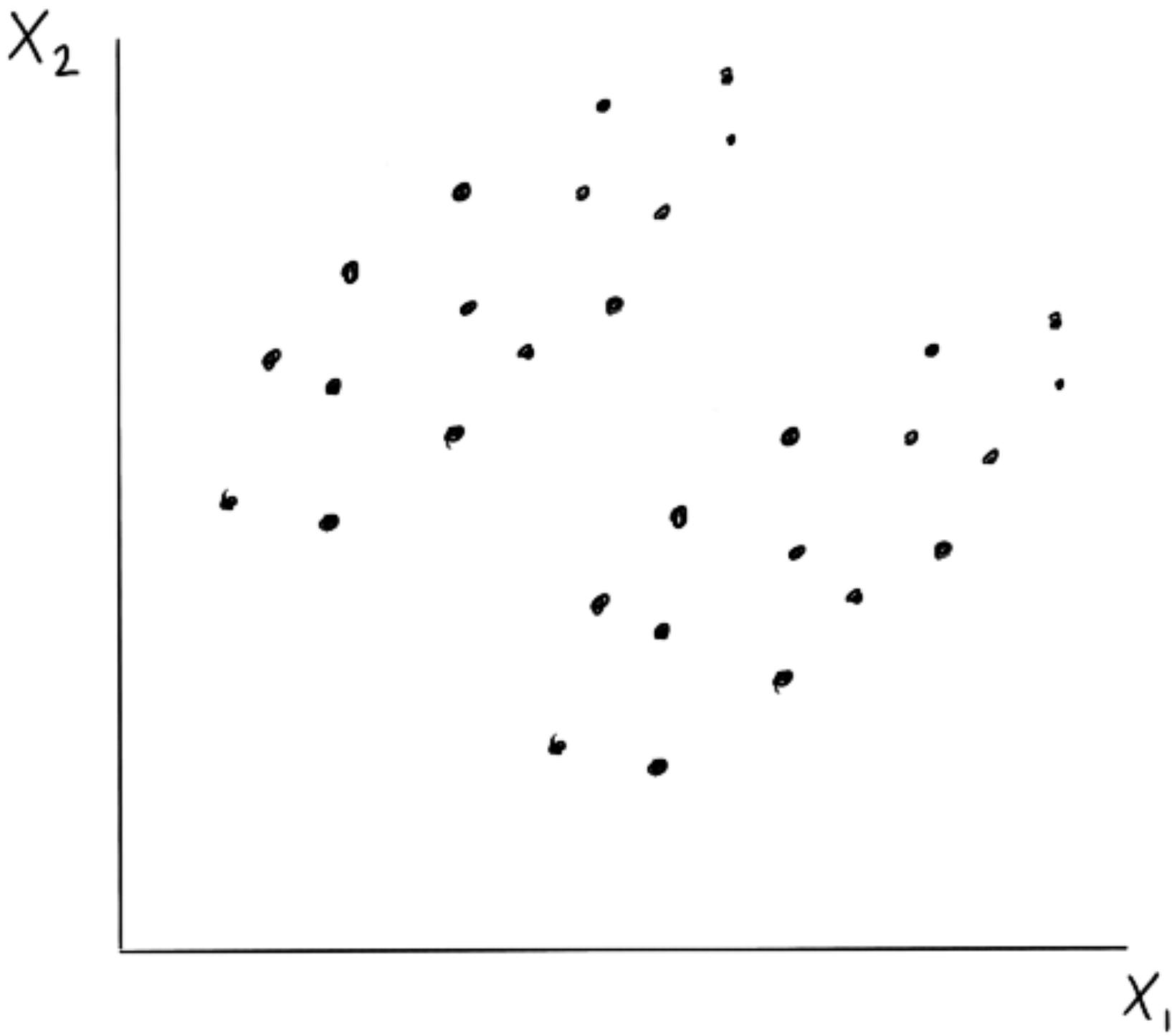


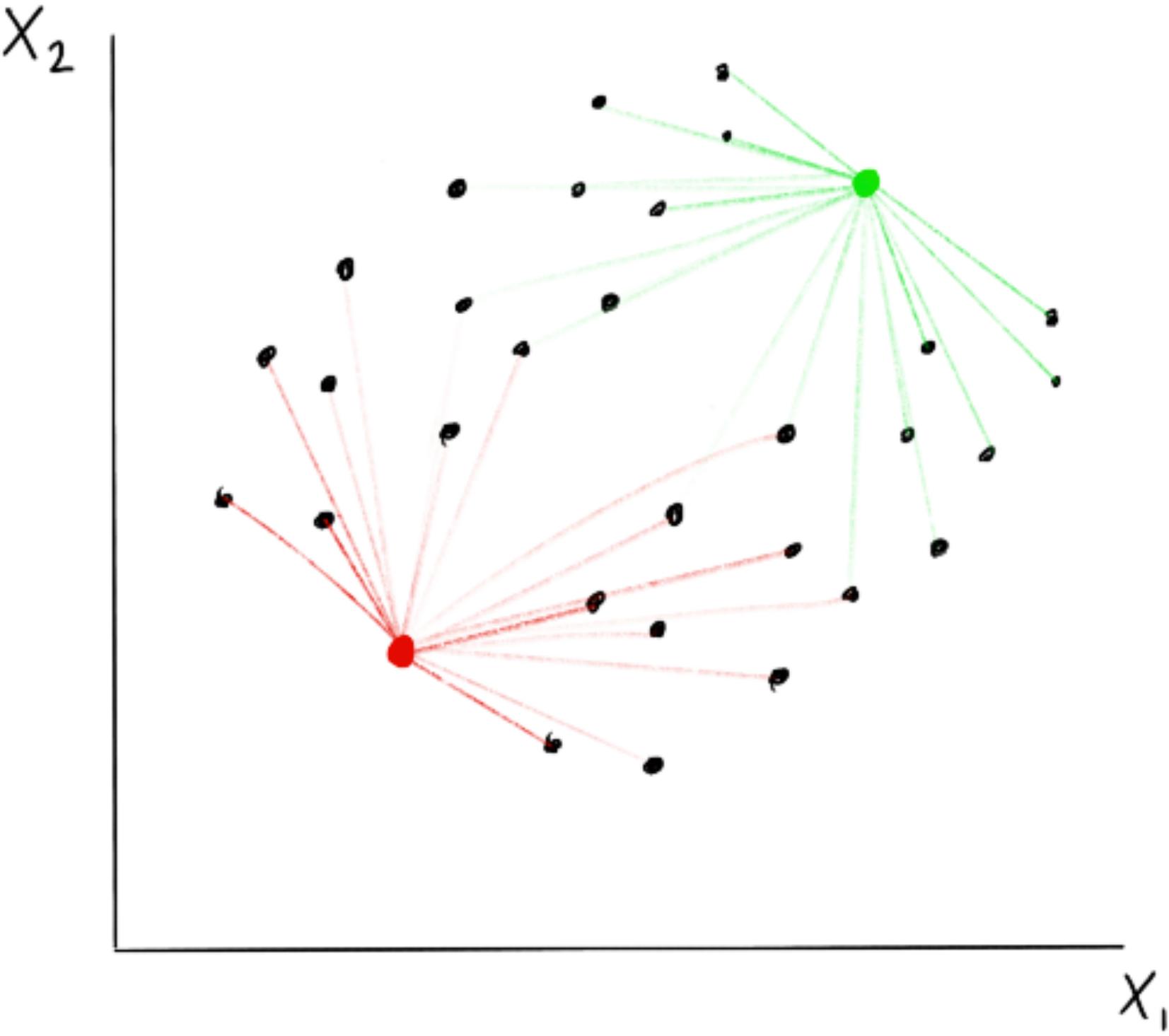
The algorithm has two steps:

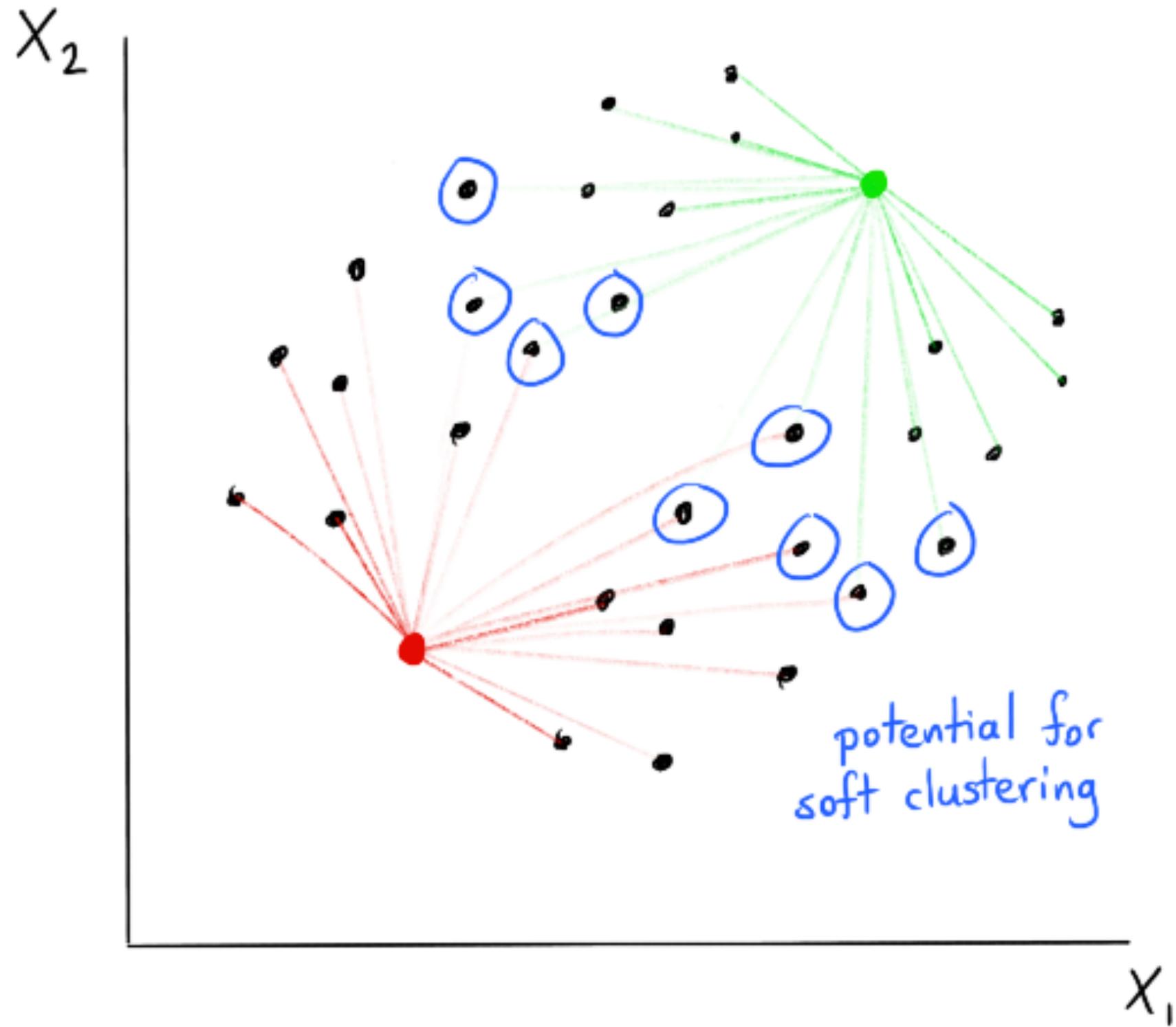
1. a step where we check distances (E)
2. a step where we move centers (M)

This (EM-algorithm) gives it a nice property: it must converge! The next step will always have less total distance between centroids and points (otherwise the centroids would not move).

But there's something "iffy" about all this.

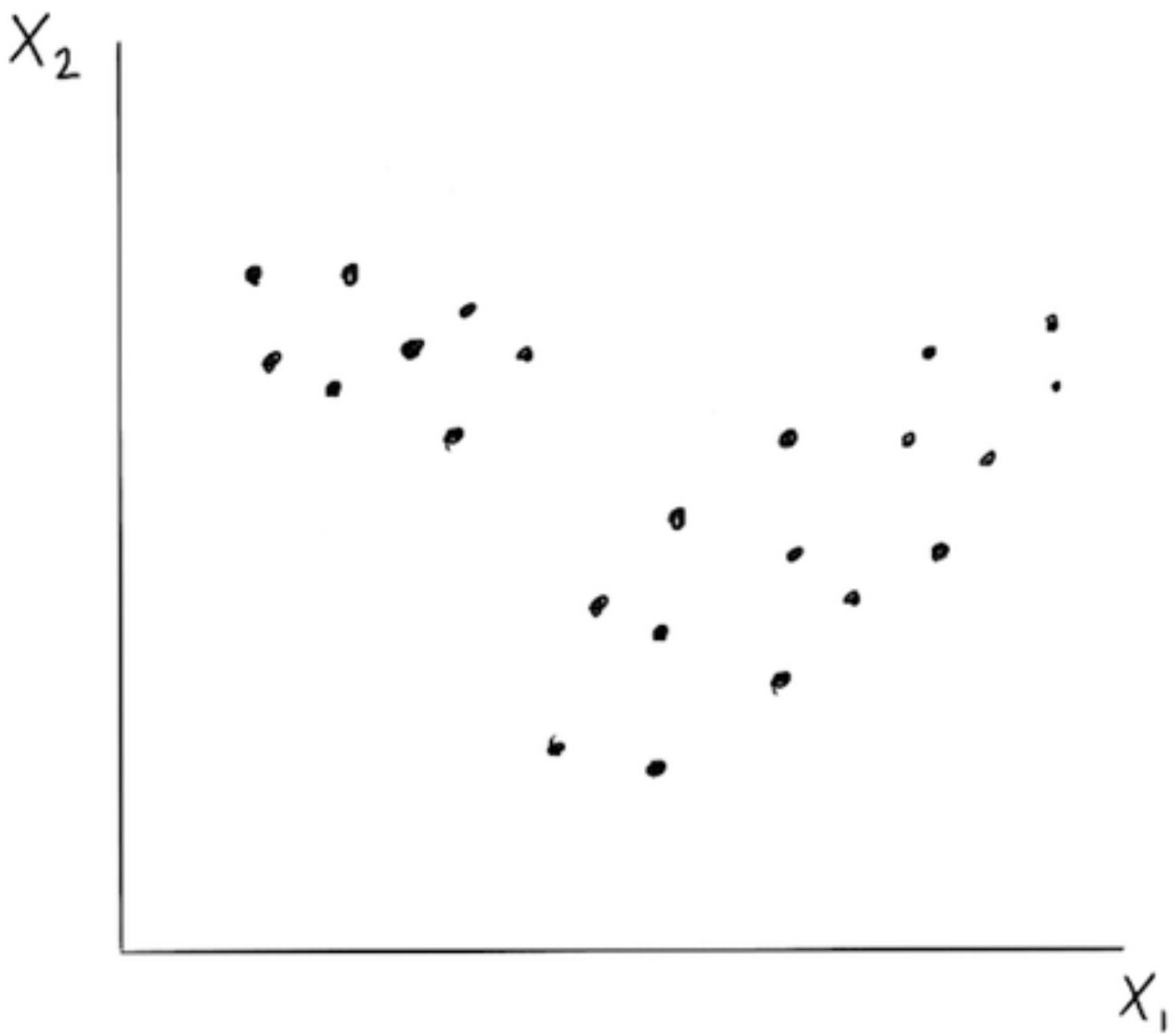


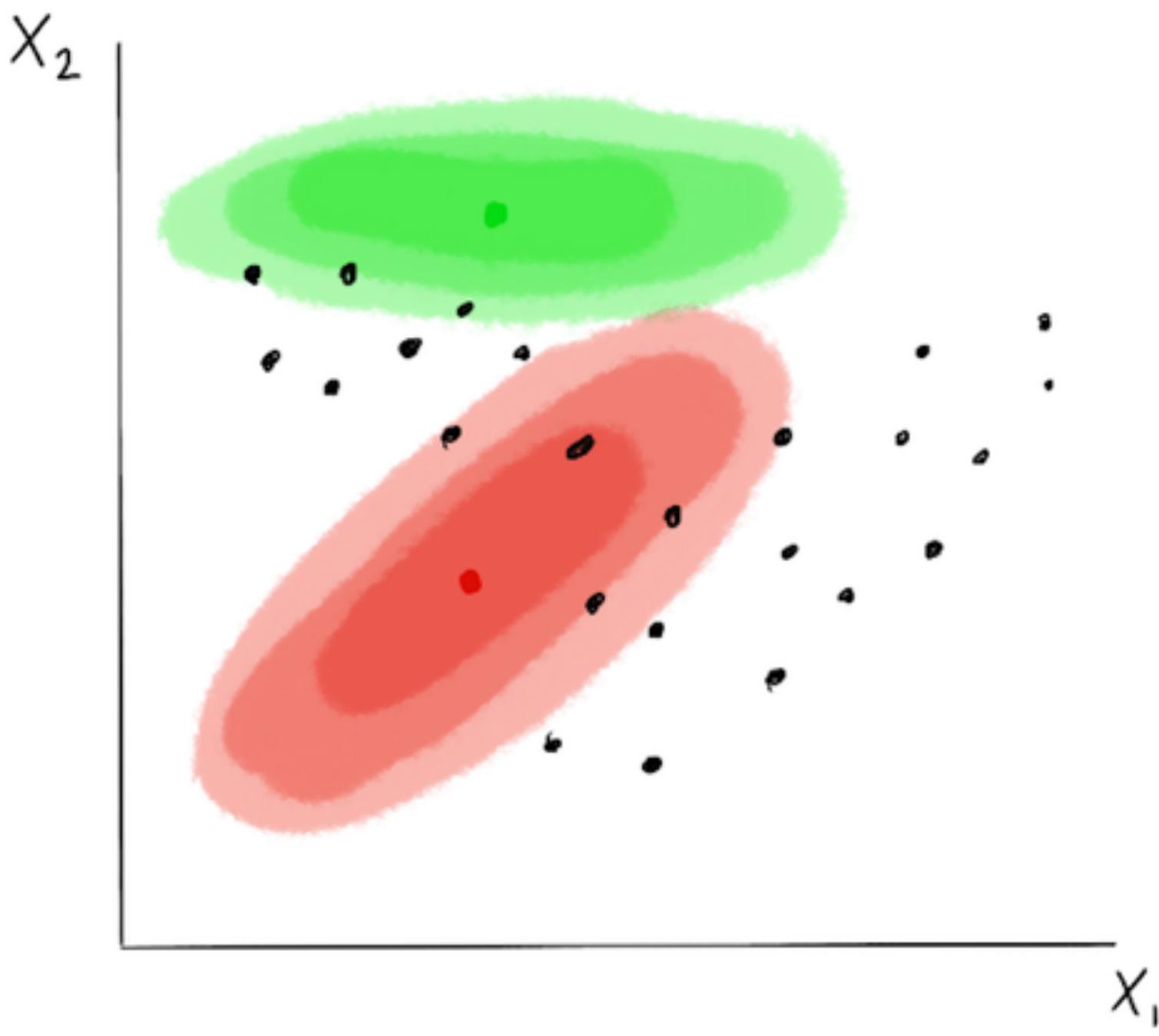


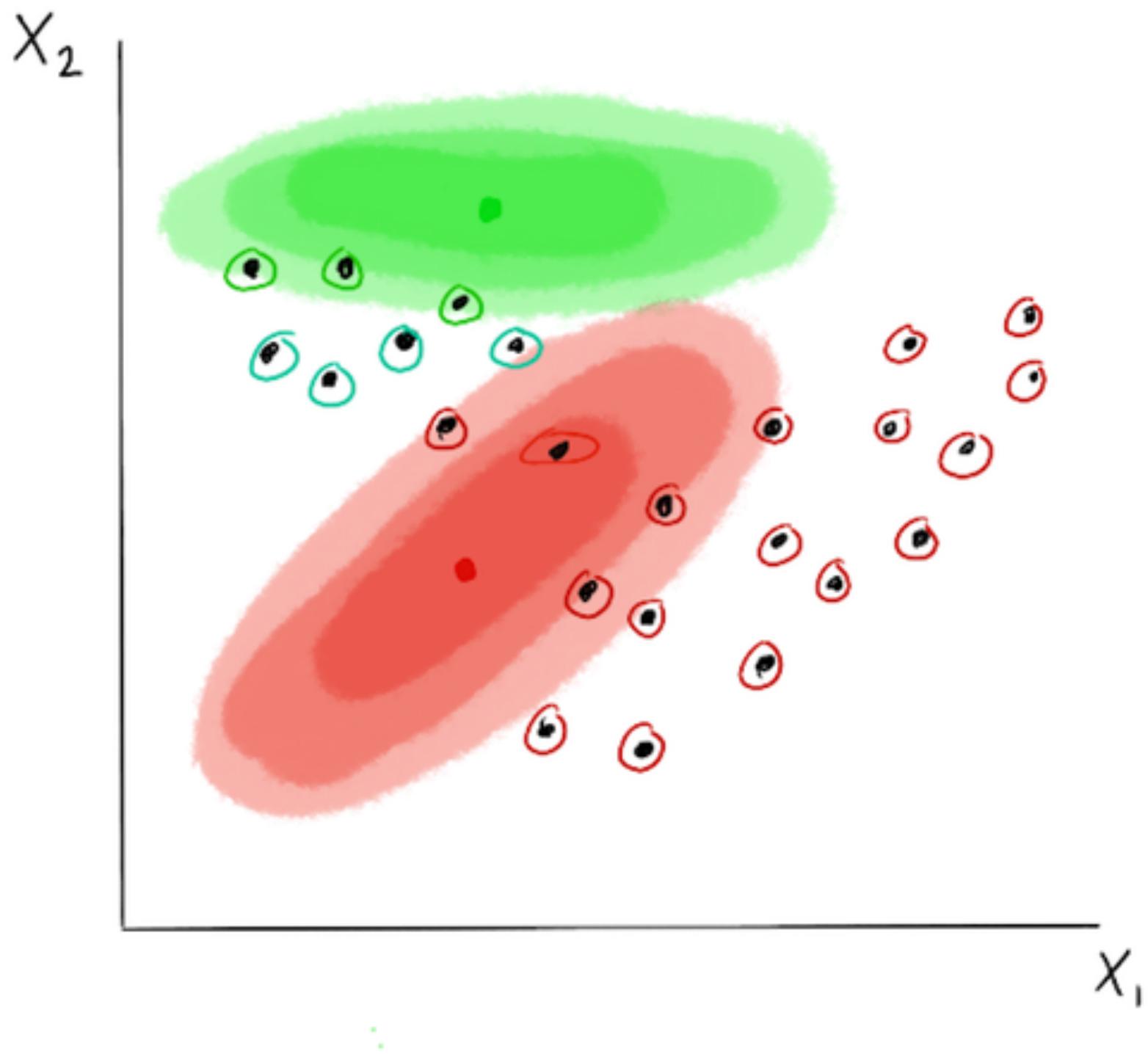


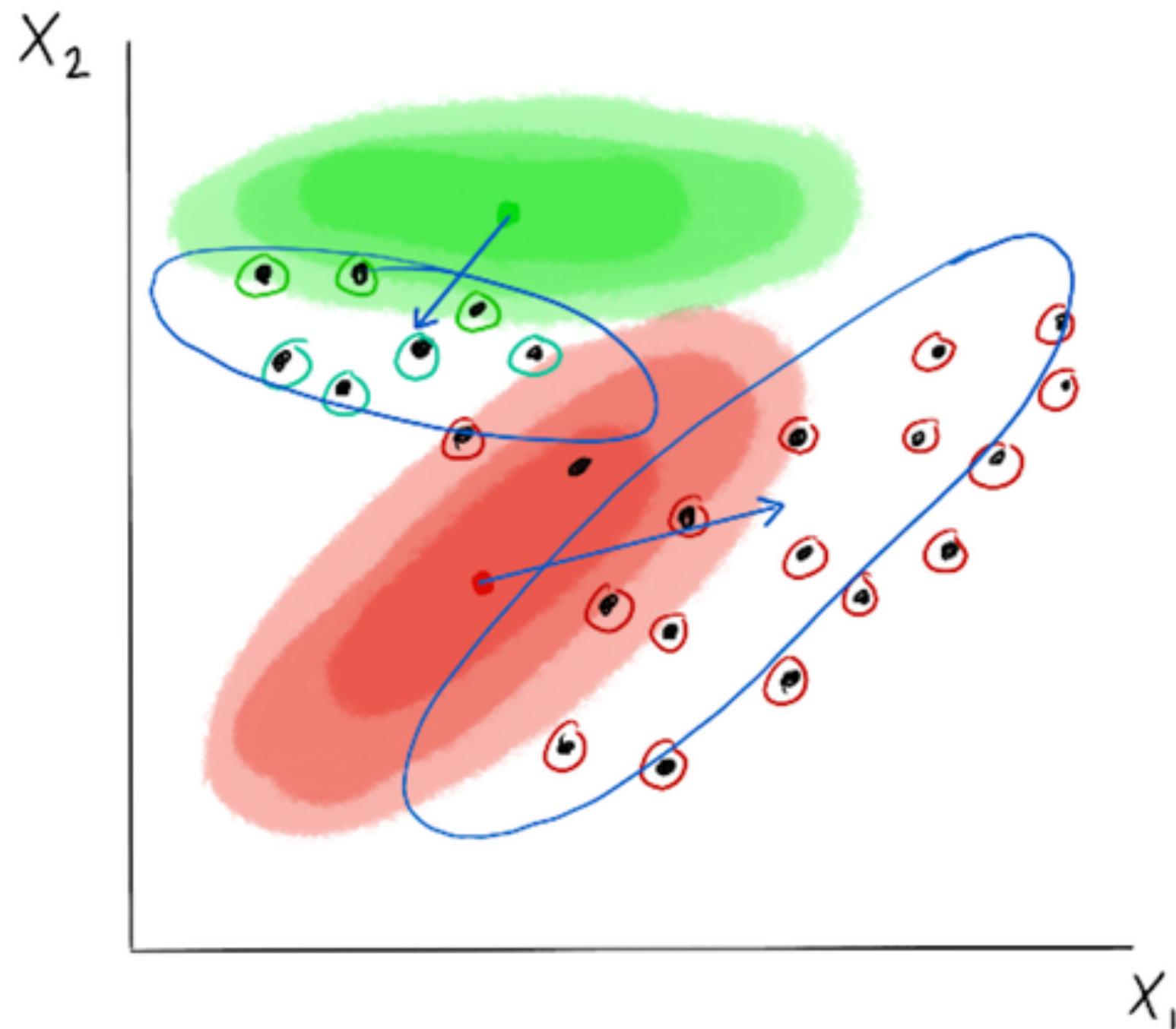
So the idea is; we'll try to do "kmeans" but instead of making hard clusters we'll try to fit multiple gaussians that can perform "soft" clusters. There will be a probability that a point belongs to a certain cluster.

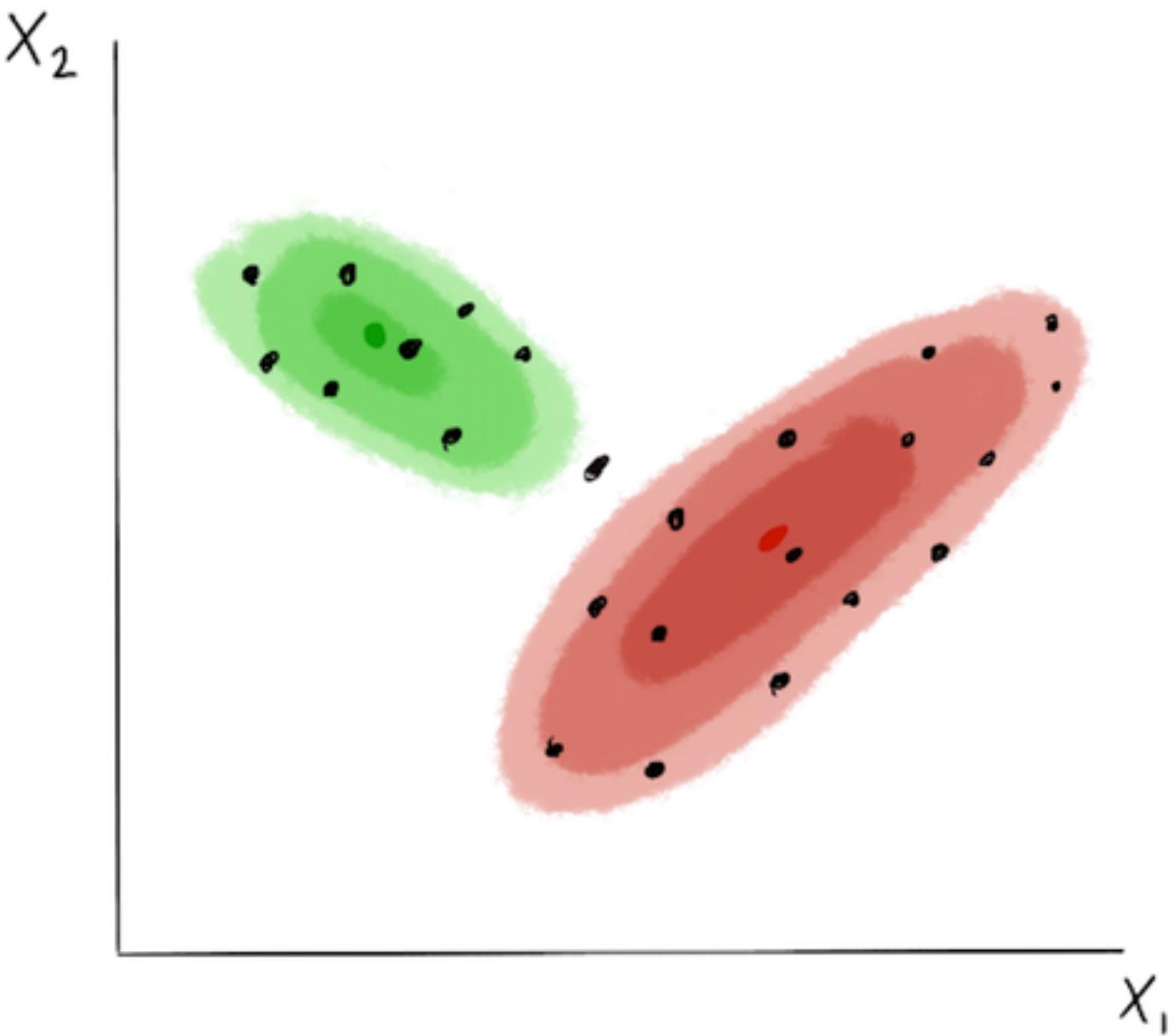
The further a point is from a center the less it influences it. Not a bad idea.

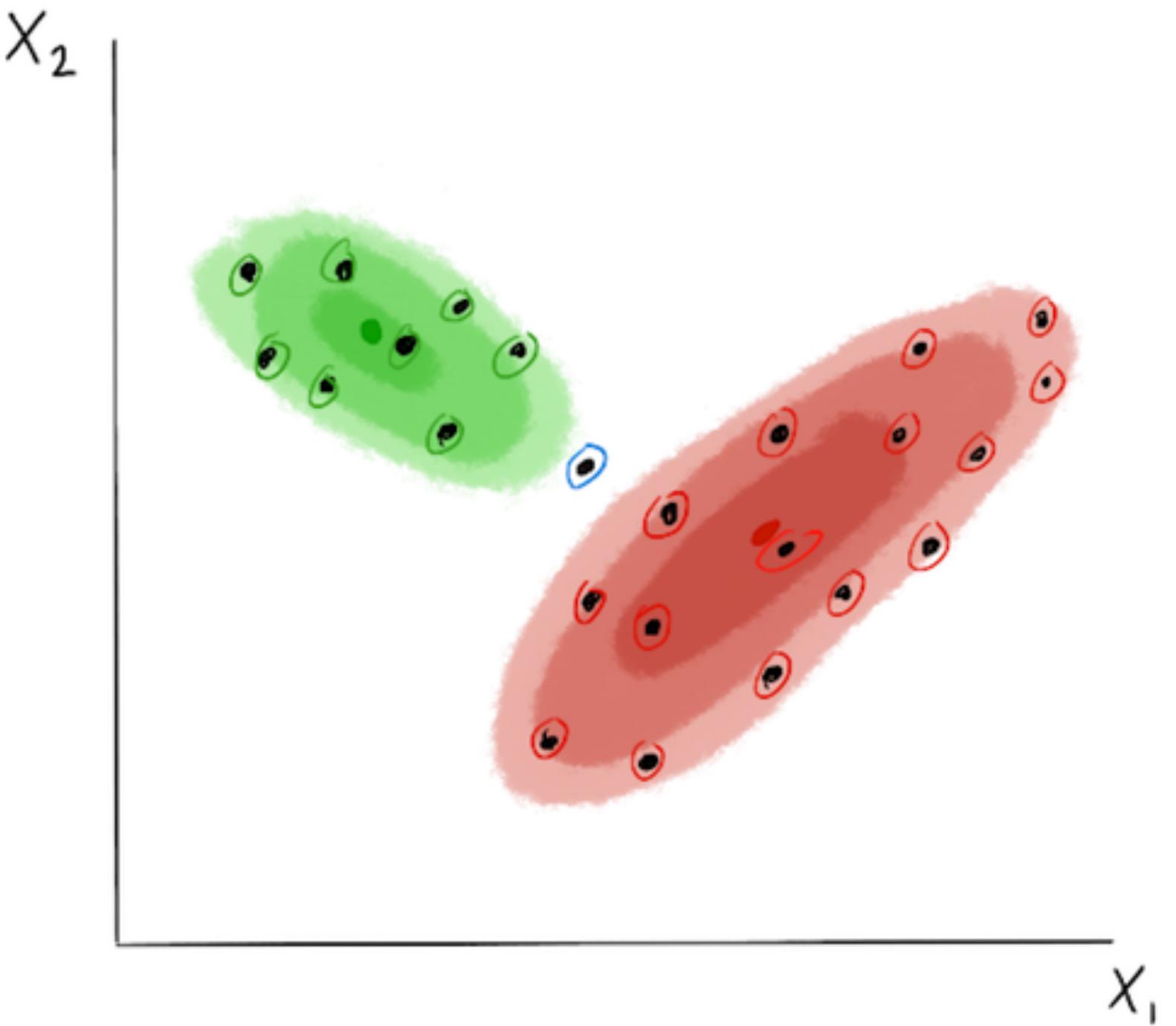


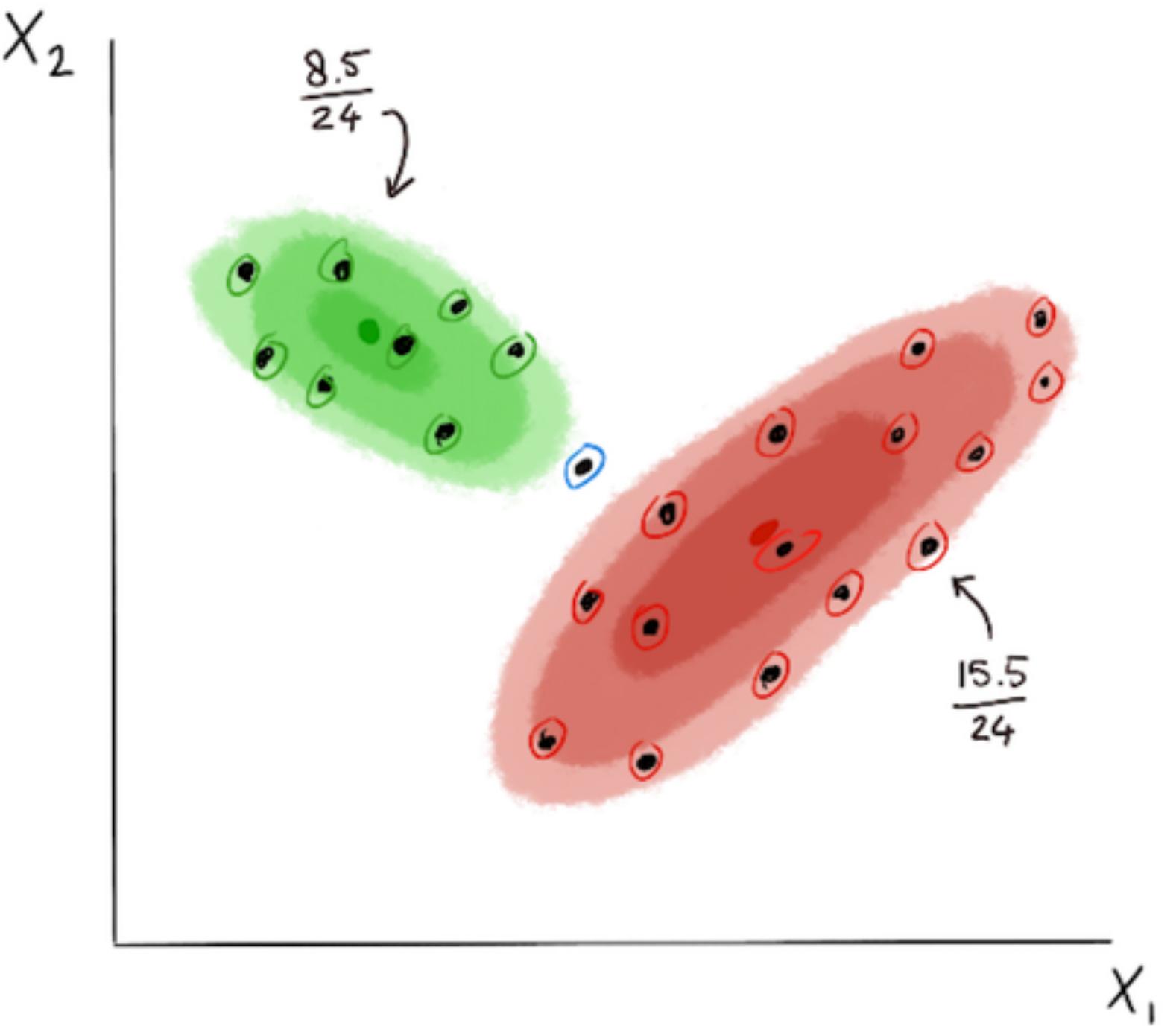




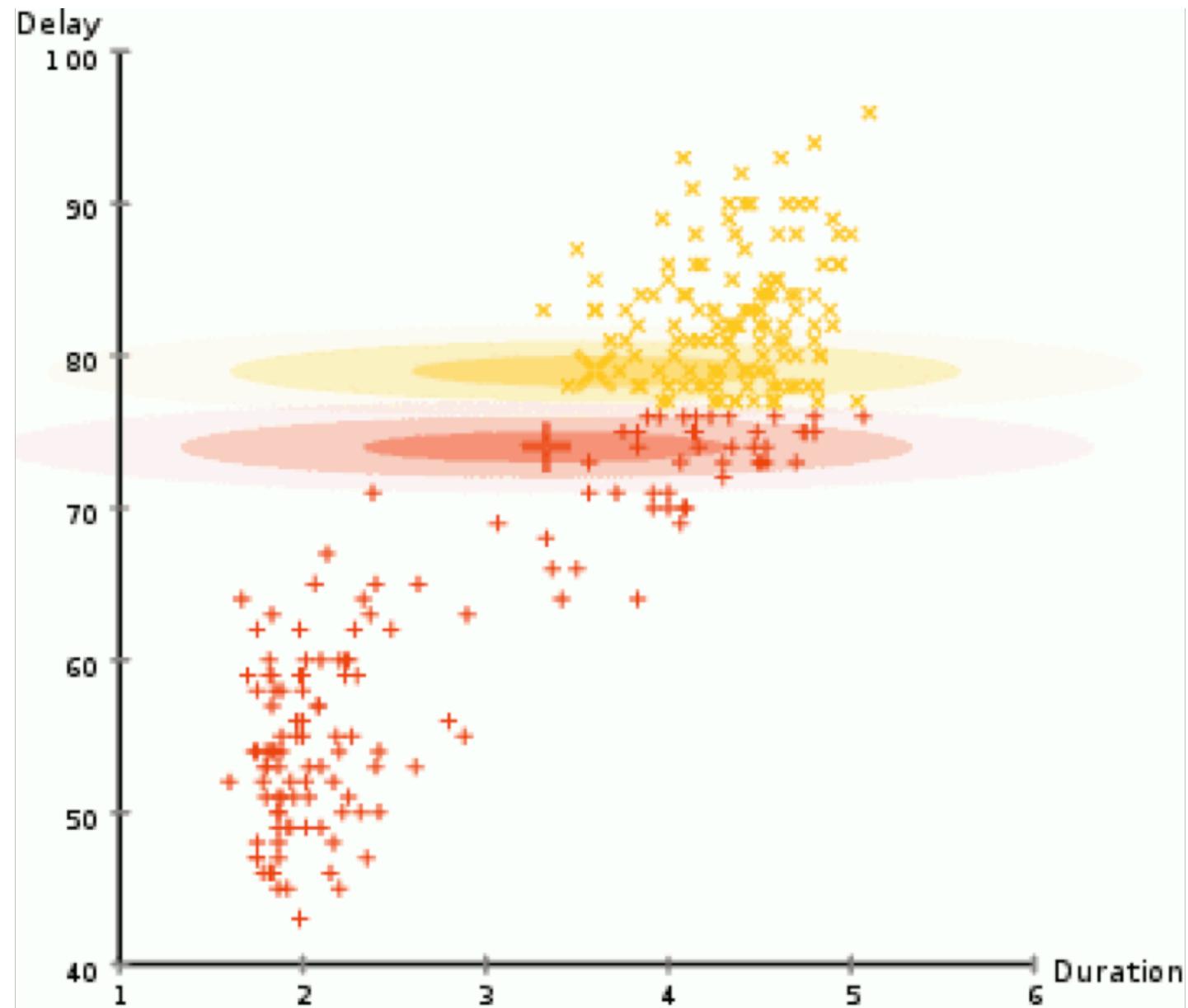








The GIF from Wikipedia



The Math from Wikipedia

And again, this story is usually summarised like this;

$$p(\mathbf{x}) = \sum_{i=1}^k \pi_i \times N(\mu_i, \sigma_i)$$

There k gaussians that we'll try to fit and each gaussian gets an associated weight π_i .

For every centroid k we'll calculate;

$$n_k = \sum_{i=1}^n p_k(i)$$

$$\mu_k = \frac{1}{n_k} \sum_{i=1}^n p_k(i)x_i$$

$$\sigma_k^2 = \frac{1}{n_k} \sum_{i=1}^n p_k(i)(x_i - \mu_k)^2$$

$$\pi_k = \frac{n_k}{n}$$

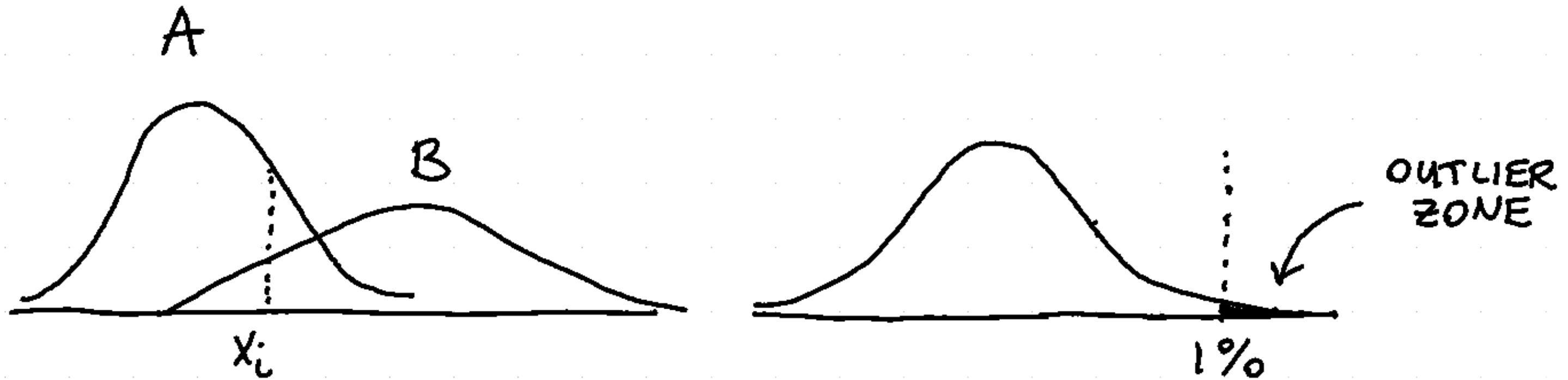
The downside of this approach of explaining things is that it distracts you from the coolest observation:

The downside of this approach of explaining things is that it distracts you from the coolest observation:

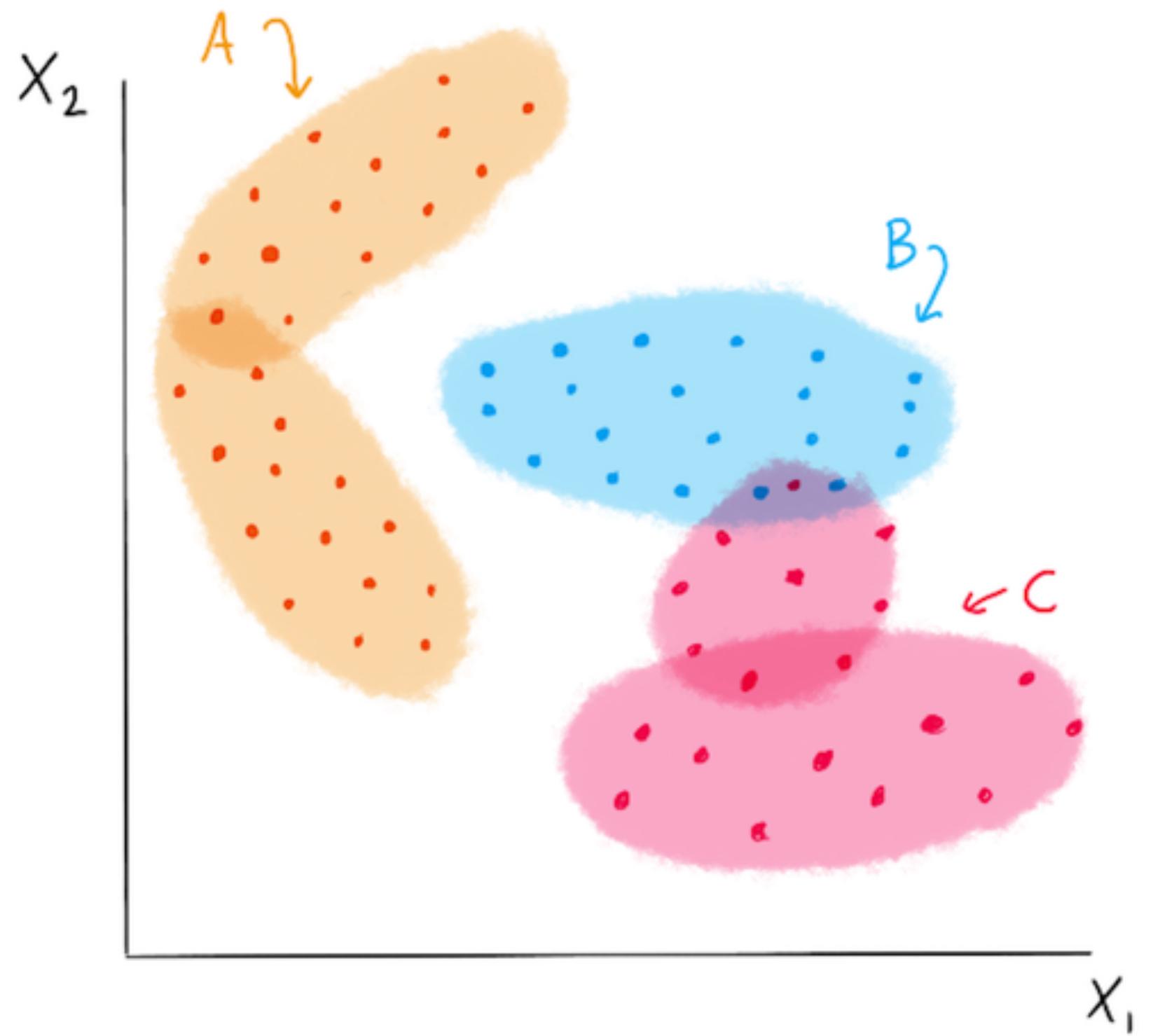
**THIS HAS SO MANY
APPLICATIONS**

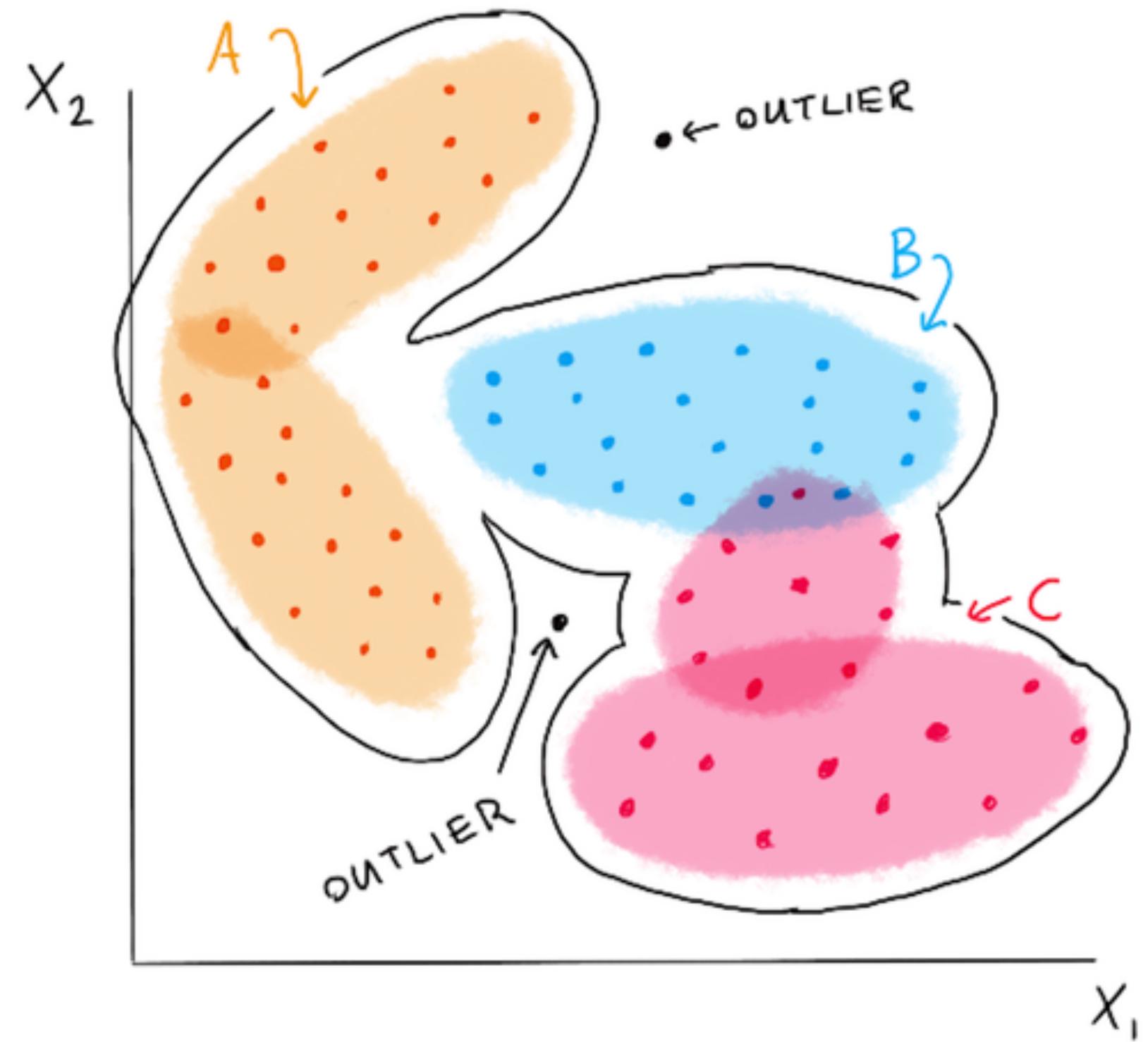
Applications beyond clustering help motivate learning.

Let's do this idea but now in higher dimensions.

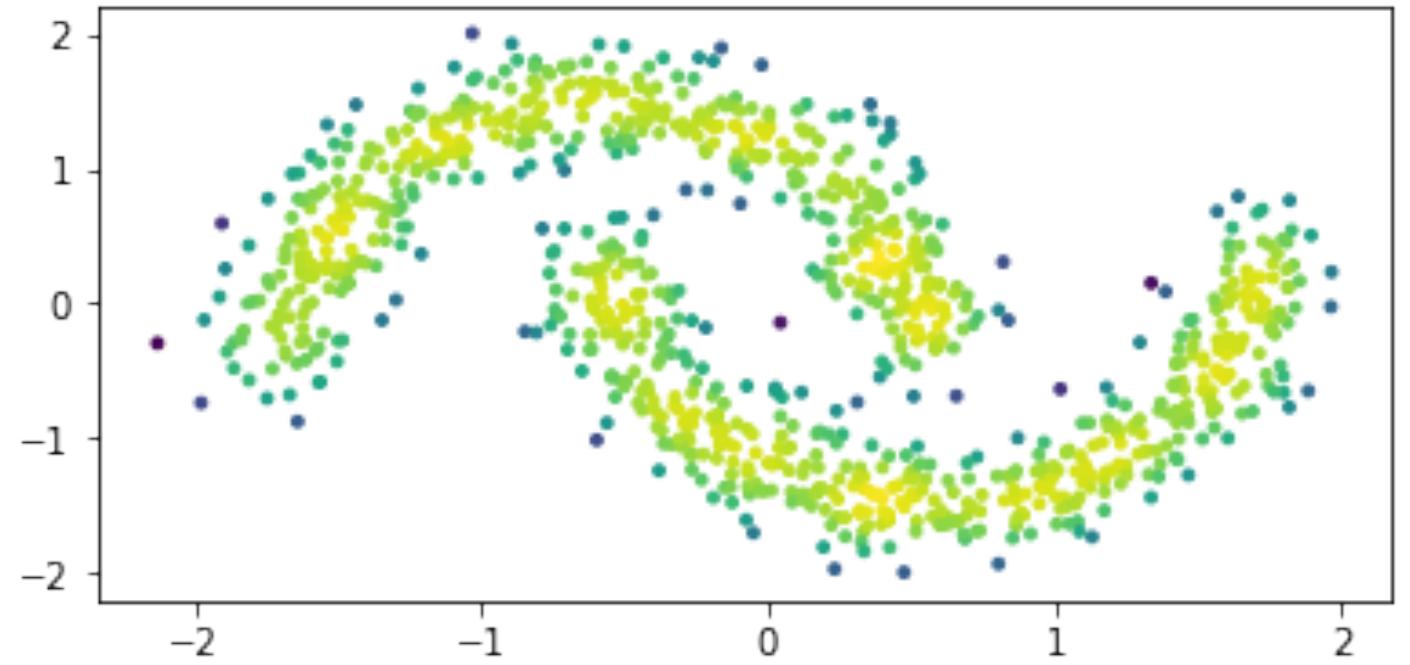


$$p(x_i = A) = \frac{f_A(x_i)}{f_A(x_i) + f_B(x_i)}$$

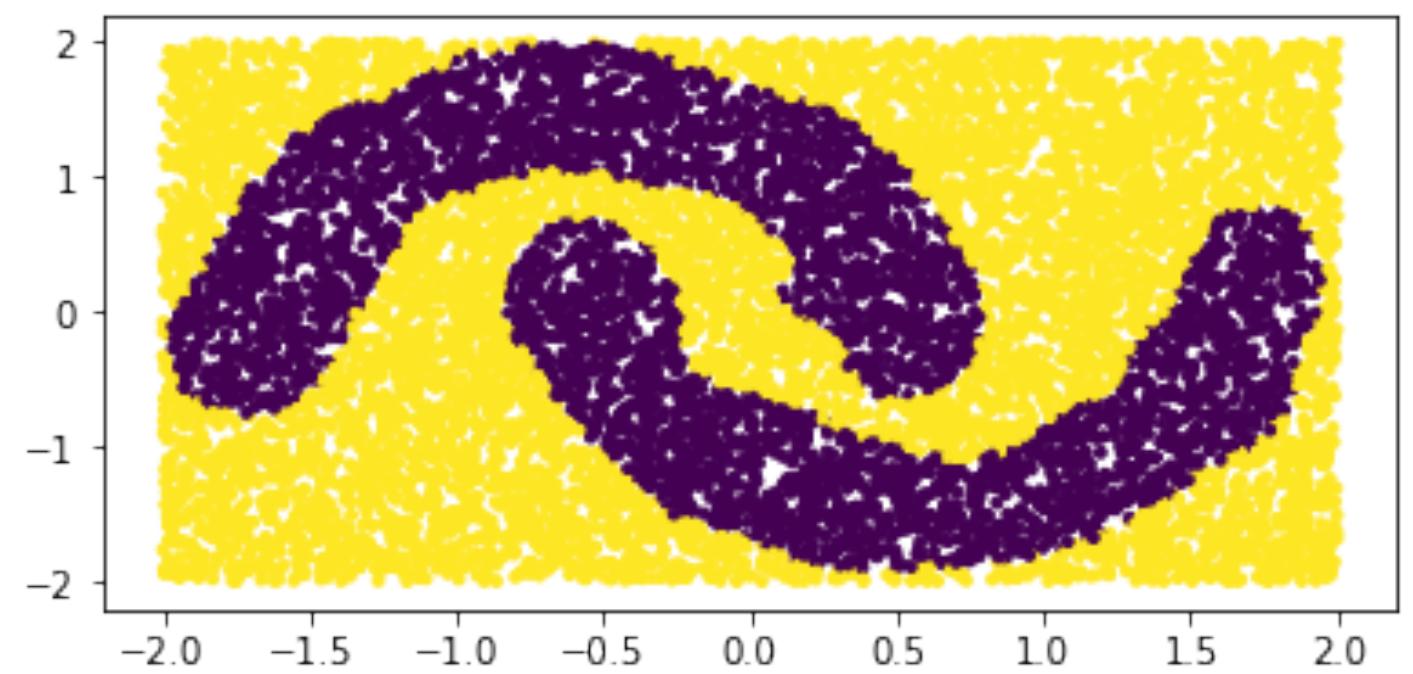


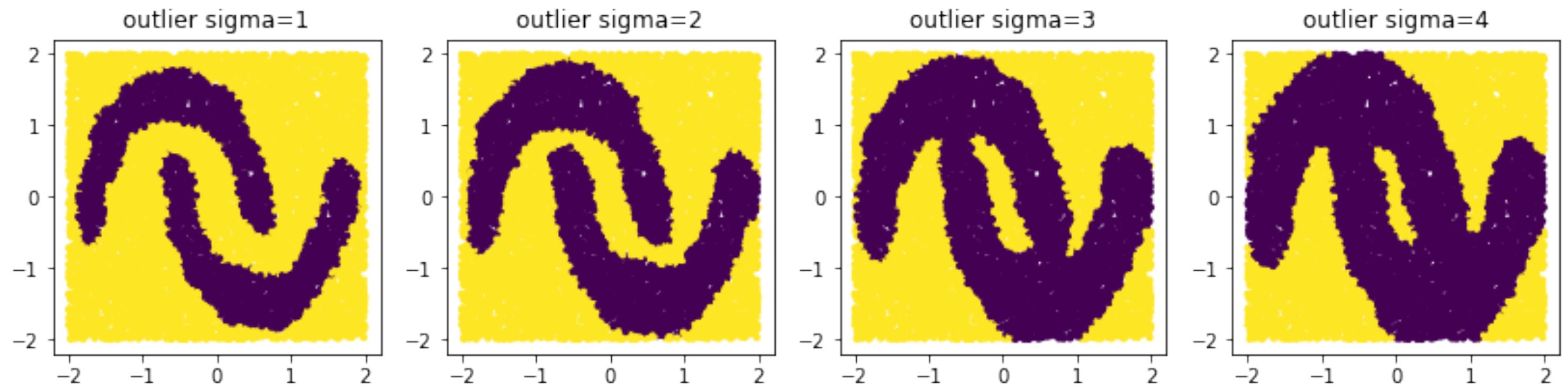


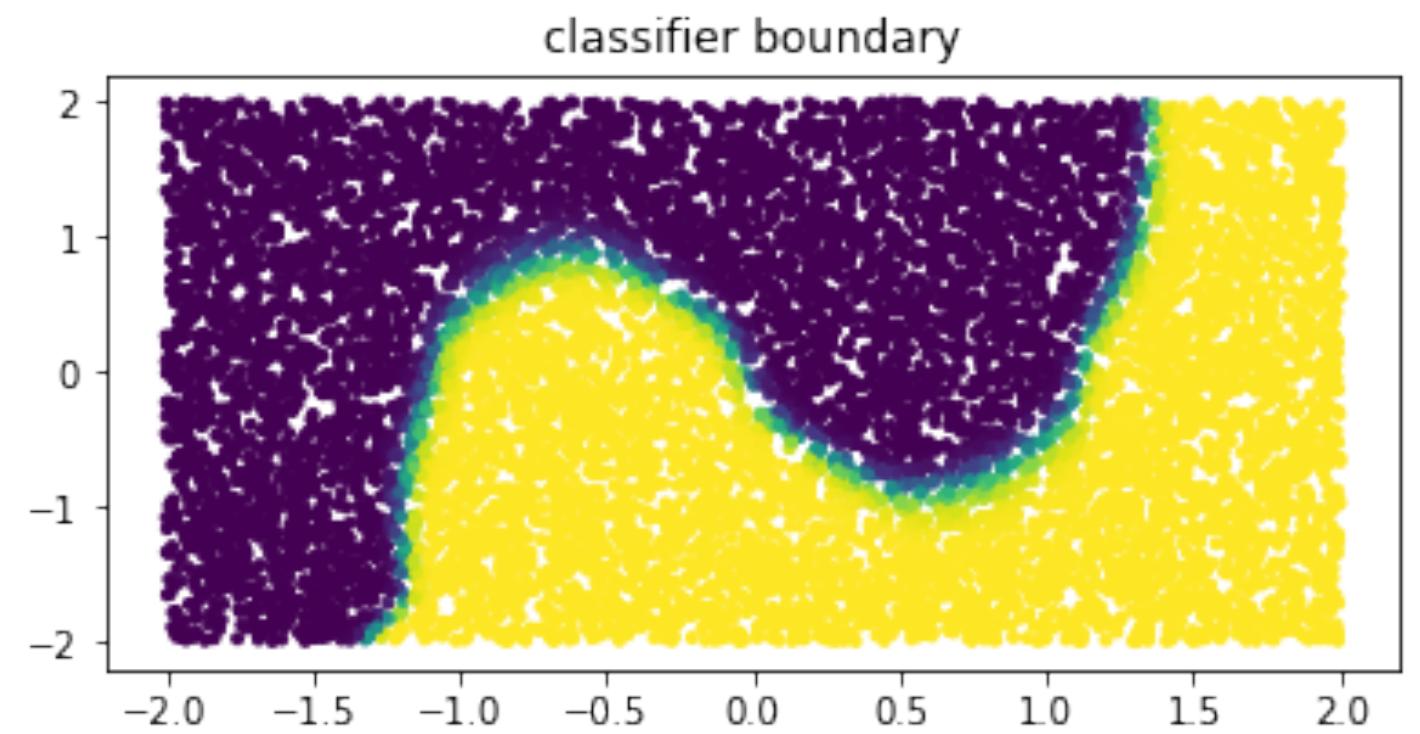
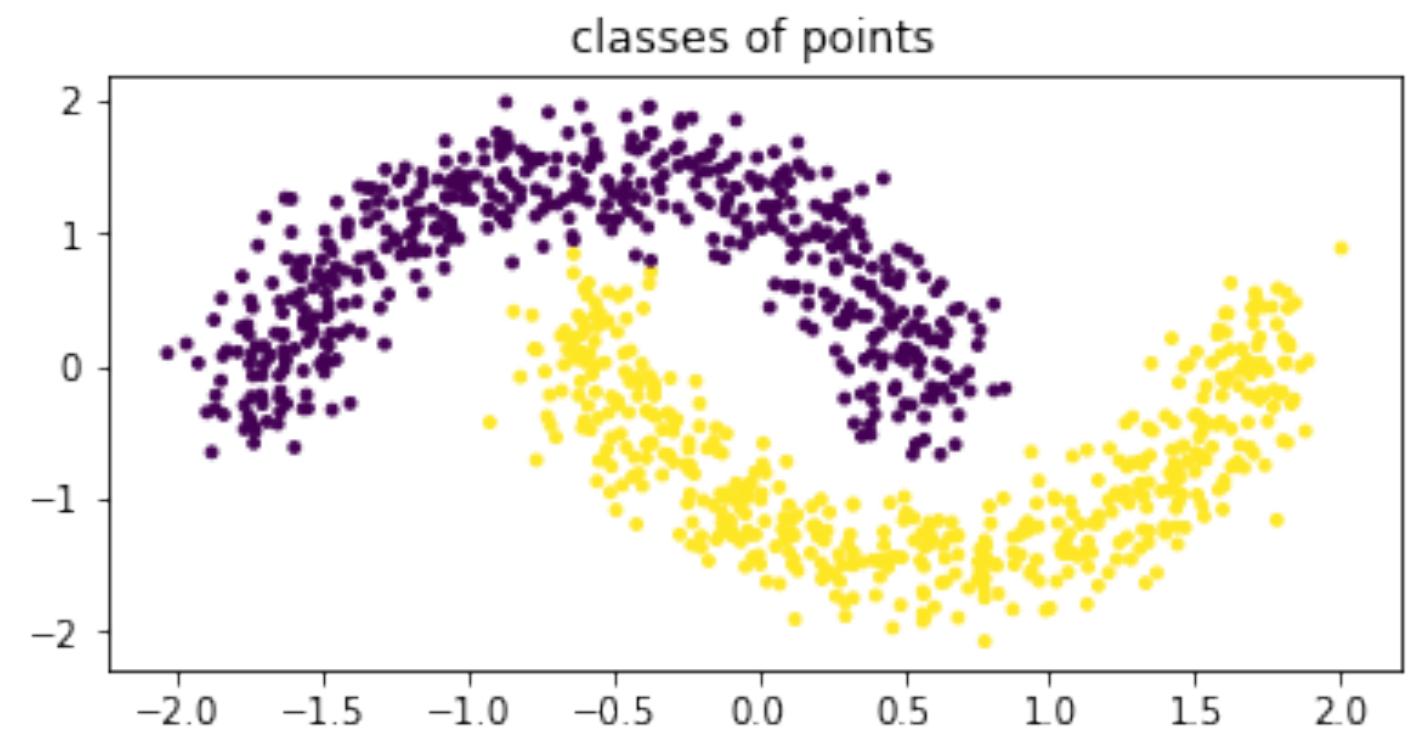
likelihood of points given mixture of 16 gaussians



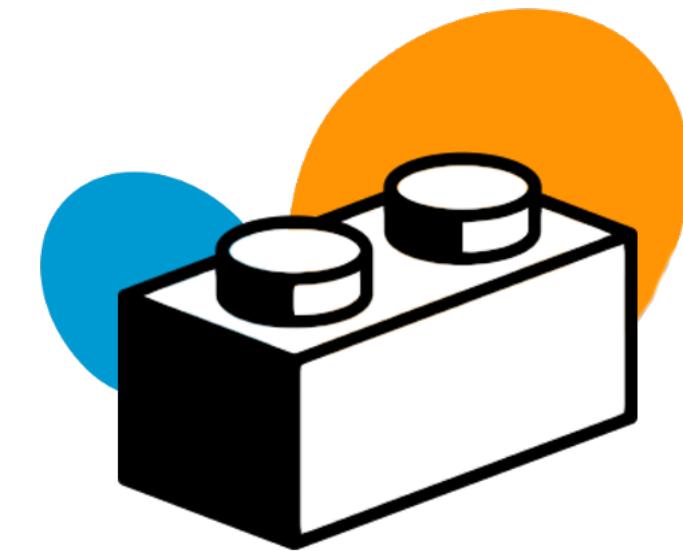
outlier selection







scikit lego

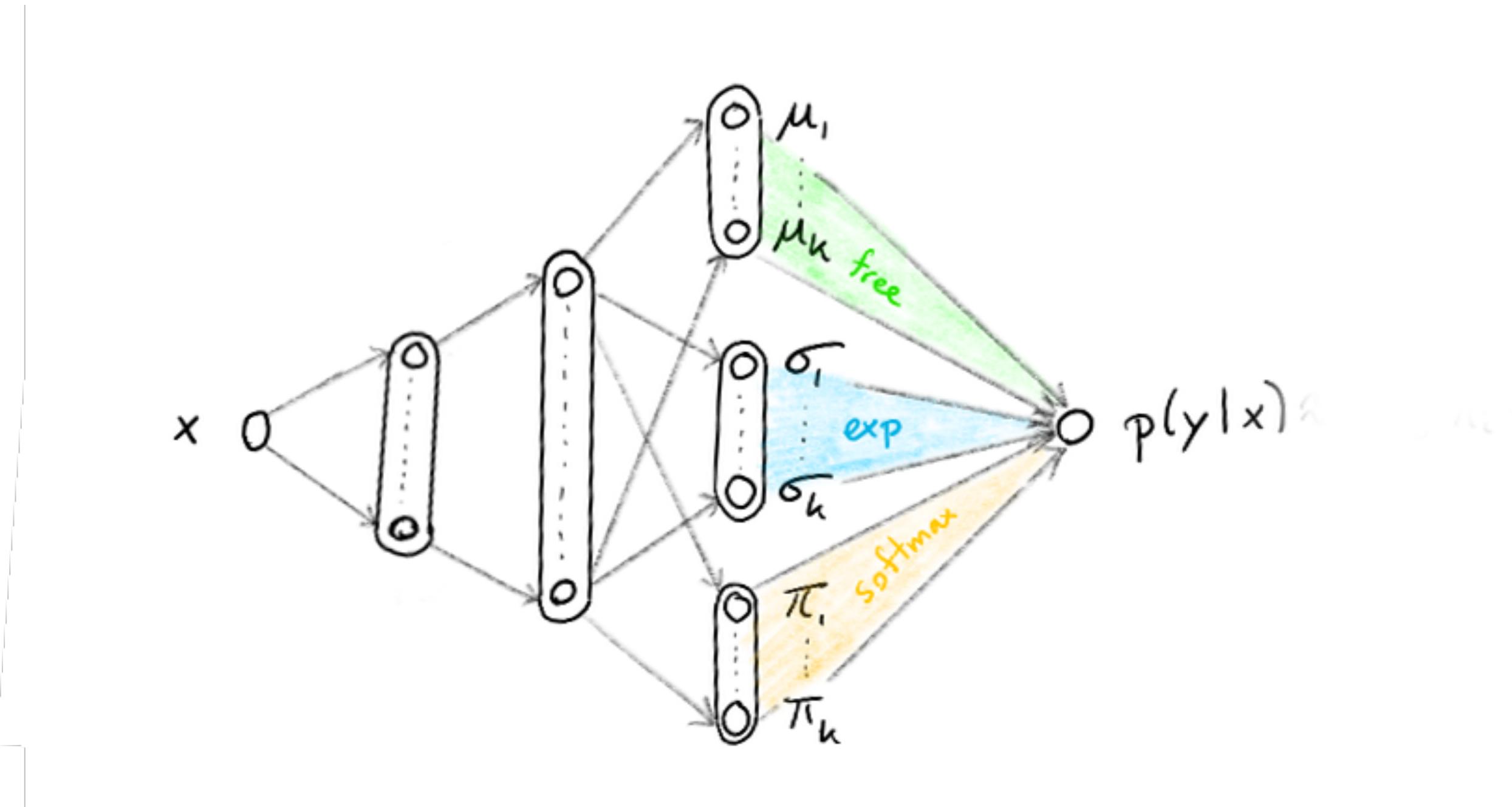


```
from sklego.mixture import GMMClassifier  
from sklego.mixture import GMMOutlierDetector
```

Open sourced! Feel free to find us on github.

But what about deep learning?

Well ... we can add useful properties to neural networks by glueing a gaussian on top of it.

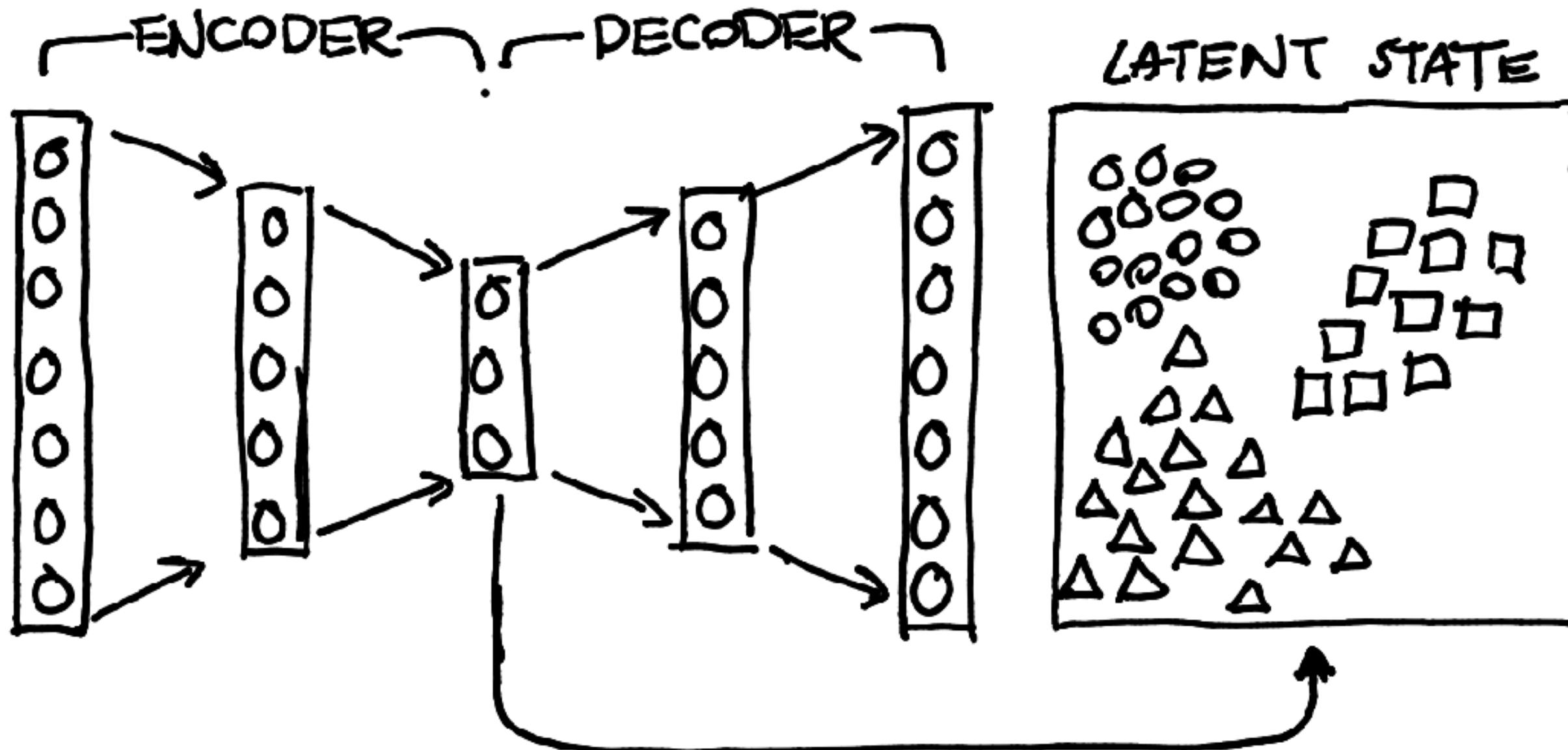


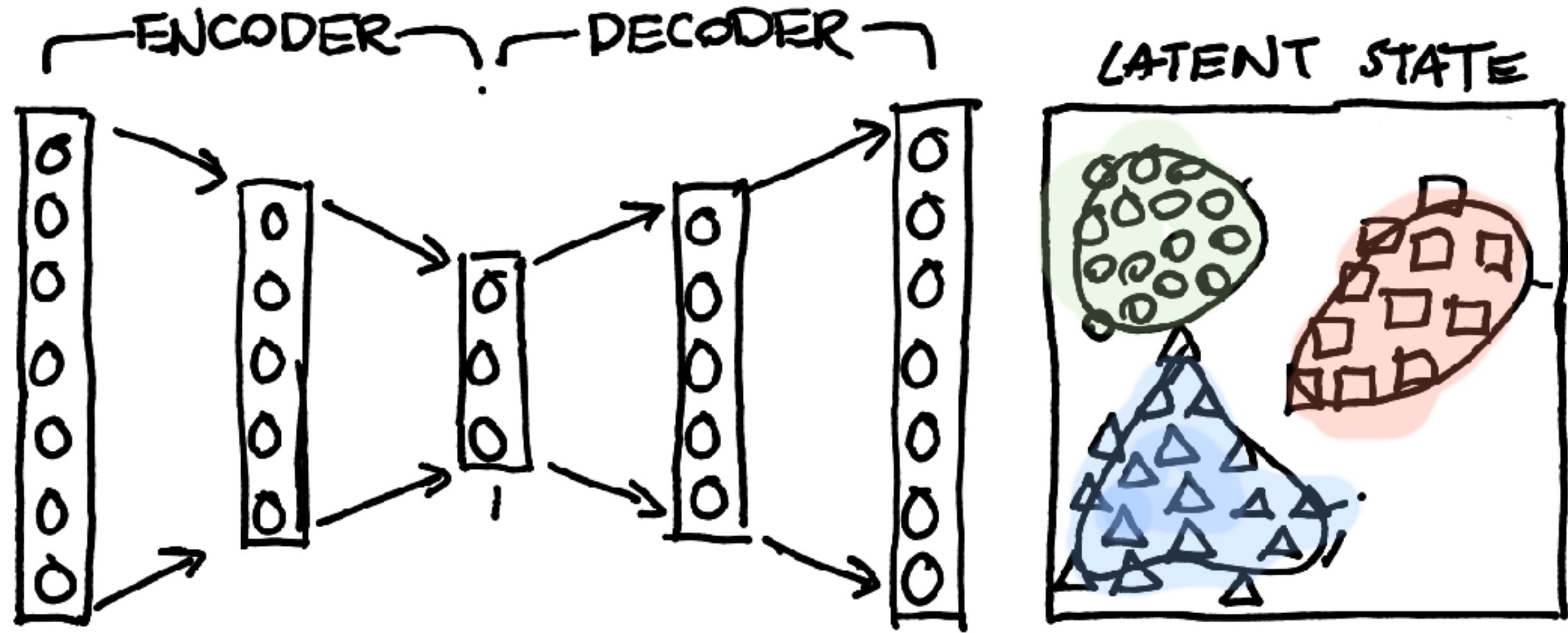
$$p(\mathbf{y}|\mathbf{x}) = \sum_{i=1}^k \pi_i \times N(\mu_i, \sigma_i)$$

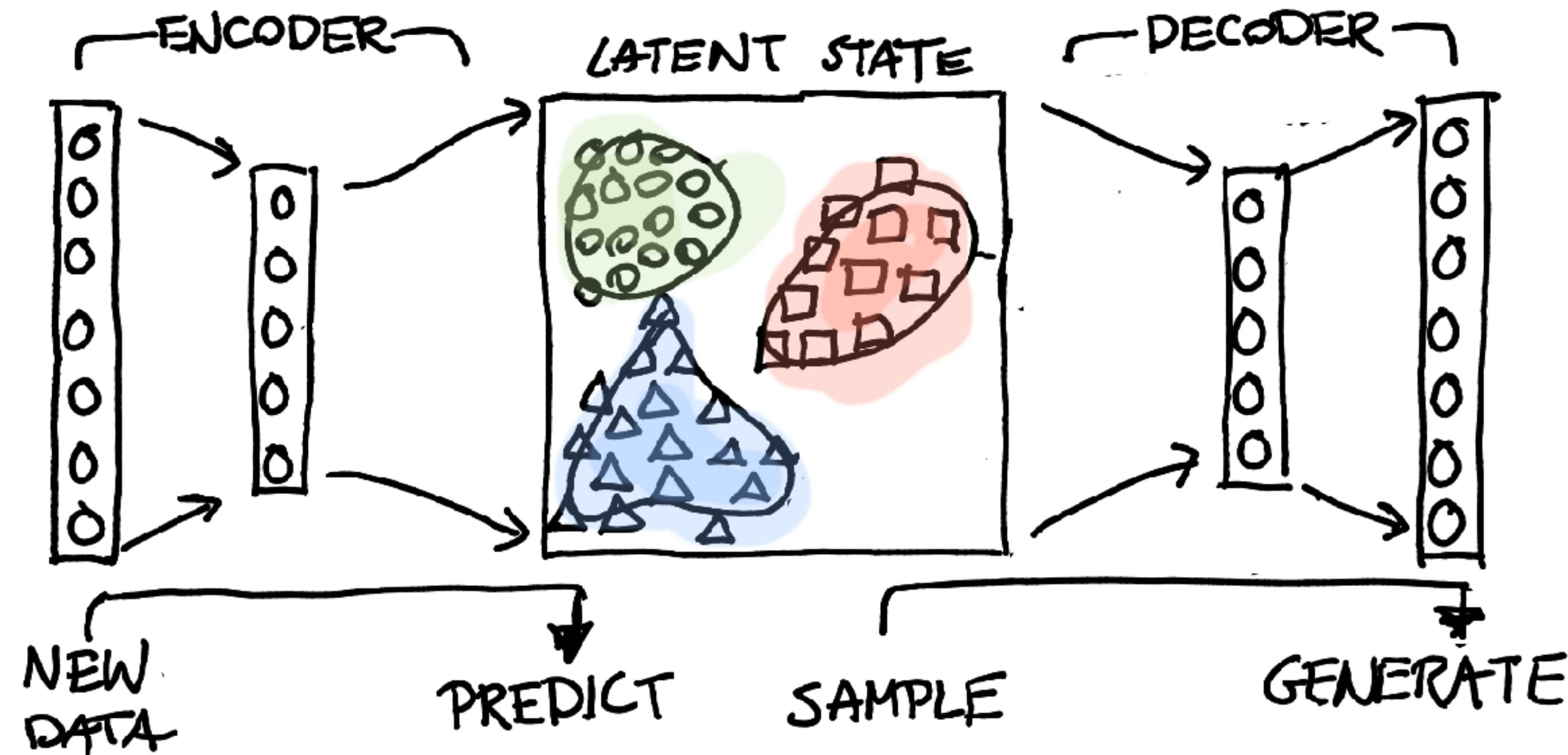
- a μ layer with nodes $\{\mu_1 \dots \mu_k\}$ which will denote the mean of a normal distribution
- a σ layer with nodes $\{\sigma_1 \dots \sigma_k\}$ which will denote the deviation of a normal distribution
- a π layer with nodes $\{\pi_1 \dots \pi_k\}$ which will denote the weight of the associated normal distribution in the resulting prediction

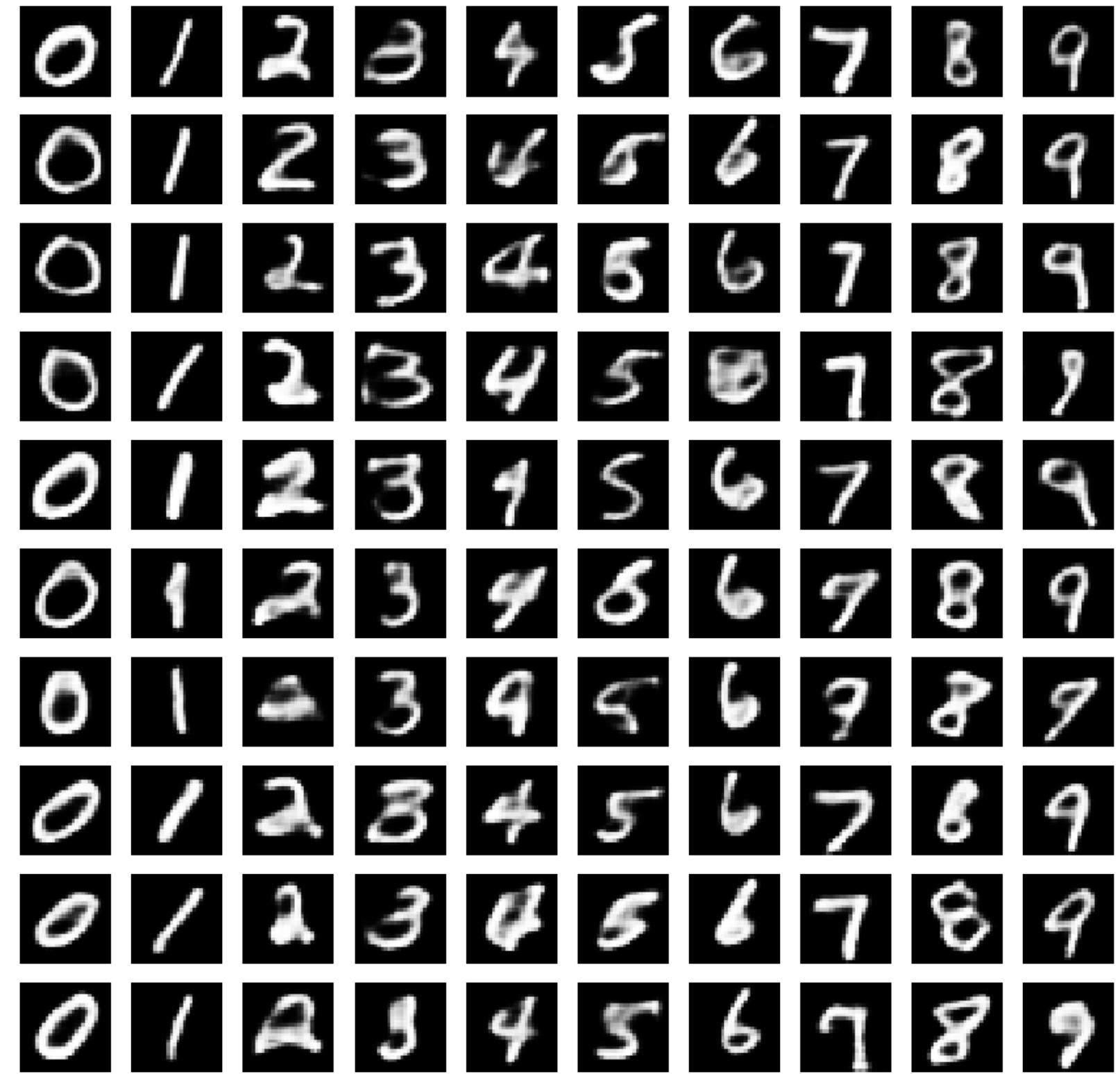
But what about deep learning?

But we can also add different properties.

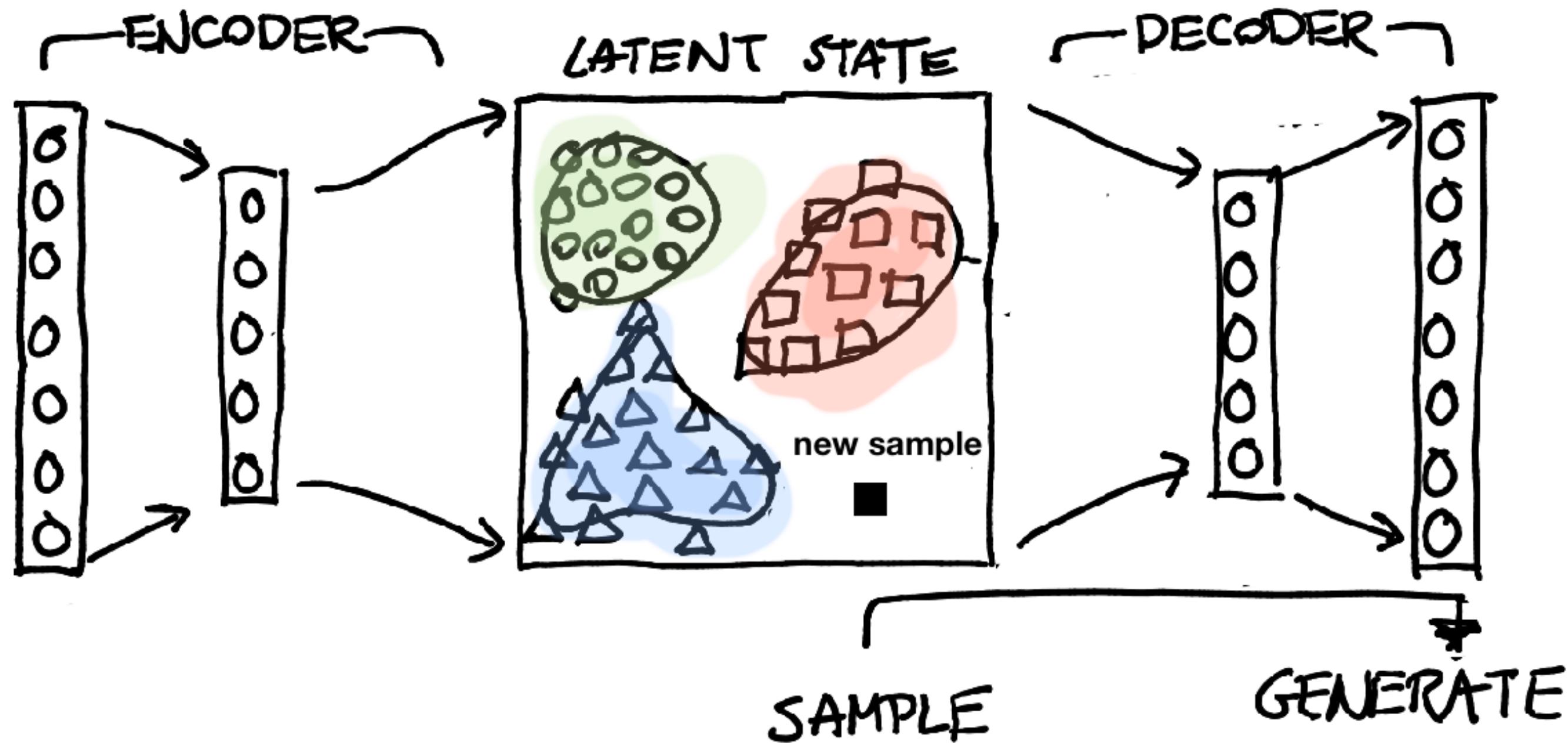


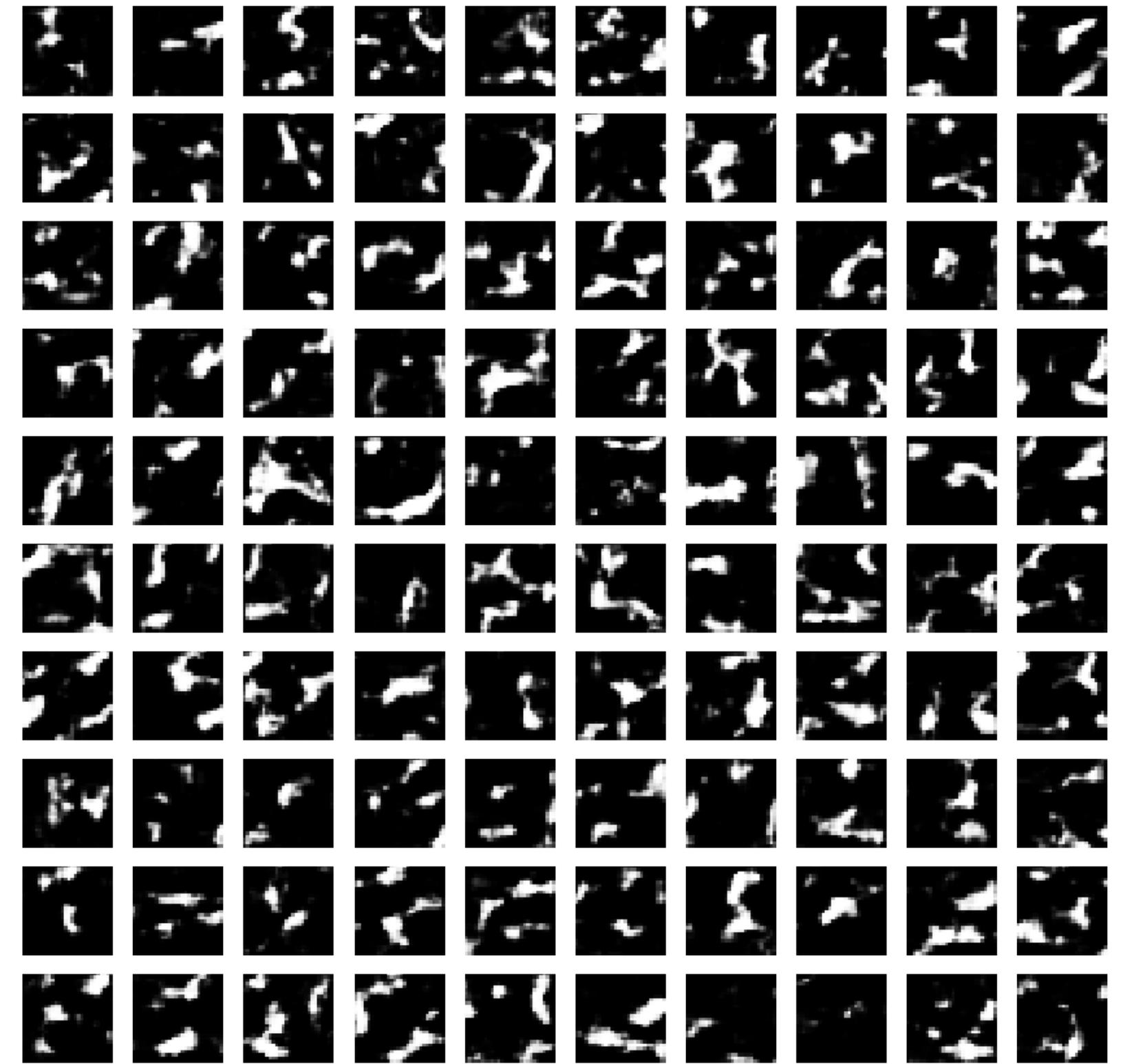




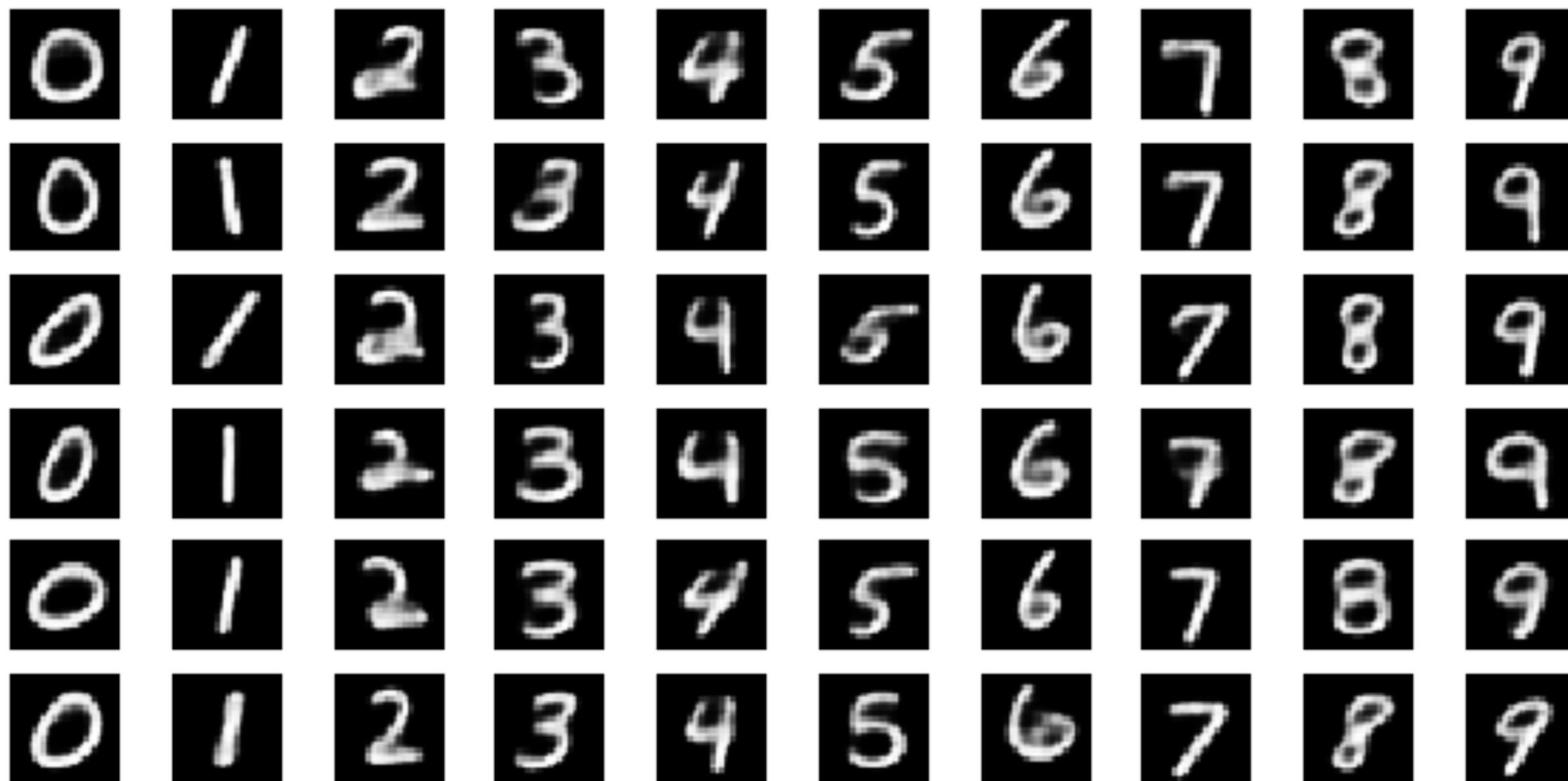








Articulating Styles

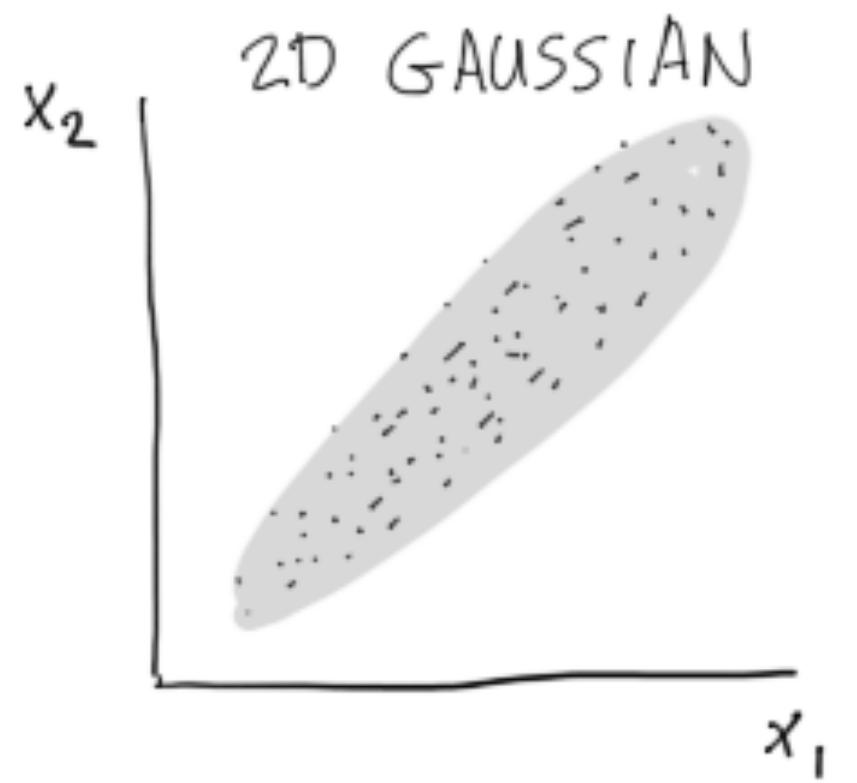


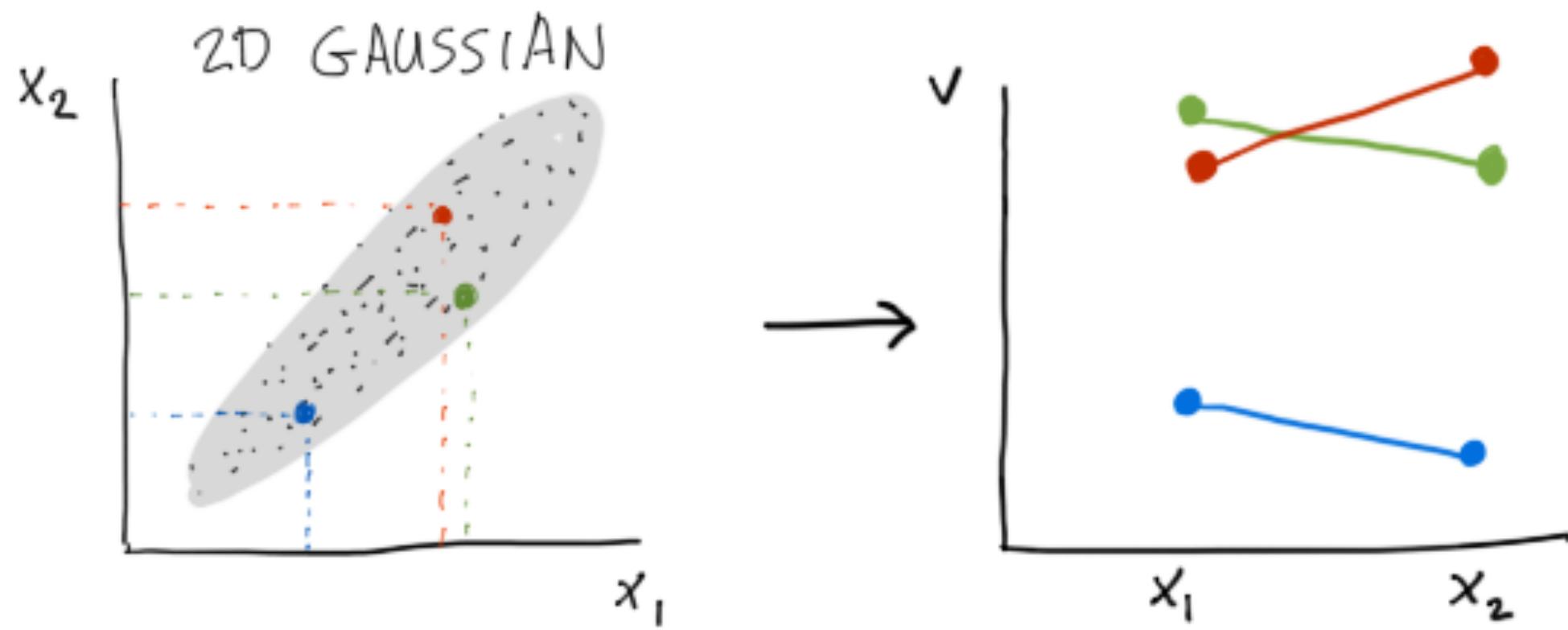
Gaussian Processes

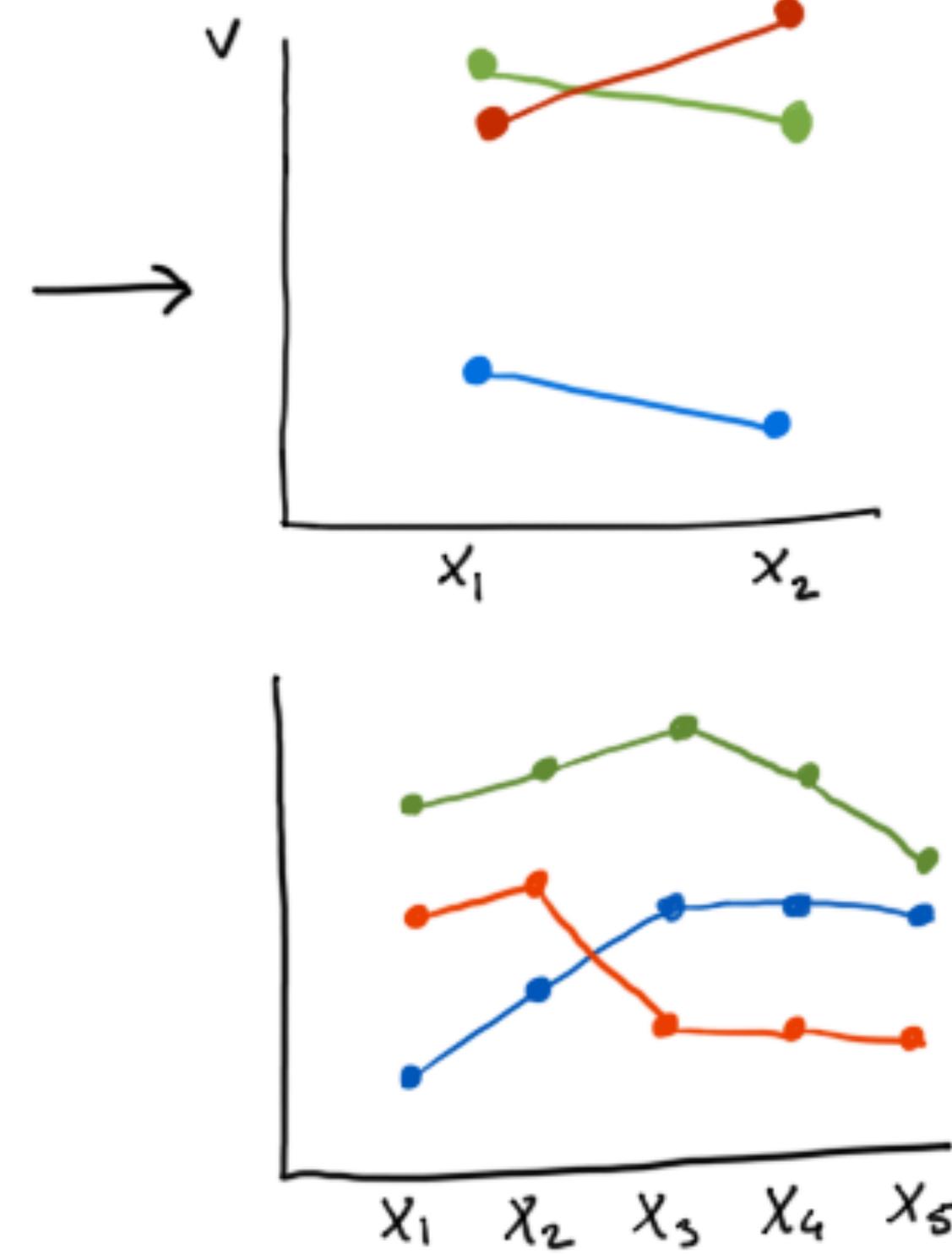
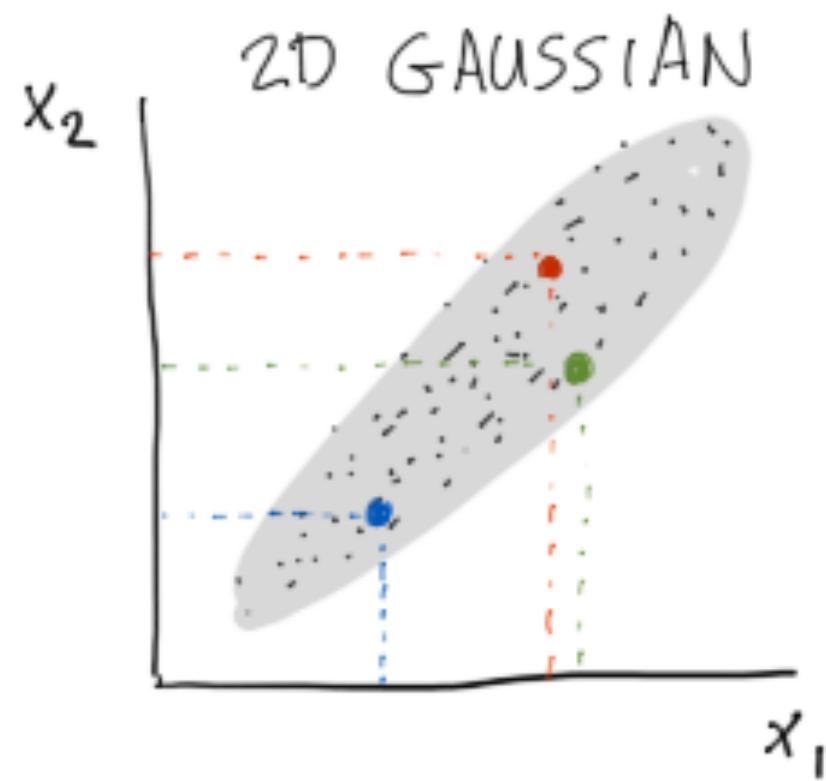
We're nearing the live coding part.

But I am going to attempt to get some intuition in you before I do.

If you've got some coffee, now is the time to drink it.
We need to increase the dimensionality to understand this next part.



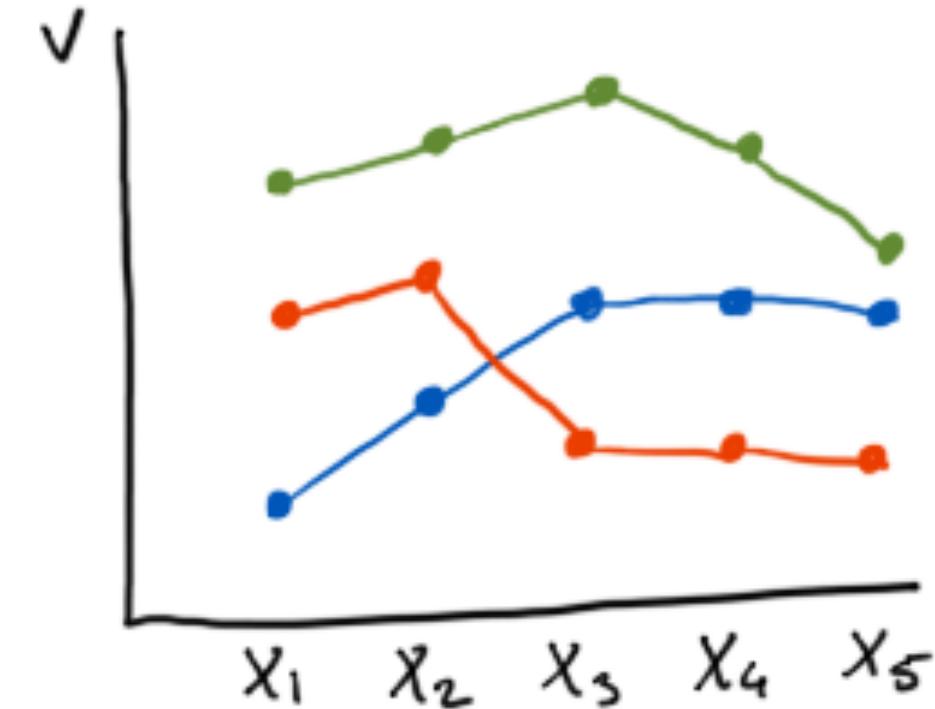




5D GAUSSIAN

$$p(x) \sim N(\mu, \Sigma)$$

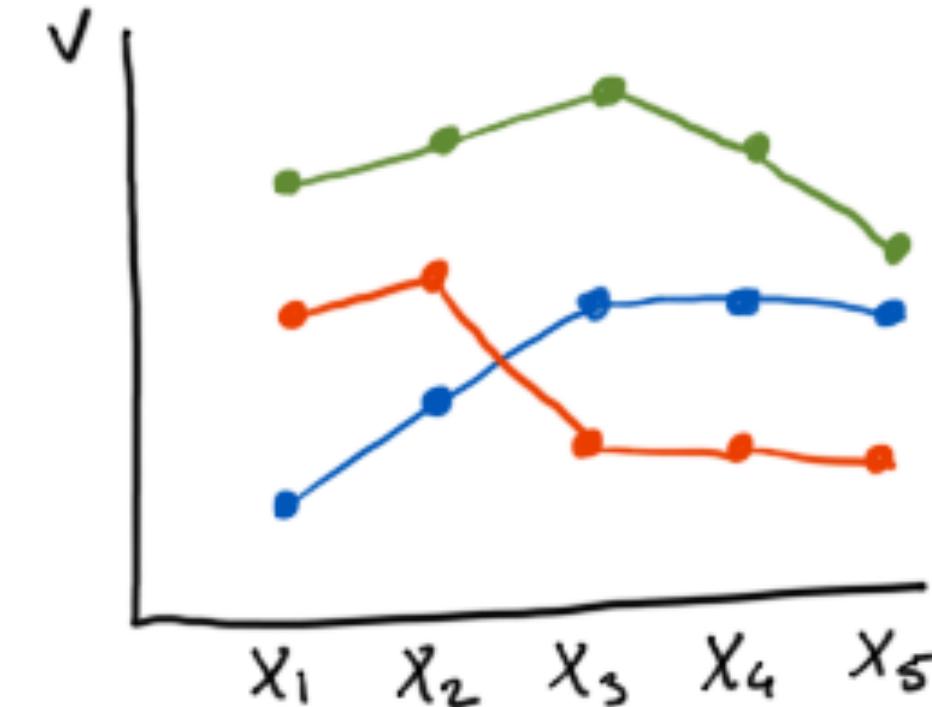
$$\Sigma = \begin{bmatrix} & 5 \times 5 \end{bmatrix}$$



$$p(x) \sim N(\mu, \Sigma)$$

$$\Sigma = \begin{bmatrix} & 5 \times 5 \\ & \end{bmatrix}$$

$$\mathcal{D} = \begin{bmatrix} 0 & 1 & 2 & 3 & 4 \\ 1 & 0 & 1 & 2 & 3 \\ 2 & 1 & 0 & 1 & 2 \\ 3 & 2 & 1 & 0 & 1 \\ 4 & 3 & 2 & 1 & 0 \end{bmatrix}$$

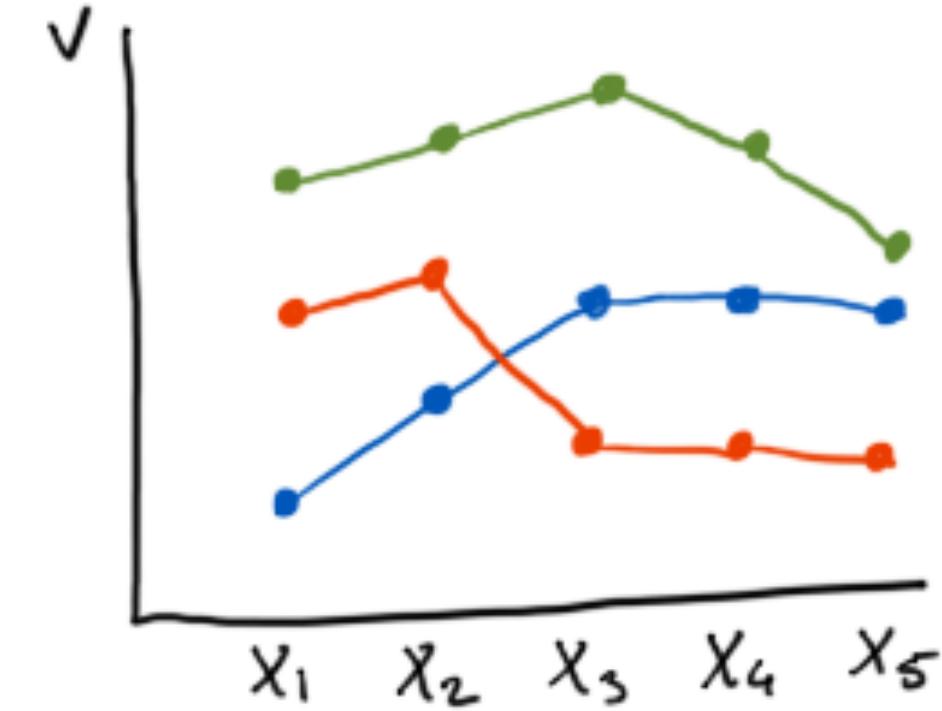


Idea : maybe the covariance between points should depend on the distance between them ?

$$p(x) \sim N(\mu, \Sigma)$$

$$\Sigma = \begin{bmatrix} & 5 \times 5 \\ & \end{bmatrix}$$

$$\mathcal{D} = \begin{bmatrix} 0 & 1 & 2 & 3 & 4 \\ 1 & 0 & 1 & 2 & 3 \\ 2 & 1 & 0 & 1 & 2 \\ 3 & 2 & 1 & 0 & 1 \\ 4 & 3 & 2 & 1 & 0 \end{bmatrix}$$



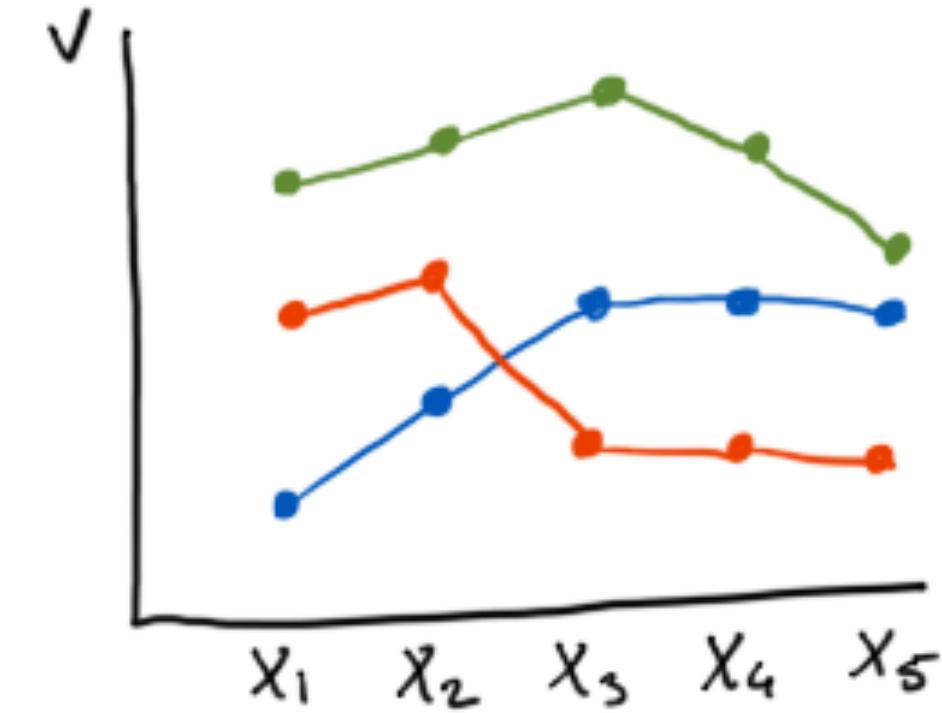
$$f(d) =$$

↑ Gaussian! if dist
is small then I'd
like covariance!

$$p(x) \sim N(\mu, \Sigma)$$

$$\Sigma = \begin{bmatrix} & & & f(2) \\ f(1) & & & \\ & & & \text{etc.} \\ & f(4) & & \end{bmatrix}$$

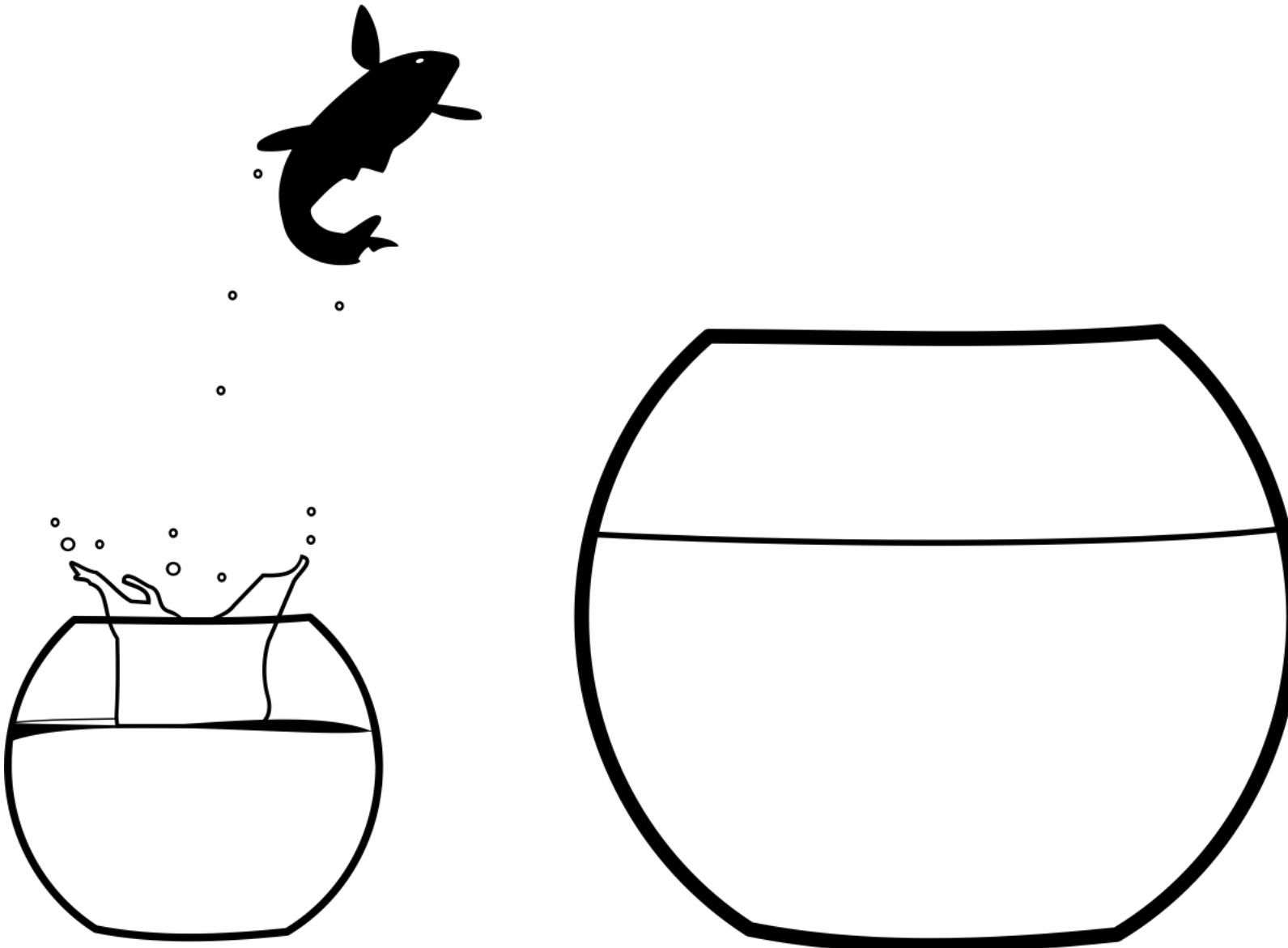
$$\mathcal{D} = \begin{bmatrix} 0 & 1 & 2 & 3 & 4 \\ 1 & 0 & 1 & 2 & 3 \\ 2 & 1 & 0 & 1 & 2 \\ 3 & 2 & 1 & 0 & 1 \\ 4 & 3 & 2 & 1 & 0 \end{bmatrix}$$



$$f(d) =$$

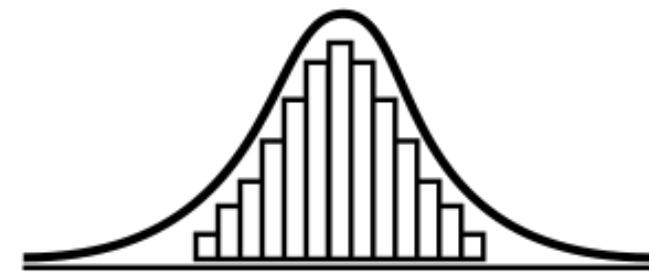
↑ Gaussian! if dist
is small then I'd
like covariance!

Live Coding Time



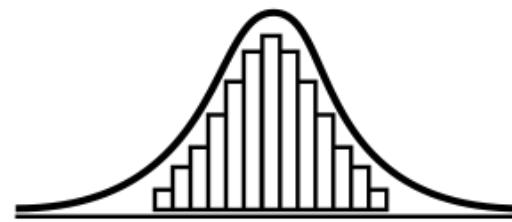
What did I just do?

Hopefully I was able to give some intuition on some things. The gaussian is everywhere in data science.



What did I just do?

It can be mind-blowing to witness the full applicability.
Understanding this mathematical lego brick allows you
to recognise algorithmic parts.



What did I just do?

If you still feel like you totally don't grasp it:
ThatsNormal[tm]. I've skipped details in favor of
intuition.

The best argument for this approach is considering the alternatives.

How to Demotivate: Wikipedia

Definition [edit]

A time continuous stochastic process $\{X_t; t \in T\}$ is Gaussian if and only if for every finite set of indices t_1, \dots, t_k in the index set T

$$\mathbf{X}_{t_1, \dots, t_k} = (X_{t_1}, \dots, X_{t_k})$$

is a multivariate Gaussian random variable.^[2] That is the same as saying every linear combination of $(X_{t_1}, \dots, X_{t_k})$ has a univariate normal (or Gaussian) distribution.

Using characteristic functions of random variables, the Gaussian property can be formulated as follows: $\{X_t; t \in T\}$ is Gaussian if and only if, for every finite set of indices t_1, \dots, t_k , there are real-valued $\sigma_{\ell j}, \mu_\ell$ with $\sigma_{jj} > 0$ such that the following equality holds for all $s_1, s_2, \dots, s_k \in \mathbb{R}$

$$E\left(\exp\left(i \sum_{\ell=1}^k s_\ell \mathbf{X}_{t_\ell}\right)\right) = \exp\left(-\frac{1}{2} \sum_{\ell,j} \sigma_{\ell j} s_\ell s_j + i \sum_\ell \mu_\ell s_\ell\right).$$

where i denotes the imaginary unit such that $i^2 = -1$.

The numbers $\sigma_{\ell j}$ and μ_ℓ can be shown to be the covariances and means of the variables in the process.^[3]

How to Demotivate: Books

An alternative and equivalent way of reaching identical results to the previous section is possible by considering inference directly in function space. We use a Gaussian process (GP) to describe a distribution over functions. Formally:

Definition 2.1 A Gaussian process is a collection of random variables, any finite number of which have a joint Gaussian distribution. □

Gaussian process

A Gaussian process is completely specified by its mean function and covariance function. We define mean function $m(\mathbf{x})$ and the covariance function $k(\mathbf{x}, \mathbf{x}')$ of a real process $f(\mathbf{x})$ as

$$\begin{aligned} m(\mathbf{x}) &= \mathbb{E}[f(\mathbf{x})], \\ k(\mathbf{x}, \mathbf{x}') &= \mathbb{E}[(f(\mathbf{x}) - m(\mathbf{x}))(f(\mathbf{x}') - m(\mathbf{x}'))], \end{aligned} \tag{2.13}$$

covariance and
mean function

and will write the Gaussian process as

$$f(\mathbf{x}) \sim \mathcal{GP}(m(\mathbf{x}), k(\mathbf{x}, \mathbf{x}')). \tag{2.14}$$

Usually, for notational simplicity we will take the mean function to be zero, although this need not be done, see section 2.7.

How to Demotivate: Paper

Algorithm 1 Infinite-horizon Gaussian process (IHGP) inference. The GP prior is specified in terms of a state space model. After the setup cost on line 2, all operations are at most $\mathcal{O}(m^2)$.

```

1: Input:  $\{y_i\}, \{\mathbf{A}, \mathbf{Q}, \mathbf{h}, \mathbf{P}_0\}, p(y | f)$  targets, model, likelihood
2: Set up  $\mathbf{P}^p(\gamma)$ ,  $\mathbf{P}^s(\gamma)$ , and  $\mathbf{G}(\gamma)$  for  $\gamma_{1:K}$  solve DAREs for a set of likelihood variances, cost  $\mathcal{O}(K m^3)$ 
3:  $\mathbf{m}_0^f \leftarrow \mathbf{0}; \quad \mathbf{P}_0^p \leftarrow \mathbf{P}_0; \quad \gamma_0 = \infty$  initialize
4: for  $i = 1$  to  $n$  do
5:   Evaluate  $\mathbf{P}_i^p \leftarrow \mathbf{P}^p(\gamma_{i-1})$  find predictive covariance
6:    $\tilde{\mu}_{f,i} \leftarrow \mathbf{h}^\top \mathbf{A} \mathbf{m}_{i-1}^f; \quad \tilde{\sigma}_{f,i}^2 = \mathbf{h}^\top \mathbf{P}_i^p \mathbf{h}$  latent
7:   if Gaussian likelihood then
8:      $\eta_i \leftarrow y_i; \quad \gamma_i \leftarrow \sigma_{n,i}^2$  if  $\sigma_{n,i}^2 := \sigma_n^2$ ,  $\mathbf{k}_i$  and  $\mathbf{P}_i^f$  become time-invariant
9:   else
10:    Match  $\exp(\nu_i f_i - \tau_i f_i^2/2) N(f_i | \tilde{\mu}_{f,i}, \tilde{\sigma}_{f,i}^2) \stackrel{\text{mom}}{=} p(y_i | f_i) N(f_i | \tilde{\mu}_{f,i}, \tilde{\sigma}_{f,i}^2)$  match moments
11:     $\eta_i \leftarrow \nu_i / \tau_i; \quad \gamma_i \leftarrow \tau_i^{-1}$  equivalent update
12:   end if
13:    $\mathbf{k}_i \leftarrow \mathbf{P}_i^p \mathbf{h} / (\tilde{\sigma}_{f,i}^2 + \gamma_i)$  gain
14:    $\mathbf{m}_i^f \leftarrow (\mathbf{A} - \mathbf{k}_i \mathbf{h}^\top \mathbf{A}) \mathbf{m}_{i-1}^f + \mathbf{k}_i \eta_i; \quad \mathbf{P}_i^f \leftarrow \mathbf{P}_i^p - \mathbf{k}_i \gamma_i \mathbf{k}_i^\top$  mean and covariance
15: end for
16:  $\mathbf{m}_n^s \leftarrow \mathbf{m}_n^f; \quad \mathbf{P}_n^s \leftarrow \mathbf{P}^s(\gamma_n)$  initialize backward pass
17: for  $i = n - 1$  to  $1$  do
18:    $\mathbf{m}_i^s \leftarrow \mathbf{m}_i^f + \mathbf{G}(\gamma_i) (\mathbf{m}_{i+1}^s - \mathbf{A} \mathbf{m}_i^f); \quad \mathbf{P}_i^s \leftarrow \mathbf{P}^s(\gamma_i)$  mean and covariance
19: end for
20: Return:  $\mu_{f,i} = \mathbf{h}^\top \mathbf{m}_i^s, \sigma_{f,i}^2 = \mathbf{h}^\top \mathbf{P}_i^s \mathbf{h}, \log p(\mathbf{y})$  mean, variance, evidence

```

What Might Help

Mathematics is useful when there's context and intuition. It's downright terrible when it's just a bunch of symbols dumped on your retina.

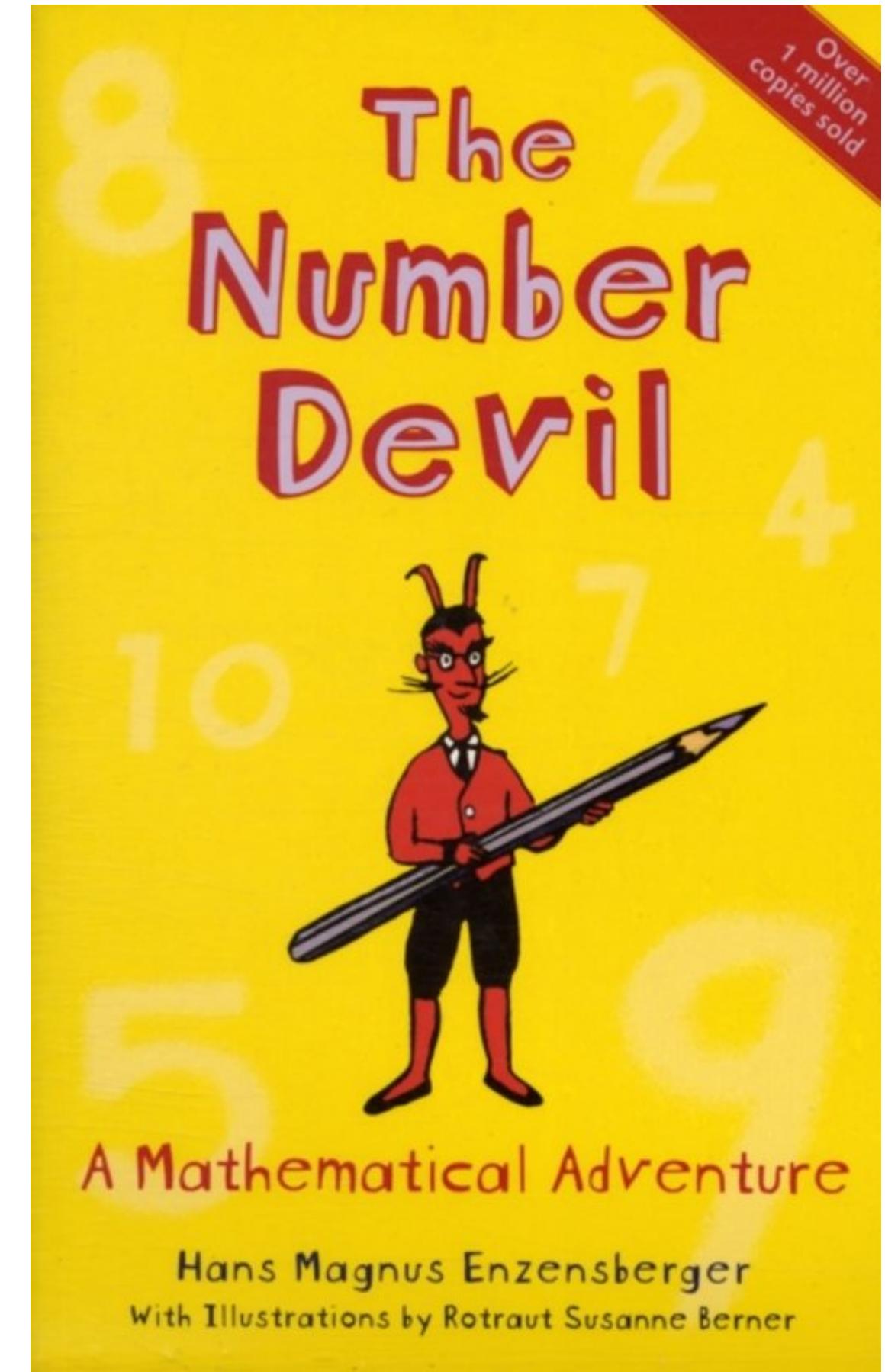
It's even worse if this occurs to you when you're young. I was lucky as a kid because my parents found a way to prevent this intellectual paralysis to me.

My Favourite Book

Remember the story in the beginning of this talk?

This book is about a boy who is afraid of maths. But as luck would have it, the number devil appears in his dreams. It's fairy tales about maths.

All who read this book when they're young ended up getting the best grades in my



Thanks for Listening