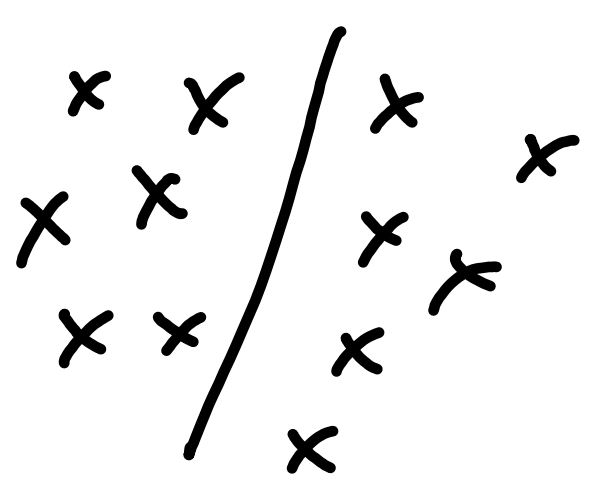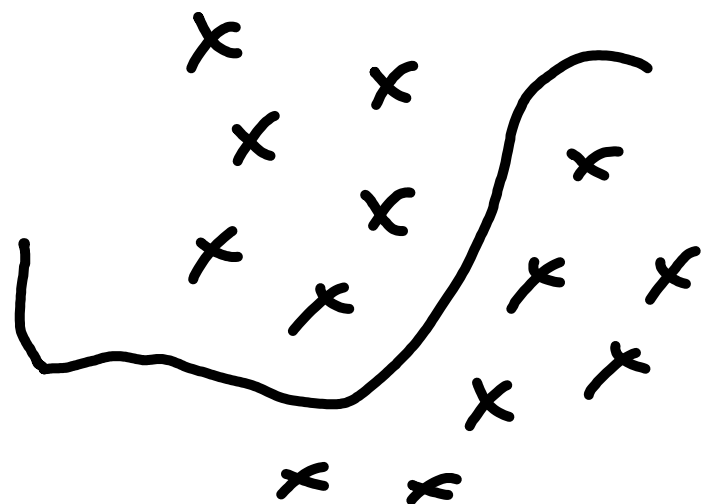# Clustering

In supervised learning ( given colour of pts ):



linear separ.

two-mean
(ML lay cond.)

We can do some kind of such classification.

Next, we consider unsupervised learning and do some clustering algorithm.

The aim is to assign new test data pt $x \in X$ to one of $k$ clusters. i.e. Define:
$\hat{j} : x \in X \longmapsto \{1, 2, \ldots k\}$. So : $X = \sum_{i}^{k} \hat{j}^{-1}(i)$.

We will use center-based clustering. i.e. fix $k$ center pts $\{m_i\}_i^k \subseteq X$. $d$ is metric on $X$. Then : $\hat{j}(x) = \arg\min_{j \leq k} d(X, m_j)$.

Assume $X = \mathbb{R}^d$. Consider Gaussian mixture model $X = \{ dv(x) = \sum_{i}^{k} s_j f_j(x) dx$ . $f_j = (\sqrt{2\pi}\sigma)^{-d}$ $e^{-\frac{1}{2} \| (x - m_j) / \sigma \|^2}$ $m_j \in [a, b]^d$. $s_j \in [\varepsilon, 1]$. $\sum_{i}^{k} s_j = 1 \}$.

And its EMF $\hat{\mathcal{L}}(v, \mathcal{X}_m) = -\sum_{t=1}^{m} \log(\sum_{j=1}^{k} f_j(x_t) g_j)$

$= -\sum_{z=1}^{m} \log(\sum_{j=1}^{k} g_j \exp(-\frac{1}{2}\|x_t - m_j\|_2^2)) + m\frac{d}{2}\log(2\pi\sigma)$

Consider $\textcircled{w} = \{(g, m) \mid g \in [\varepsilon,1]^k, m \in \mathbb{R}^{d \times k}\}$.

But this model isn't identifiable by:

$\mathcal{L} v_{\pi} = \sum_{i}^{k} g_{\pi(j)} N(\cdot | M_{\pi(j)}, \sigma^2 I) = \mathcal{L} v$, for $\forall$

$\pi \in S_k$. (k-permutation)

$J_1: \theta \in \textcircled{w} \longmapsto \mathcal{V}_\theta \in \mathcal{H}$ isn't injective.

$\Rightarrow \textcircled{w}$ isn't real parametric model.

$\underline{Rmp}$: We can use EM algo. to find EMR for $\hat{\mathcal{L}}_n$. But it's not unique here.

We will rebuild the parameter space:

$\overline{\textcircled{n}} := \{(g, m) \subset ([\varepsilon,1]^k \cap \{\mathbb{1}^T = 1\}) \times [a,b]^{d \times k}:$

$m_{1,1} \leq m_{2,1} - \varepsilon \leq m_{3,1} - \varepsilon \leq \cdots \leq m_{k,1} - \varepsilon\}$. i.e. adding

some order into $\textcircled{w}$. $\Rightarrow \overline{\textcircled{n}}$ is cpt and $\theta$

$\in \overline{\textcircled{n}} \longmapsto \mathcal{V}_\theta$ is injective

$\underline{prop}$. For para. model $\overline{\textcircled{n}}$ induced above. For $v$

$= \mathcal{V}_{\theta_0} \in Im(\mathcal{V}), \theta_0 \in \overline{\textcircled{n}}$. If $\hat{\theta}_n$ is $\hat{\mathcal{L}}_n$-MLE

i.i.d. $\hat{\theta}_n = \arg\min_{\bar{\theta}} \bar{I}_n(\nu_\theta))$ Then: $\hat{\theta}_n \xrightarrow{pr} \theta_0$

Pf: Since $\bar{\Theta}$ is cpt. We want to apply uniform LLN before.

Note $\ell(x|\theta) = -\log(\sum_j s_j e^{-\frac{1}{2}\|\frac{x-m_j}{\sigma}\|^2}) + const.$

is conti. in $(s, m)$.

It remains to prove: $|\ell(x|\theta)| \leq k(x) \in L'(\nu_{\theta_0})$

i) $\sum_j^k s_j e^{-\frac{1}{2}\square} \leq \sum s_j = 1.$

2) $\sum_j^k s_j e^{-\frac{1}{2}\square} \geq \sum \exp(-\frac{1}{2}\max_j \|\frac{x-m_j}{\sigma}\|^2)$

$\geq \sum \exp(-\frac{1}{2}\max_{[a,b]^d} \|X - c\|^2/\sigma^2)$

$\max_{[a,b]^d} \|X-c\|^2 \geq \|X-X_0\|^2 - const. \quad X_0 \in [a,b]^d.$

$\sum_j = \exists k(x) \in L'(\nu_{\theta_0}).$

Next, we see how the data space $X = \mathbb{R}^d$
be partitioned basing on GMM:

We define decision rule by maximum a-post erior (MAP) principle:

i) View data $X$ as $X$-valued r.v. following

c.d.f. $p(X|\theta)$. $\theta$ is $\mathcal{H}$-valued r.v.

ii) Model r.v. $J \in [1, \cdots k]$ which encodes the
affiliation of dist. of $X$. (label)

$$\Rightarrow p(X|\theta) = \sum_j^k p(X, j|\theta) = \sum_j^k p(X|j, \theta) p(J=j|\theta)$$

$$= \sum_j^k f_j(X|\theta) \, \xi_j$$

So: $p(j|X, \theta) = p(X|j, \theta) \cdot p(j|\theta) / p(X|\theta)$

$$= f_j(X|\theta) \, \xi_j / \sum_j^k \xi_j f_j(X|\theta).$$

MAP principle demands $X \in \mathcal{X}$ to be assigned

$\hat{j}$ if $\hat{j} = \arg\max_{j=1 \cdots k} \xi_j f_j(X|\theta)$.

Rmk: The MAP decision rule will converge

for $X_k \overset{i.i.d}{\sim} f_j$ if $\theta \mapsto p(\cdot|\theta)$ is conti.

by prop. above. $\forall j$.

For Gauss we maximize $\log(\xi_j f_j(X|\theta)) =$

$\log(\xi_j (\sqrt{2\pi} \sigma^2 j^k)) - \|X - m_j\|_2^2 / 2.$

When $\xi_j = 1/k$. $\forall j$. equal frequency. $\Rightarrow$ MAP

maximize $-\|X - m_j\|_2^2 / 2$ by choosing $j \in [1, \cdots k]$.

But in practical, we process as below:

a) Pick random center pts $m_j$, $j \leq k$.

b) Apply MAP rule to get clusters of data

c) Update centers by minimizing:

$$-\log \prod_{\hat{j}(t)=\ell}^{k} f(x_t \mid m) = -\frac{1}{2} \sum_{t:1, \hat{j}(t)=\ell}^{k} \left\| \frac{x_t - m}{\sigma} \right\|_2^2$$

$$\Rightarrow \text{choose } \tilde{m}_\ell = N_\ell^{-1} \sum_{t:1, \hat{j}(t)=\ell}^{n} x_t, \text{ where } N_\ell := $$

$$|\{ t : \hat{j}(t) = \ell, \ t = 1, 2, \cdots n \}|$$

d) Iterate by applying MAP basing on new centers $\{\hat{\tilde{m}}_\ell\}_1^k$.

e) Stop when no association change clusters.

Rmk: It's called k-means Algorithm. But it sometimes fails to converge.

(It also depends on initialization of $m_j$)