CSC591- HW1   (Group: ALL_12)
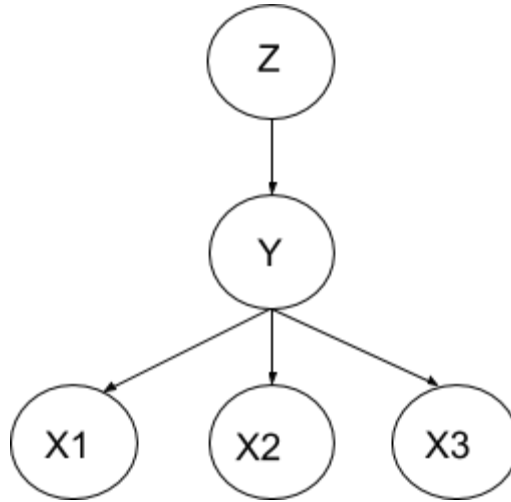Anusha (amanur)
Kartik (kcshah)
Shailaja (smallic)

1. With the given scenario,
   a. The decision tree may vary depending on the interpretation of the feature variable, if specifically applied to nuclear plant scenario. The tree below is a generalised version as per the given question. Making Z as the parent of Y makes sense as if the sensor i.e., the core temperature (as given by Z) is abnormal (high), then Y value can be predicted to malfunction. The feature variables are accordingly the child nodes of Y.



   b. $P(Y=1| Z= \text{missing}) = \dfrac{P(Z=missing|\ Y=1)\ P(Y=1)}{P(Z=missing\ |\ Y=1)\ P(Y=1) + P(Z=missing|\ Y=0)\ P(Y=0)}$

   By substituting the given values,

   $P(Y=1| Z= \text{missing}) = \dfrac{0.15\theta_{\ Y=1}}{0.15\theta_{\ Y=1} + 0.03(1-\theta_{\ Y=1})}$

   c. The likelihood is given by:
   $$P(X_1, X_2, X_3| \lambda\ ) = (1/2)^{20}\ (\lambda)^{24}\ (1 - 3\lambda)^4$$

   The log likelihood is given by:
   $$\text{Log } P(X_1, X_2, X_3| \lambda) = 20 \log \tfrac{1}{2} + 24 \log \lambda + 4 \log (1 - 3\lambda)$$

   In order to calculate the value of $\lambda$:

$\frac{d \log P(X_1,X_2,X_3|)}{d\lambda}$ = 24/ $\lambda$ + 4/( 1 − 3$\lambda$ )(-3) = 0

$\Rightarrow$   24/ $\lambda$ - 12/( 1 − 3$\lambda$ ) = 0

$\Rightarrow$  24( 1 − 3$\lambda$ ) -12 $\lambda$ = 0

$\Rightarrow$  $\lambda$ = 2/7

d. For single example, as all the variables are Bernoulli, (for example, in coin toss, it is $p^k(1-p)^{(1-k)}$ where *k* can take value of either 0,1 for head or tail), similarly this is applied to our case as follows-

e.

$$\text{E-step} - Q(y=1 \mid z, x)$$

$$P(y=1 \mid z, x) = \frac{P(y=1, zx \mid \theta)}{P(y=1, zx \mid \theta) + P(y=0, zx \mid \theta)}$$

$$= \frac{0.15\,\theta_{y=1}}{0.15\theta_{y=1} + 0.03(1-\theta_{y=0})}$$

M-step —

$$P(data \mid \theta) = \prod_{i=1}^{n} P(y^i \mid z^i) P(x_1^i \mid y^i) P(x_2^i \mid y^i) P(x_3^i \mid y^i) P(z^i)$$

$$\log P(data \mid \theta) = \sum_{i=1}^{n} \log P(y^i \mid z^i) + \log P(x_1^i \mid y^i) + \log P(x_2^i \mid y^i) + \log P(x_3^i \mid y^i) + \log P(z^i)$$

$$\frac{\partial \log P(data \mid \theta)}{\partial \theta_{y \mid z, x}} = \sum_{i=1}^{n} \frac{\partial \log P(y^i \mid z^i)}{\partial \theta_{y \mid z}}$$

$$\rightarrow \theta_{y \mid z, x} = \frac{\sum_{i=1}^{n} \sum_{y} \delta(z^i) \, Q(y^i = y \mid z, x)}{\sum_{i=1}^{n} \delta(z^i)}$$
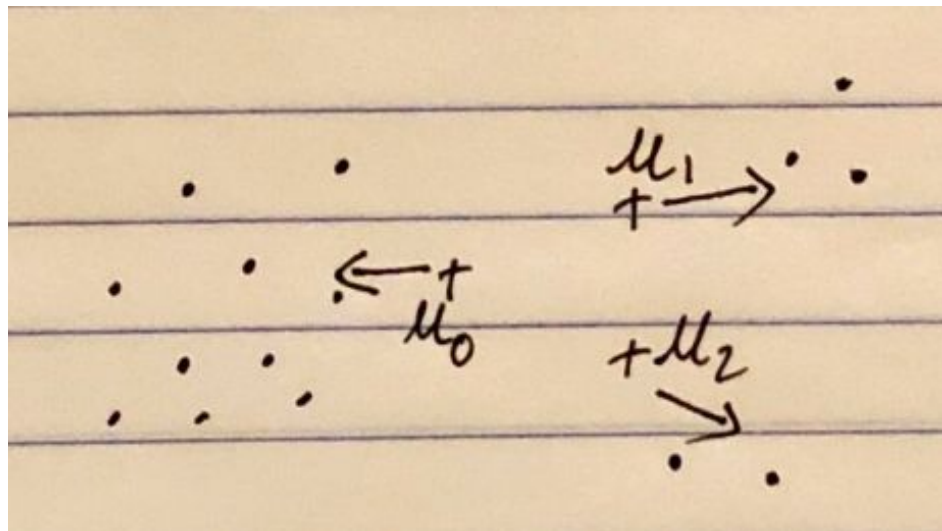
2. GMM
   a. No, the same cluster centers will not be produced as K-Means performs hard clustering, meaning, every datapoint is assigned to a single cluster. The K-Means cluster center is the means of the the data points belonging to that cluster.

Whereas GMM trained with EM performs soft clustering. In GMM, every data point is given a probability of the point belonging to every cluster. Each point is assigned to all clusters with some probability or weight. This affects the mean value of the cluster which is the center.
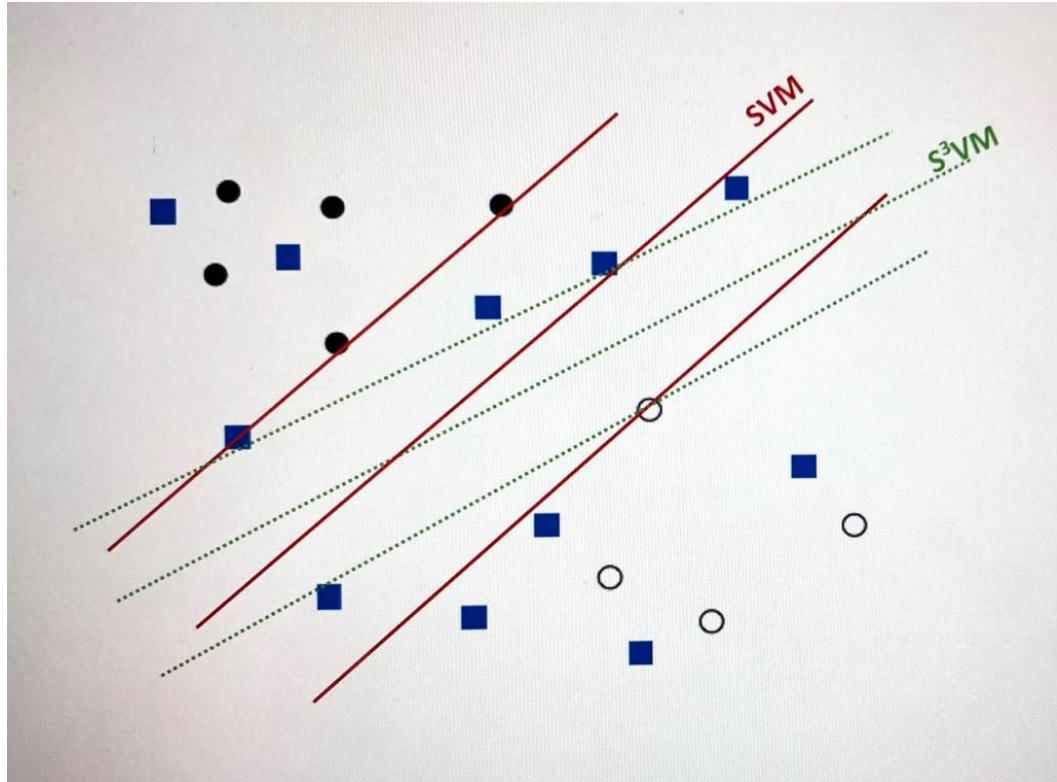
This difference in hard and soft clustering in the 2 methods leads to different cluster centers.

b. EM

   i.



   ii.    Marginal likelihood increases for every iteration of EM until it reaches local optimum.

   iii.   The estimate of $\pi_0$ will increase on the next EM iteration because each point is assigned to all the clusters and $\pi_0$ is calculated considering the sum of all the probabilities of all the points in cluster 1. Since, $\mu_1$ & $\mu_2$ are close to $\mu_0$ it wasn't getting a higher probability. So, in every EM iteration $\pi_0$ will increase.

3. With the given labeled and unlabeled data points,

   a. The marginal boundaries and separating hyper-plane for both *SVM* and $S^3VM$ are shown below:
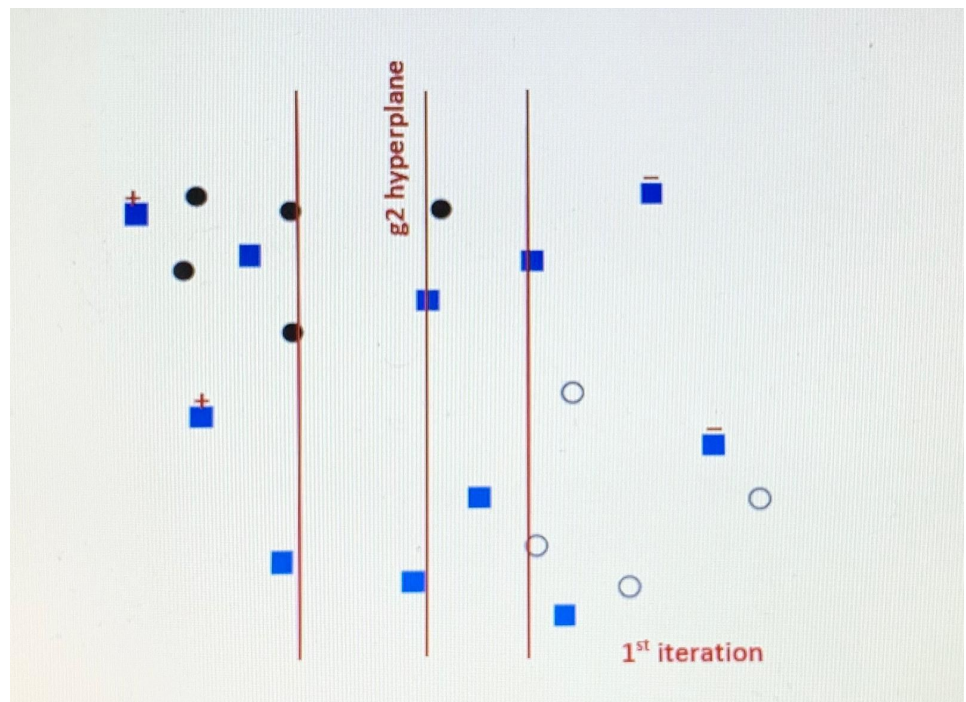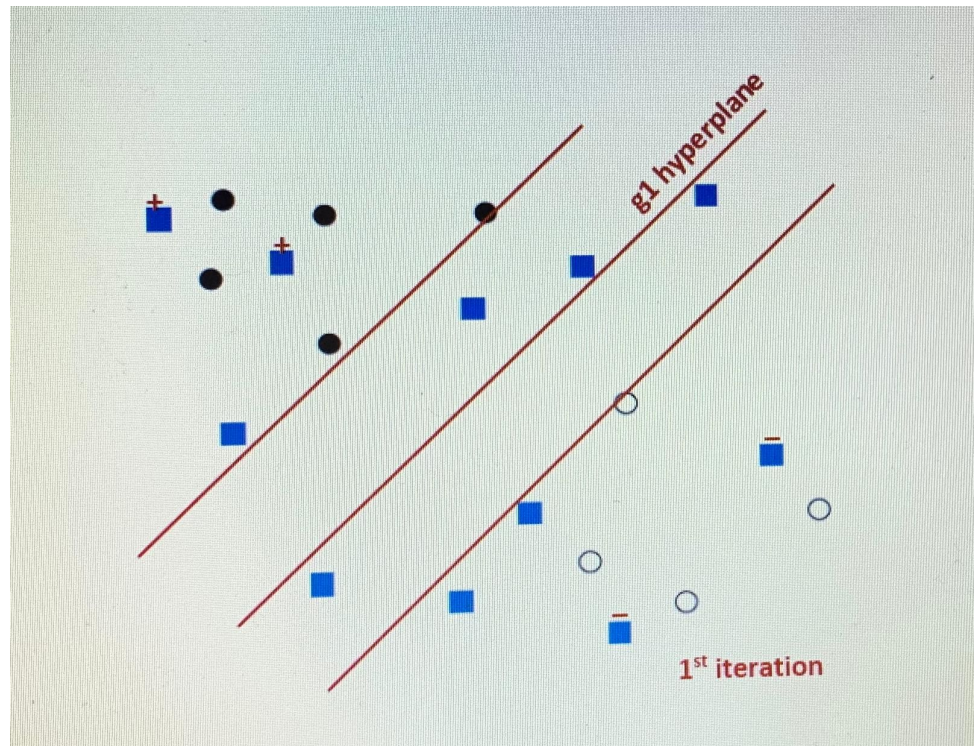
b.

i. Datasets whose features naturally partition into two sets, and algorithms that use this division, fall into the co-training and the main underlying assumptions are (1) each set of features is sufficient for classification, and (2) the two feature sets of each instance are conditionally independent and not redundant.
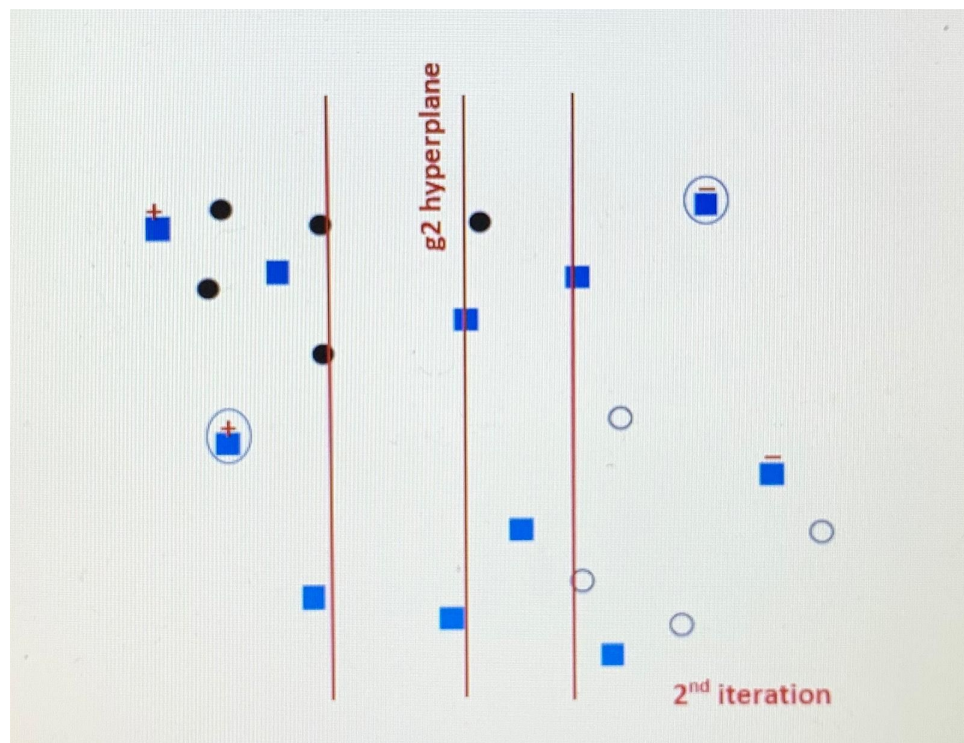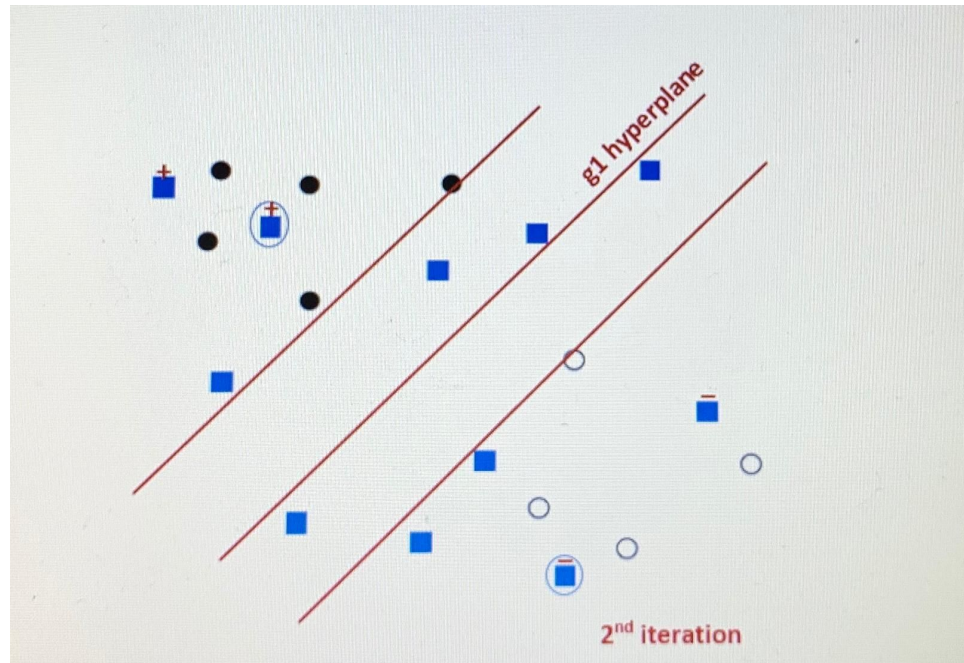
Idea is that if two algorithms use different representations of their hypothesis, they can learn two diverse models that can complement each other by labeling some unlabeled data & helping the training set of others. (Blum and Mitchel (1998)).

The two classifiers that can be used are SVM classifiers with different views i.e. different hyperplane separating the data points.

ii.

g1 hyperplane

1st iteration



g2 hyperplane

1st iteration

g1 hyperplane

2nd iteration



g2 hyperplane

2nd iteration

In the 2nd iteration, only the circled data points are labelled.

c. The hyper plane won't change as the co-training is performed based on the farthest distance from the hyper-plane, the nearer unlabeled points won't change and the hyper plane would still remain the same as $S^3VM$.

4. Data sequence: S= ({A,B} {B,C} {D,E} {A,D} {B,E,F})

   With $n_g = 0$, $x_g = 2$, $m_s = 5$, $w_s = 1$,

   a. No, it is not a subsequence of S as it violates the min_gap constraint for the event {A}{B}{C}. If we choose {A}{B} to be in the same element in the data sequence, e1, this leads to a min gap of 0 and violates the constraint. Similarly , for {B}{C}, If we choose {B}{C} to be in the same element in the data sequence, e2, this leads to a min gap of 0 and violates the constraint.

   b. Yes, it is a subsequence of S as it satisfies all the time constraints.

   c. Yes, it is a subsequence of S as it satisfies all the time constraints.

   d. No, it is not a subsequence of S as it violates the max_gap constraint for the event {C,D,E} {A,F} in w as $x_g = 3$.

   e. No, it is not a subsequence of S as it violates the max_gap constraint for the event {A,B,C,D} {E,F} in w as $x_g = 4$.

5. From the given frequent 3-sequence,

   a. Candidate 4-sequences are:

      i.   <{1,2,3} {4}>
      ii.  <{1,2,3} {5}>
      iii. <{1,3} {4,5}>
      iv.  <{1,3} {4} {5}>
      v.   <{1} {2,4} {4}>
      vi.  <{2,4} {4,5}>
      vii. <{2,3} {4,5}>
      viii. <{2,3} {4} {5}>

   b. Candidate 4-sequences that are pruned are:

      i.   <{1,2,3} {4}> as <{1,2} {4}> not there.
      ii.  <{1,2,3} {5}> as <{1,2} {5}> not there.
      iii. <{1,3} {4,5}> as <{1,3} {5}> not there.
      iv.  <{1,3} {4} {5}> as <{1,3} {5}> not there.
      v.   <{1} {2,4} {4}> as <{1} {2} {4}> not there.
      vi.  <{2,4} {4,5}> as <{2,4} {5}> not there.
      vii. <{2,3} {4} {5}> as <{2} {4} {5}> not there

      Only <{2,3} {4,5}> is not pruned

   c. With the given time constraint, the candidates that are pruned are:

      i.   <{1,2,3} {4}> as {1,2} {4} contiguous subsequence is not frequent.
      ii.  <{1,2,3} {5}> as {1,2} {5} contiguous subsequence is not frequent.

iii.    <{1,3} {4,5}> as {1,3} {5} contiguous subsequence is not frequent.
iv.    <{1} {2,4} {4}> as <{1} {2} {4}> contiguous subsequence is not
       frequent.
v.     <{2,4} {4,5}> as <{2,4} {5}> contiguous subsequence is not
       frequent.
vi.    <{2,3} {4} {5}> as <{2} {4} {5}> contiguous subsequence is not
       frequent.

6.



U(with no expert) => (0.05*10000K - 600 * 0.95) => -70K

| Type | Expert Pass | Expert Fail |
|---|---|---|
| Company success | 0.75*0.05 | 0.25*0.05 |
| Company fail | 0.25*0.95 | 0.75*0.95 |

P(Expert Pass) => 0.75*0.05 + 0.25*0.95 => 0.275
P(Company success|expert pass) => 0.75*0.05/0.275 => 0.13636
P(Company fail | exp pass) => 0.25*0.95/0.275 => 0.8636
U(Buy Company given exp pass) => 0.13636*9800K-800K*0.8636 => $645.448K

P(Expert fail) => 0.25*0.05+0.75*0.95 => 0.725
P(Company success | exp fail) => 0.25*0.05/0.725 => 0.0172
P(Company fail | exp fail) => 0.75*0.95/0.725 => 0.9827
U(Buy Company given exp fail) => 0.0172*9800K - 0.9827*800K => -$617.6K

U(with expert) => 0.275*645.448K - 0.725*200K => 32.4982K

Conclusion: We would only buy the company by hiring the experts and taking their suggestion since the chance of getting profit is more but, if they suggest no then we won't buy. If we don't take suggestions then it will be a higher loss.