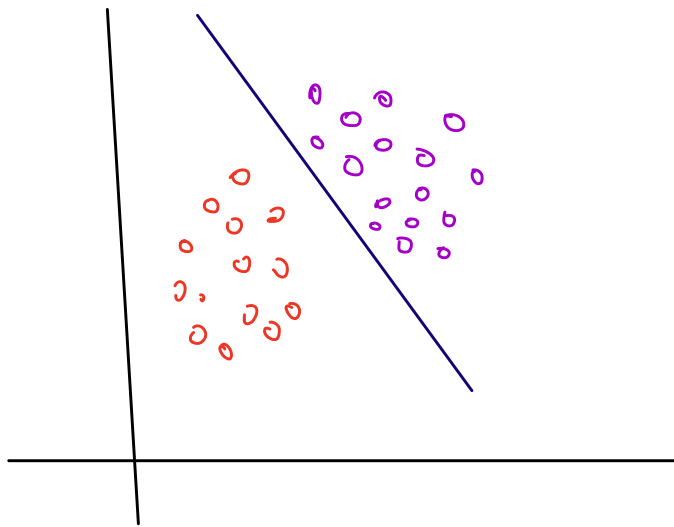


UNSUPERVISED LEARNING

$$S = \{x^{(1)}, x^{(2)}, \dots, x^{(n)}\}$$

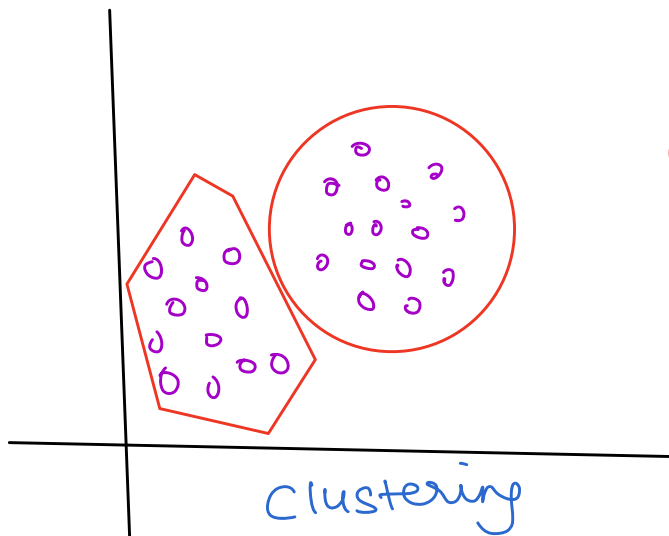
$x^{(1)} \in \mathbb{R}^d$



→ Logistic Regression

We had correct answers for each input for supervised

However in unsupervised setting



Here you have to find the interesting structure.

clustering

K-Means Algorithm

$$S = \{x^{(1)} \dots x^{(n)}\} \quad x^{(i)} \in \mathbb{R}^d$$

k-clusters

1. initialize cluster centroids

$$\mu_1, \mu_2, \mu_3 \dots \mu_k \in \mathbb{R}^d \text{ randomly}$$

2. Repeat till convergence,

$$\text{For every } i, \text{ set } c^{(i)} = \arg \min_j \|x^{(i)} - \mu_j\|_2^2$$

$$\text{For every } j, \text{ set } \mu_j = \frac{\sum_{i=1}^n 1\{c^{(i)} = j\} \cdot x^{(i)}}{\sum_{i=1}^n 1\{c^{(i)} = j\}}$$

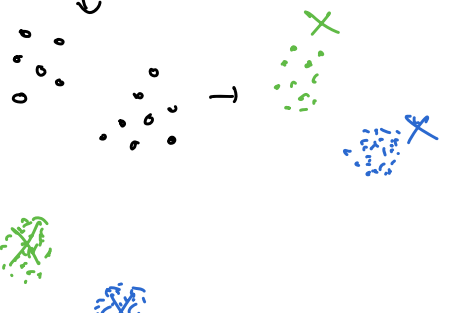
→ Initialize centroids randomly

→ c is array of length n

→ Set to identity of nearest means → $\arg \min$

→ μ_j means of all $x^{(i)}$'s for which $c^{(i)} = j$

can be thought
as labelling
clusters.



After labelling
we find μ_i

and readjust it
such that it
is at center of
clusters.

20

$$J(C, \mu) = \sum_{i=1}^n \|x^{(i)} - \mu_{c^{(i)}}\|_2^2$$

↓

Distortion Function

We use Co-ordinate Descent. It is a variant of gradient descent wherein we minimize loss wrt few variables holding others constant

One step optimize wrt C

Second step optimize wrt μ

$J \rightarrow$ non-convex, end up with different solⁿs depending on initialization but we always get a solⁿ i.e. it converges.

(Q) Why use labelled identities, when we have K-Means?

\rightarrow With K Means, you end up with diff solⁿs i.e. different cluster identities depending on initialization

* K-value has to be chosen by us based on domain knowledge, we can use Elbow method

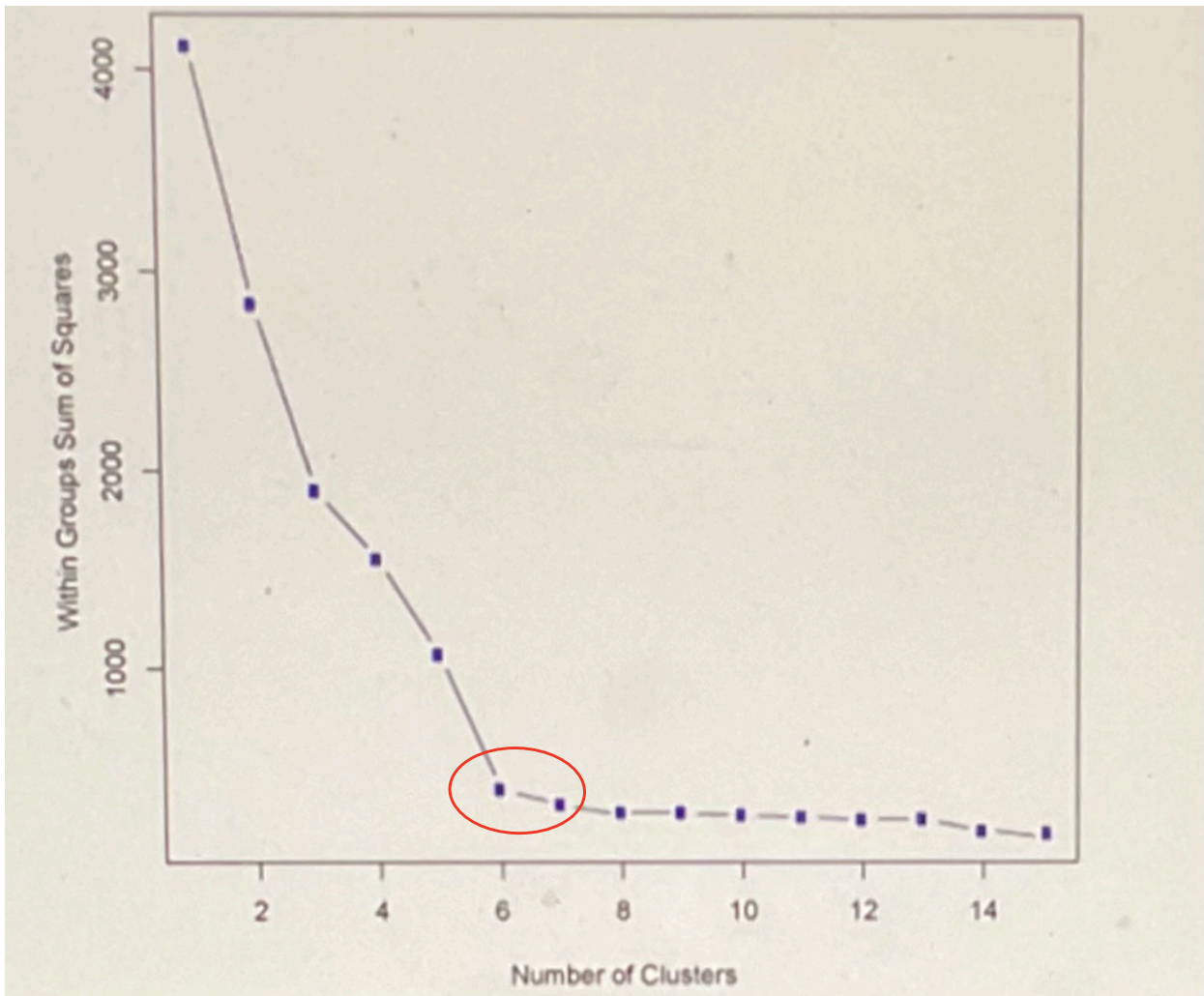
Elbow Method: Choose k at which SSE decreases abruptly.

We compute SSE i.e. sum of squared distance between each member of cluster and its centroid.

Do it for $k=2, 4, 6, \dots$

If we plot k against SSE, you will see error \downarrow as k gets larger as cluster increase in number & \downarrow in size.

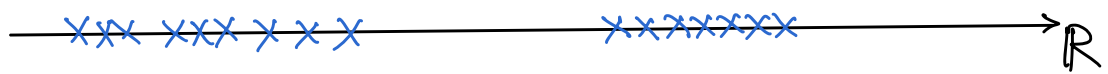
Choose a k at which SSE decreases abruptly



Choose k , where on T_k , you don't minimize within group S.S.E.

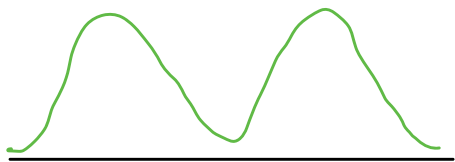
In this example, elbow occurs around 6-7

Density Estimation

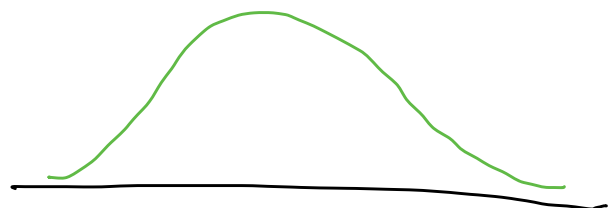


These points are sampled from a probability distribution
Since on \mathbb{R} line, it is sampled from P.D.F.

We have to estimate the PDF (To fully
fit we would have to use dirac delta)



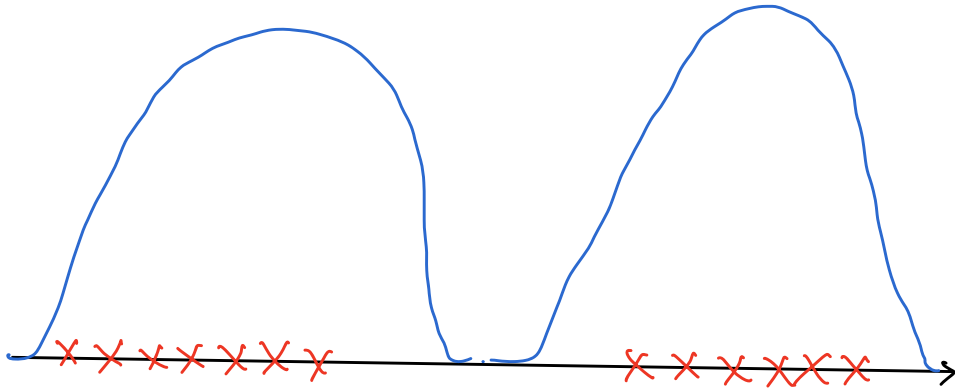
OR-



... Many possibilities

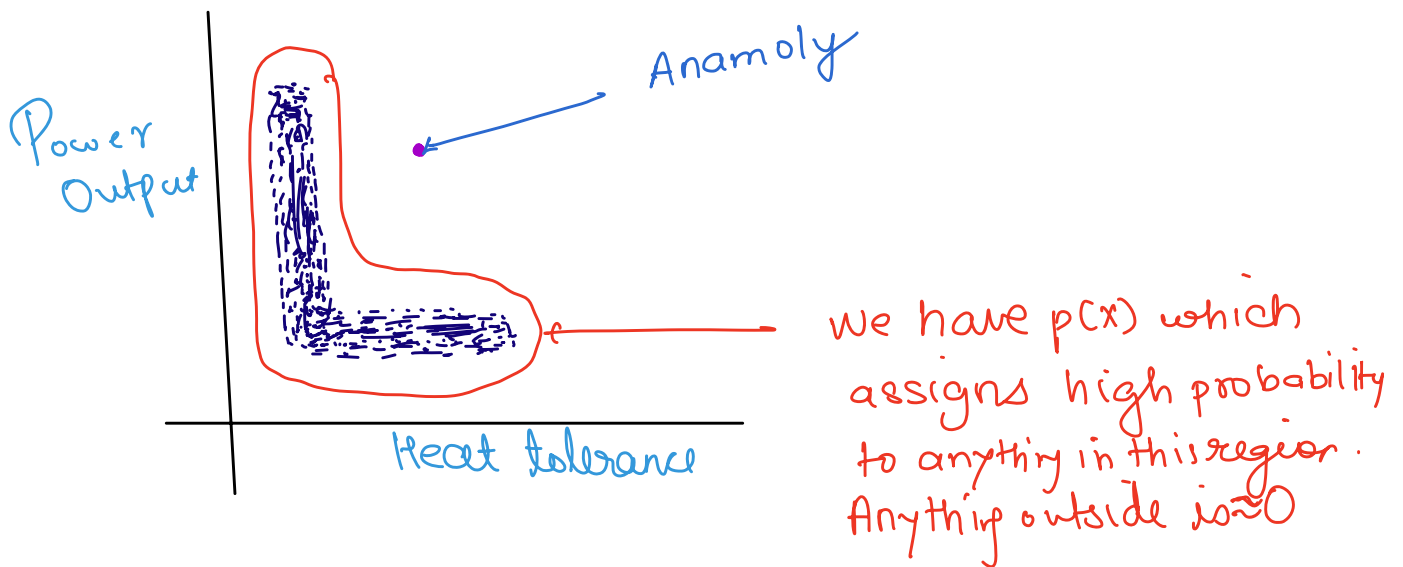
We have to estimate a smooth curve over a fixed
number of observations.

Gaussian Mixture Models



We assume above points sampled from 2 Gaussians
This is in unsupervised setting.

In Supervised setting, we had G.D.A. Here, we don't have the y labels and covariance can be different as opposed to G.D.A where it was same
Application: Anomaly Detection



Choice of k is once again on user

GMM \rightarrow Soft K-Means

$$S = \{x^{(1)}, x^{(2)}, x^{(3)}, x^{(4)}, x^{(5)} \dots x^{(n)}\}$$

$z^{(i)} \sim \text{Multinomial}(\phi) \leftarrow \text{Class identity}$

\downarrow which one of K

$\phi_j \geq 0 \quad \sum_{j=1}^K \phi_j = 1$

$\phi_j = p(z^{(i)} = j)$

$\xrightarrow{\text{Class prior tells how many examples of cluster } j}$

$x^{(i)} | z^{(i)} = j \sim N(\mu_j, \Sigma_j)$

$z^{(i)}$ similar to $y^{(i)}$ in G-D-A

- latent variable as it is not observed

$$\begin{aligned} \log p(x; \phi, \mu, \Sigma) &= \mathcal{L}(\mu, \Sigma, \phi) \\ &= \log \sum_z p(x, z; \phi, \mu, \Sigma) \end{aligned}$$

$p(z) \rightarrow$ Class prior

$p(x, z) \rightarrow$ model

$p(z|x) \rightarrow$ posterior

$p(x) \rightarrow$ evidence

$z \rightarrow$ latent variable

For each point assign a weight to cluster centroid i.e. posterior distribution of $p(z)$

[Inspired by K-Means]

Randomly initialize μ, ϕ, Σ

Repeat until convergence

E-Step: For each i, j set:

$$\omega_j^{(i)} = p(z^{(i)} = j | x^{(i)}; \phi, \mu, \Sigma)$$

Use Bayes Rule

M-Step: Update Parameters

$$\phi_j = \frac{1}{n} \sum_{i=1}^n \omega_j^{(i)}$$

$$\mu_j = \frac{\sum_{i=1}^n \omega_j^{(i)} x^{(i)}}{\sum_{i=1}^n \omega_j^{(i)}} \quad \leftarrow \text{every } x \text{ contributes}$$

$$\Sigma_j = \frac{\sum_{i=1}^n \omega_j^{(i)} (x^{(i)} - \mu_j)(x^{(i)} - \mu_j)^T}{\sum_{i=1}^n \omega_j^{(i)}}$$

eg.

$K=3$

$$p(z^{(i)}=j | x^{(i)})$$

$$\begin{bmatrix} 0.1 \\ 0.7 \\ 0.2 \end{bmatrix} \begin{matrix} K=1 \\ K=2 \\ K=3 \end{matrix}$$

} K-Means

$$\begin{bmatrix} 0 \\ 1 \\ 6 \end{bmatrix}$$

Here we do a soft assignment instead of K-means where there is hard assignment.

Gaussian

Multinomial

$$p(Z|X) = \frac{\overbrace{p(X|Z)}^{\text{Gaussian}} \overbrace{p(Z)}^{\text{Multinomial}}}{\sum_Z p(X|Z) p(Z)}$$

↳ used in E-step

↓
[You get softmax with quadratic features]

GDA → if covariances same, you get logistic
if covariances diff, you get curved.

Here diff covariances, so curved, hence softmax with quadratic features