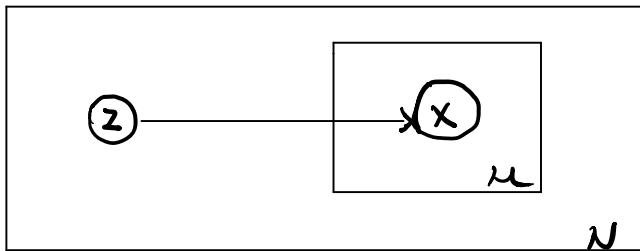


1. 1.0
(a)

- (b) learning parameters are $\theta_A, \theta_B, \theta_C$
- (c) latent variables are all $Z^{(i)}$ which indicate city assignment

1.1 (a) $\log p(\bar{X}^{(i)}, \bar{Z}^{(i)}; \bar{\theta}) = \log \frac{1}{3} \cdot (\theta_{Z^{(i)}})^{s_i} \cdot (1 - \theta_{Z^{(i)}})^{m - s_i}$
 $= \log \frac{1}{3} + s_i \cdot \log \theta_{Z^{(i)}} + (m - s_i) \log (1 - \theta_{Z^{(i)}})$

- Because for each city, we have $p = \frac{1}{3}$ to choose it and rain or not problem would be a Binomial($m, \theta_{Z^{(i)}}$) for each city.

- (b) Just multiply all of the trials together
- $\log P(\bar{X}, \bar{Z}; \bar{\theta}) = \log \prod_{i=1}^N p(\bar{X}^{(i)}, \bar{Z}^{(i)}; \bar{\theta})$
 $= \sum_{i=1}^N \log p(\bar{X}^{(i)}, \bar{Z}^{(i)}; \bar{\theta})$
 $= N \log \frac{1}{3} + s_i \cdot \log \theta_{Z^{(i)}} + (m - s_i) \log (1 - \theta_{Z^{(i)}})$

- (C) Denote the loglikelihood function as L

$$L = \log P(\lambda, Z; \bar{\theta})$$

\Rightarrow define $I_{Ai} = \begin{cases} 1 & \text{if } Z^{(i)} = A \\ 0 & \text{if } Z^{(i)} \neq A \end{cases} \quad \forall i$

$$I_{Bi} = \begin{cases} 1 & \text{if } Z^{(i)} = B \\ 0 & \text{if } Z^{(i)} \neq B \end{cases} \quad \forall i$$

$$I_{Ci} = \begin{cases} 1 & \text{if } Z^{(i)} = C \\ 0 & \text{if } Z^{(i)} \neq C \end{cases} \quad \forall i$$

$\therefore I_{Ai}, I_{Bi}, I_{Ci}$ are indicator functions for A, B, C

$$\cdot \frac{d}{d\theta_A} L = \sum_{i=1}^N [S_i \cdot (\theta_A)^{-1} - (\mu - S_i) \cdot (1 - \theta_A)^{-1}] \cdot I_{Ai} = 0$$

$$\sum_{i=1}^N [S_i \cdot (1 - \theta_A) - (\mu - S_i) \theta_A] \cdot I_{Ai} = 0$$

$$\sum_{i=1}^N [S_i - S_i \cdot \theta_A - \mu \theta_A + S_i \cdot \theta_A] \cdot I_{Ai} = 0$$

$$\hat{\theta}_A = \frac{\sum_{i=1}^N S_i \cdot I_{Ai}}{\sum_{i=1}^N \mu I_{Ai}}$$

\Rightarrow exactly same for B and C

$$\hat{\theta}_B = \frac{\sum_{i=1}^N S_i \cdot I_{Bi}}{\sum_{i=1}^N \mu I_{Bi}} ; \quad \hat{\theta}_C = \frac{\sum_{i=1}^N S_i \cdot I_{Ci}}{\sum_{i=1}^N \mu I_{Ci}}$$

- \cdot Intuitively, these 3 mle's $\hat{\theta}_A, \hat{\theta}_B, \hat{\theta}_C$ are just use the total number of sunny days of A, B, C in the sample data devived by the total number of observed sample total days of A, B, C to get a fraction. $\hat{\theta}_i$ is $\frac{\text{total sunny days of } i}{\text{total days of } i}$

$$\cdot \hat{\theta}_A = \frac{8+7}{10 \times 2} = 0.75 ; \quad \hat{\theta}_B = \frac{2+3}{2 \times 10} = 0.25$$

$$\hat{\theta}_C = \frac{11}{2 \times 10} = 0.55$$

$$\begin{aligned}
 1.2 \text{ (a)} \quad p(Z^{(i)} = k | \bar{x}^{(i)}, \bar{\theta}^{(i)}) &= \frac{P(Z^{(i)} = k, \bar{x}^{(i)}; \bar{\theta}^{(i)})}{P(\bar{x}^{(i)}; \bar{\theta}^{(i)})} \\
 &= \frac{P(\bar{x}^{(i)} | Z^{(i)} = k; \bar{\theta}^{(i)}) \cdot P(Z^{(i)} = k)}{\sum_{j=A}^C P(\bar{x}^{(i)} | Z^{(i)} = j; \bar{\theta}^{(i)}) \cdot P(Z^{(i)} = j)} \\
 \Rightarrow P(Z^{(i)} = k) &= P(Z^{(i)} = j) = \frac{1}{3} \quad \forall k, j \\
 \therefore &= \frac{P(\bar{x}^{(i)} | Z^{(i)} = k; \bar{\theta}^{(i)})}{\sum_{j=A}^C P(\bar{x}^{(i)} | Z^{(i)} = j; \bar{\theta}^{(i)})} \\
 \Rightarrow P(\bar{x}^{(i)} | Z^{(i)} = A; \bar{\theta}^{(i)}) &= \binom{n}{s_i} \cdot (\theta_A^{(i)})^{s_i} (1 - \theta_A^{(i)})^{n-s_i} \\
 \Rightarrow P(\bar{x}^{(i)} | Z^{(i)} = B; \bar{\theta}^{(i)}) &= \binom{n}{s_i} \cdot (\theta_B^{(i)})^{s_i} (1 - \theta_B^{(i)})^{n-s_i} \\
 \Rightarrow P(\bar{x}^{(i)} | Z^{(i)} = C; \bar{\theta}^{(i)}) &= \binom{n}{s_i} \cdot (\theta_C^{(i)})^{s_i} (1 - \theta_C^{(i)})^{n-s_i} \\
 \therefore p(Z^{(i)} = k | \bar{x}^{(i)}, \bar{\theta}) &= \frac{(\theta_K^{(i)})^{s_i} (1 - \theta_K^{(i)})^{n-s_i}}{(\theta_A^{(i)})^{s_i} (1 - \theta_A^{(i)})^{n-s_i} + (\theta_B^{(i)})^{s_i} (1 - \theta_B^{(i)})^{n-s_i} + (\theta_C^{(i)})^{s_i} (1 - \theta_C^{(i)})^{n-s_i}}
 \end{aligned}$$

(b) Write a python fn to calculate, see appendix.

$$\cdot \bar{\theta}^{(0)} = 0.7, 0.4, 0.6$$

$$① \quad \hat{Z}_1^{(1)}: \quad P(Z^{(1)} = A | \bar{x}; \bar{\theta}) = 0.6396$$

$$P(Z^{(1)} = B | \bar{x}; \bar{\theta}) = 0.0298 \quad \therefore \hat{Z}_1^{(1)} = A$$

$$P(Z^{(1)} = C | \bar{x}; \bar{\theta}) = 0.3393$$

$$Z_1^{(2)}: \quad P(Z^{(2)} = A | \bar{x}; \bar{\theta}) = 0.0378$$

$$P(Z^{(2)} = B | \bar{x}; \bar{\theta}) = 0.8068 \quad \therefore \hat{Z}_1^{(2)} = B$$

$$P(Z^{(2)} = C | \bar{x}; \bar{\theta}) = 0.1593$$

$$Z_1^{(3)}: \quad P(Z^{(3)} = A | \bar{x}; \bar{\theta}) = 0.2041$$

$$P(Z^{(3)} = B | \bar{x}; \bar{\theta}) = 0.397945 \quad \therefore B \& C \text{ tie} \Rightarrow B$$

$$P(Z^{(3)} = C | \bar{x}; \bar{\theta}) = 0.397945 \quad \therefore \hat{Z}_1^{(3)} = B$$

$$Z_1^{(4)}: \begin{aligned} P(Z=A|\bar{x}; \bar{\theta}) &= 0.0109 \\ P(Z=B|\bar{x}; \bar{\theta}) &= 0.9093 \end{aligned} \quad \therefore \hat{Z}_1^{(4)} = B$$

$$P(Z=C|\bar{x}; \bar{\theta}) = 0.0798$$

$$Z_1^{(5)}: \begin{aligned} P(Z=A|\bar{x}; \bar{\theta}) &= 0.3558 \\ P(Z=B|\bar{x}; \bar{\theta}) &= 0.1982 \end{aligned}$$

$$P(Z=C|\bar{x}; \bar{\theta}) = 0.4459$$

$$Z_1^{(6)}: \begin{aligned} P(Z=A|\bar{x}; \bar{\theta}) &= 0.5089 \\ P(Z=B|\bar{x}; \bar{\theta}) &= 0.081 \end{aligned}$$

$$P(Z=C|\bar{x}; \bar{\theta}) = 0.4101$$

$$\hat{\theta}_{A1} = \frac{8+7}{20} = 0.75; \quad \hat{\theta}_{B1} = \frac{3+5+2}{30} = \frac{1}{3}$$

$$\hat{\theta}_{C1} = \frac{6}{10} = 0.6$$

$$\therefore \bar{\theta} = [0.75, \frac{1}{3}, 0.6]$$

$$\bar{z}_1 = [A, B, B, B, C, A]$$

$$\textcircled{2} \quad Z_2^{(1)}: \begin{aligned} P(Z=A|\bar{x}; \bar{\theta}) &= 0.6942 \\ P(Z=B|\bar{x}; \bar{\theta}) &= 0.0075 \end{aligned}$$

$$\therefore \hat{Z}_2^{(1)} = A$$

$$P(Z=C|\bar{x}; \bar{\theta}) = 0.2981$$

$$Z_2^{(2)}: \begin{aligned} P(Z=A|\bar{x}; \bar{\theta}) &= 0.0101 \\ P(Z=B|\bar{x}; \bar{\theta}) &= 0.8509 \end{aligned}$$

$$\therefore \hat{Z}_2^{(2)} = B$$

$$P(Z=C|\bar{x}; \bar{\theta}) = 0.1389$$

$$Z_2^{(3)}: \begin{aligned} P(Z=A|\bar{x}; \bar{\theta}) &= 0.1476 \\ P(Z=B|\bar{x}; \bar{\theta}) &= 0.3452 \end{aligned}$$

$$\therefore \hat{Z}_2^{(3)} = C$$

$$P(Z=C|\bar{x}; \bar{\theta}) = 0.5072$$

$$\begin{aligned}
 Z_2^{(4)}: P(Z=A|\bar{x}; \bar{\theta}) &= 0.001 \\
 P(Z=B|\bar{x}; \bar{\theta}) &= 0.9466 \quad \therefore \hat{Z}_2^{(4)} = B \\
 P(Z=C|\bar{x}; \bar{\theta}) &= 0.0515 \\
 Z_2^{(5)}: P(Z=A|\bar{x}; \bar{\theta}) &= 0.3248 \\
 P(Z=B|\bar{x}; \bar{\theta}) &= 0.1254 \quad \therefore \hat{Z}_2^{(5)} = C \\
 P(Z=C|\bar{x}; \bar{\theta}) &= 0.5528 \\
 Z_2^{(6)}: P(Z=A|\bar{x}; \bar{\theta}) &= 0.5197 \\
 P(Z=B|\bar{x}; \bar{\theta}) &= 0.0337 \quad \therefore \hat{Z}_2^{(6)} = A \\
 P(Z=C|\bar{x}; \bar{\theta}) &= 0.4464 \\
 \therefore \hat{\theta}_{A2} &= \frac{8+7}{20} = 0.75 ; \quad \hat{\theta}_{B2} = \frac{3+2}{20} = 0.25 \\
 \hat{\theta}_{C2} &= \frac{11}{20} = 0.55 \\
 \therefore \bar{\theta}^{(2)} &= [0.75, 0.25, 0.55] \\
 \hat{Z}_2 &= [A \ B \ C \ B \ C \ A]
 \end{aligned}$$

$$\begin{aligned}
 Z_3^{(1)}: P(Z=A|\bar{x}; \bar{\theta}) &= 0.7859 \\
 P(Z=B|\bar{x}; \bar{\theta}) &= 0.001 \quad \therefore \hat{Z}_3^{(1)} = A \\
 P(Z=C|\bar{x}; \bar{\theta}) &= 0.2129 \\
 Z_3^{(2)}: P(Z=A|\bar{x}; \bar{\theta}) &= 0.009 \\
 P(Z=B|\bar{x}; \bar{\theta}) &= 0.76313 \quad \therefore \hat{Z}_3^{(2)} = B \\
 P(Z=C|\bar{x}; \bar{\theta}) &= 0.2274 \\
 Z_3^{(3)}: P(Z=A|\bar{x}; \bar{\theta}) &= 0.16645 \\
 P(Z=B|\bar{x}; \bar{\theta}) &= 0.16645 \quad \therefore \hat{Z}_3^{(3)} = C \\
 P(Z=C|\bar{x}; \bar{\theta}) &= 0.667
 \end{aligned}$$

$$\begin{aligned}
 Z_3^{(4)}: P(Z=1 | \bar{x}; \bar{\theta}) &= 0.00126 \\
 P(Z=2 | \bar{x}; \bar{\theta}) &= 0.923 \quad \therefore \hat{Z}_3^{(4)} = B \\
 P(Z=3 | \bar{x}; \bar{\theta}) &= 0.07509 \\
 Z_3^{(5)}: P(Z=1 | \bar{x}; \bar{\theta}) &= 0.36446 \\
 P(Z=2 | \bar{x}; \bar{\theta}) &= 0.0405 \quad \therefore \hat{Z}_3^{(5)} = C \\
 P(Z=3 | \bar{x}; \bar{\theta}) &= 0.595 \\
 Z_3^{(6)}: P(Z=1 | \bar{x}; \bar{\theta}) &= 0.5961 \\
 P(Z=2 | \bar{x}; \bar{\theta}) &= 0.007 \quad \therefore \hat{Z}_3^{(6)} = A \\
 P(Z=3 | \bar{x}; \bar{\theta}) &= 0.3965
 \end{aligned}$$

- we can find that Z assignment does not change
 $\Rightarrow \bar{\theta}$ remains unchanged
 $\therefore \bar{\theta}^{(3)} = [0.75, 0.25, 0.55]$
- $\hat{Z}_3 = [A \ B \ C \ B \ C \ A]$

1.3. (a) to find θ_K that maximize $\theta_K^{(t+1)}$

$$\begin{aligned}
 \therefore \frac{d}{d\theta_K} \theta_K^{(t+1)} &= \sum_{i=1}^N p(Z^{(i)}=k | \bar{x}^{(i)}, \bar{\theta}^{(t)}) \cdot \\
 &\quad (s_i \cdot \theta_K^{(t)} - (N-s_i)(1-\theta_K^{(t)})) = 0 \\
 \sum_{i=1}^N p(Z^{(i)}=k | \bar{x}^{(i)}, \bar{\theta}^{(t)}) \cdot (s_i - N\theta_K^{(t)}) &= 0
 \end{aligned}$$

$$\sum_{i=1}^N p(Z^{(i)}=k | \bar{x}^{(i)}, \bar{\theta}^{(t)}) \cdot s_i = \sum_{i=1}^N p(Z^{(i)}=k | \bar{x}^{(i)}, \bar{\theta}^{(t)}) \cdot N \cdot \hat{\theta}_K$$

$$\therefore \hat{\theta}_K^{(t+1)} = \hat{\theta}_K = \frac{\sum_{i=1}^N p(Z^{(i)}=k | \bar{x}^{(i)}, \bar{\theta}^{(t)}) \cdot s_i}{N \cdot \sum_{i=1}^N p(Z^{(i)}=k | \bar{x}^{(i)}, \bar{\theta}^{(t)})}$$

- (b) · The update is fraction which
=> numerator is sum of the number of sunny
time the probability of being this city.
we can understand this as an "advanced"
sum of sunny days altogether considering
all sample.
- => denominator is same type as numerator
except that we are now considering total
of days using weight of each city on
each trial.
- The equation is just updating by "expected"
ratio of sunny days on total days for
each city which make huge sense to me.
 - Yes. correspond to my intuition. reasonable !

$$\begin{aligned}
 2. (a) \cdot \text{posterior } P(Z=1 | w_1, w_2) &= \frac{P(Z=1, w_1, w_2)}{P(w_1, w_2)} \\
 \Rightarrow P(Z, w_1, w_2) &= P(Z, w_2 | w_1) P(w_1) \\
 &= P(w_2 | Z, w_1) \cdot P(Z | w_1) P(w_1) \\
 (w_2 \perp\!\!\!\perp w_1 | Z) &= P(w_2 | Z) \cdot P(Z | w_1) P(w_1) \\
 &= \theta_{w_2|Z} \phi_{Z|w_1} P(w_1) \\
 \Rightarrow P(w_1, w_2) &= \sum_{z=1}^k P(Z, w_1, w_2)
 \end{aligned}$$

$$\begin{aligned}
 \therefore P(Z=1 | w_1 = \text{"machine"}, w_2 = \text{"learning"}) &= \frac{\theta_{w_2=\text{"learning"} | Z=1} \phi_{Z=1 | w_1=\text{"machine"}} P(w_1)}{\sum_{z=1}^k \theta_{w_2=\text{"learning"} | Z} \phi_{Z | w_1=\text{"machine"}} - P(w_1)} \\
 &= \frac{\theta_{w_2=\text{"learning"} | Z=1} \phi_{Z=1 | w_1=\text{"machine"}}}{\sum_{z=1}^k \theta_{w_2=\text{"learning"} | Z} \phi_{Z | w_1=\text{"machine"}}}
 \end{aligned}$$

$$\begin{aligned}
 (b) P(w_2, w_3, \dots, w_T | w_1) &= \frac{P(w_1, w_2, \dots, w_T)}{P(w_1)} \\
 &= \frac{P(w_1) \cdot P(w_2 | w_1) \cdot P(w_3 | w_2) \cdots P(w_T | w_{T-1})}{P(w_1)} \\
 &= P(w_1 | w_1) P(w_2 | w_1) \cdots P(w_T | w_{T-1})
 \end{aligned}$$

$$\begin{aligned}
 \therefore \log P(w_2, \dots, w_T | w_1) &= \sum_{i=1}^{T-1} \log P(w_{i+1} | w_i) \\
 &= \sum_{i=1}^{T-1} \log \sum_{z=1}^k P(w_{i+1} | z) P(z | w_i) \\
 &= \sum_{i=1}^{T-1} \log \sum_{z=1}^k \theta_{w_{i+1}|z} \phi_{z|w_i}
 \end{aligned}$$

- (c) • The numerator is summing probability cluster z when a pair which the next word is w appears

- The denominator is summing probability of cluster z when pairs which the next word could be all words is in the dictionary appear , which we could understand as the prob. of cluster I given the words in dictionary. (All case)
- Intuitively, the update is doing

$$\theta_{wiz} = \frac{\text{prob}(z=z | \text{next word is } w)}{\text{prob}(I=z | \text{next word can be all words in dict})}$$
 which is very clear. the probability that the next word is w given cluster I .
- Yes θ correspond to my intuition

3. (a)
- we have weights of each Gaussian τ :
 $\Rightarrow K-1$ parameters
 - for mean of each Gaussian $\bar{\mu}^{(i)}$
 $\Rightarrow k \times d$ parameters
 - And 1 common variance shared σ^2
 $\Rightarrow 1$ parameters
 - ∴ total # = $K-1 + dk + 1 = cd + 1 - K$

(b) code only

(c) 1. Because we define our BIC as

$$BIC = m \ln N - 2 \ln \hat{L}$$

and our goal is to minimize BIC

\Rightarrow higher mle will result in lower BIC

\Rightarrow lower # of variable will result in lower BIC

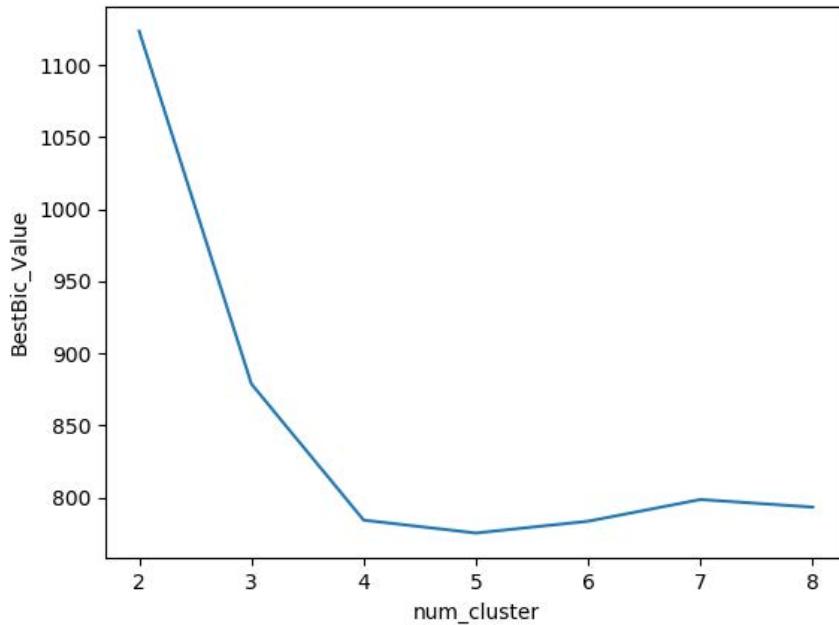
• So I'll pick a lower BIC value.

2. I'll choose the # of cluster according to

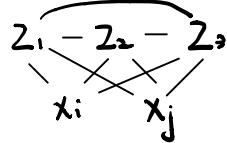
the "elbow" rule. I'll see the graph and find out the point that decreasing rate of mle becomes low. And that would be the point that our optimal number of cluster be.

- (d)
- The BIC falls on range around $700 \sim 1200$.
 - \Rightarrow The curve first going down dramatically
 - \Rightarrow Then the decreasing rate of the curve slows down
 - \Rightarrow after $k=5$, the curve rise.
 - Curve goes down at beginning because larger k induces a more complex model, result in higher log-likelihood \Rightarrow BIC goes down
 - Curve goes down because with $k \uparrow$, the increase rate of log-likelihood becomes lower than the penalty for having more clusters.
 - I would choose the point $k=5$.
Because our BIC continue decreasing greatly until we reach the point $k=5$, then increasing.
 \Rightarrow the increasing is due to the increase of mle cannot offset the penalty of the 1 more parameter, which tells us that we should stop adding more clusters.
 - Graph as below

Graph for 3.(d)



4. (a) • D-separation for check $x_i \perp\!\!\!\perp x_j \mid \{z_1, z_2\} \quad \forall i, j$



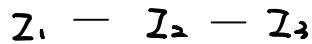
Two possible conditions to check

\Rightarrow there is obvious a path between variable we interest,
so go check next

\Rightarrow path $x_i \rightarrow z_3 \rightarrow x_j$ does not go through z_1 or z_2
 \therefore No !

$\therefore x_1, x_2, x_3$ are not independent of each other given Z_1, Z_2

• d-separation for check $Z_1 \perp\!\!\!\perp Z_3 \mid Z_2$



Two possible condition to check

$\Rightarrow \exists Z_1 \rightarrow Z_2 \rightarrow Z_3$. a path between Z_1 & Z_3

\Rightarrow ALL path between Z_1, Z_3 go through Z_2 !

\therefore Yes !

$\therefore Z_1 \perp\!\!\!\perp Z_3 \mid Z_2$

(b) $P(Z_1, Z_2, Z_3, x_1, x_2, x_3)$

$$= P(Z_1) \cdot P(Z_2 \mid Z_1) \cdot P(Z_3 \mid Z_2) \cdot P(x_1 \mid Z_1, Z_2, Z_3) \cdot P(x_2 \mid Z_1, Z_2, Z_3) \cdot P(x_3 \mid Z_1, Z_2, Z_3)$$

(c) $Z_1: 2-1 = 1 \quad ; \quad Z_2: 2 \times (2-1) = 2$

$Z_3: 2 \times (2-1) = 2$

$$x_1, x_2, x_3 : 2 \times 2 \times 2 \times (2-1) = 8$$

$$\therefore \text{Total \# of variables} = 1 + 2 + 2 + 8 \times 3 = 29$$

(d) I would use HML

- because as we can see in problem 3, even with the same log-likelihood, we will penalty a model more if it has relatively more parameters.
- HML has same log-likelihood as above model but less parameters.

\therefore HML

$$5. (a) P(S \cdot T \cdot R \cdot H) = P(S)P(T|S)P(R|S)P(H|T, R)$$

$$(b) P(R=T | H=T) = \frac{P(R=T, H=T)}{P(H=T)}$$

$$\begin{aligned} P(R, H) &= \sum_S \sum_T P(S, T, R, H) \\ &= \sum_S P(S) P(R|S) \sum_T P(T|S) P(H|T, R) \\ \therefore P(R=T, H=T) &= [0.8 \cdot 0.2 \cdot (0.6 \times 1 + 0.4 \times 0.5)] + \\ &\quad [0.2 \cdot 0.8 \cdot (0.4 \times 1 + 0.6 \times 0.5)] \\ &= 0.128 + 0.112 = 0.24 \end{aligned}$$

$$\begin{aligned} P(R=F, H=T) &= [0.8 \times 0.8 \times (0.6 \times 0.5 + 0.4 \times 0.1)] \\ &\quad + [0.2 \times 0.2 \times (0.4 \times 0.5 + 0.6 \times 0.1)] \\ &= 0.192 + 0.008 = 0.2 \end{aligned}$$

$$\therefore P(H=T) = 0.24 + 0.2 = 0.44$$

$$\therefore P(R=T | H=T) = \frac{0.24}{0.44} = 0.5454$$

- b. (a) • if we fix all $\{\bar{u}^{in}\}$, then our optimization problem becomes

$$\min_{v, \xi} I_{j=1}^m \frac{1}{2} \|\bar{v}^{in}\|^2 + C \sum_{j: Y_{ij} \neq 0} \xi_{ij}$$

$$\text{subject to } Y_{ij} (u^{in}, v^{in}) \geq 1 - \xi_{ij}$$

$$\xi_{ij} \geq 0 \quad \forall i, j \text{ where } Y_{ij} \neq 0$$

- This is a single SVM problem of classification w.r.t. \bar{v}^{in} and datapoint \bar{u}^{in} , label Y_{ij}
- Because for $\{\bar{v}^{in}\}$, there are totally $m \bar{v}^{in}$'s.
So our problem becomes m SVM problems on v

- (b) • Same as (a). our optimization problem becomes:

$$\min_{v, \xi} I_{i=1}^n \frac{1}{2} \|\bar{u}^{in}\|^2 + C \sum_{j: Y_{ij} \neq 0} \xi_{ij}$$

$$\text{subject to } Y_{ij} (u^{in}, v^{in}) \geq 1 - \xi_{ij}$$

$$\xi_{ij} \geq 0 \quad \forall i, j \text{ where } Y_{ij} \neq 0$$

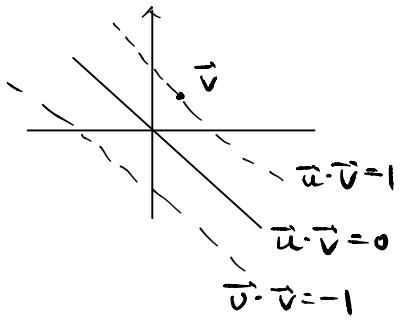
- This is a single SVM problem of classification w.r.t. \bar{u}^{in} and datapoint \bar{v}^{in} , label Y_{ij}
- Because for $\{\bar{u}^{in}\}$, there are totally $n \bar{u}^{in}$'s.
So our problem becomes n SVM problems on v

(c) • No.

- Because we only train on the observed datapoints, that is, train on all the (\vec{v}_j, y_j) that has $y_j \neq 0$. We can also see the constrain in (a)(b) SVM. we don't care the row that $y_j = 0$, no training on it so no effect.

(d)

- This would be a SVM problem on $\vec{u}^{(i)}$ with only one datapoint pair $\{\vec{v}^{(i)}, y_{ij}=1\}$
- Therefore our SVM becomes
$$\min \frac{1}{2} \|\vec{u}^{(i)}\|^2 \quad \text{s.t. } \vec{u}^{(i)} \cdot \vec{v}^{(i)} \geq 1$$
- Thus, $\vec{v}^{(i)}$ should be support vector!
$$\therefore \vec{u}^{(i)} \cdot \vec{v}^{(i)} = 1$$



- To maximize margin, we would take a decision boundary that perpendicular to \vec{v} , so \vec{u} should be parallel to \vec{v}

$$\Rightarrow \text{suppose } \vec{u} = a \cdot \vec{v}^{(i)}$$

$$\therefore a \cdot \|\vec{v}^{(i)}\|^2 = 1 \Rightarrow a = \frac{1}{\|\vec{v}^{(i)}\|^2}$$

$$\therefore \vec{u}^{(i)} = \frac{1}{\|\vec{v}^{(i)}\|^2} \cdot \vec{v}^{(i)}$$

(e) i. $\vec{v}^{(1)} = [-1.0]$; $\vec{v}^{(2)} = [0.1]$

and we only train on (1,1), (1,2) and (2,1)

- $\vec{u}^{(1)} \min \| \vec{u}^{(1)} \|^2 \text{ st. } (\vec{u}^{(1)} \cdot \vec{v}^{(i)}) \cdot Y_{i,1} \geq 1 \text{ for } i=1,2$

we have $\vec{v}^{(1)} = [-1.0]$ with $Y = -1$,

$\vec{v}^{(2)} = [0.1]$ with $Y = 1$

$$\Rightarrow \text{say } \vec{u}^{(1)} = [u_1, u_2]$$

$$\therefore (\vec{u}^{(1)} \cdot \vec{v}^{(1)}) \cdot Y_{1,1} = -u_1 \geq 1 \Rightarrow u_1 \leq 1$$

$$(\vec{u}^{(1)} \cdot \vec{v}^{(2)}) \cdot Y_{1,2} = u_2 \geq 1$$

$$\Rightarrow \text{to minimize } \| \vec{u}^{(1)} \|^2 = u_1^2 + u_2^2$$

because we have $u_1^2 \geq 1$, $u_2^2 \geq 1$

\therefore we have to take $u_1 = -1$, $u_2 = 1$

$$\Rightarrow \vec{u}^{(1)} = [-1, 1]$$

- $\vec{u}^{(2)} \min \| \vec{u}^{(2)} \|^2 \text{ st. } (\vec{u}^{(2)} \cdot \vec{v}^{(i)}) \cdot Y_{2,i} \geq 1$

we have $\vec{v}^{(2)} = [1.0]$. $Y = 1$

$$\Rightarrow \text{say } \vec{u}^{(2)} = [u_1, u_2]$$

$$\therefore (\vec{u}^{(2)} \cdot \vec{v}^{(2)}) \cdot 1 = u_1 \geq 1$$

$$\Rightarrow \text{to minimize } \| \vec{u}^{(2)} \|^2 = u_1^2 + u_2^2$$

because we have $u_1^2 \geq 1$, $u_2^2 \geq 0$

\therefore we have to take $u_1 = 1$, $u_2 = 0$

$$\therefore \vec{u}^{(2)} = [1.0]$$

$$\therefore U = \begin{bmatrix} -1 & 1 \\ 1 & 0 \end{bmatrix}$$

$$\therefore Y_{2,2} = \vec{u}^{(2)} \cdot \vec{v}^{(2)} = [1.0] \cdot [0.1] = 0$$

\therefore we cannot predict $Y_{2,2}$

$$\text{iii. } \bar{u}^{(1)} = [-1.1] ; \bar{u}^{(2)} = [1.0]$$

and we only train on (1,1), (1,2) and (2,1)

$$\cdot \bar{v}^{(1)} \min \| \bar{v}^{(1)} \|^2 \text{ st. } (\bar{u}^{(1)}, \bar{v}^{(1)}) \cdot Y_{i1} \geq 1 \text{ for } i=1,2$$

$$\text{we have } \bar{u}^{(1)} = [-1.1] \text{ w.l.o.g. } Y = -1$$

$$\bar{u}^{(2)} = [1.0] \text{ w.l.o.g. } Y = 1$$

$$\Rightarrow \text{say } \bar{v}^{(1)} = [v_1, v_2]$$

$$\therefore (\bar{u}^{(1)} \cdot \bar{v}^{(1)}) \cdot Y_{11} = -v_1 + v_2 \geq 1$$

$$(\bar{u}^{(1)} \cdot \bar{v}^{(1)}) \cdot Y_{12} = v_1 \geq 1$$

$$\Rightarrow \text{to minimize } \| \bar{v}^{(1)} \|^2 = v_1^2 + v_2^2$$

because we have

$$v_1^2 - 2v_1v_2 + v_2^2 \geq 1$$

$$v_1^2 + v_2^2 \geq 1 + 2v_1v_2$$

$$\therefore v_1 \geq 1 \Rightarrow v_2 = 0 \text{ for max}$$

$$\therefore v_1 = 1$$

$$\therefore \bar{v}^{(1)} = [1, 0]$$

$$\cdot \bar{u}^{(2)} \min \| \bar{v}^{(2)} \|^2 \text{ st. } (\bar{u}^{(2)}, \bar{v}^{(2)}) \cdot Y_{21} \geq 1$$

$$\text{we have } \bar{u}^{(2)} = [-1.1] . Y = 1$$

$$\Rightarrow \text{say } \bar{v}^{(2)} = [v_1, v_2]$$

$$\therefore (\bar{u}^{(2)} \cdot \bar{v}^{(2)}) \cdot 1 = -v_1 + v_2 \geq 1$$

$$\Rightarrow \text{to minimize } \| \bar{v}^{(2)} \|^2 = v_1^2 + v_2^2$$

$$\text{because we have } v_1^2 - 2v_1v_2 + v_2^2 \geq 1$$

$$v_1^2 + v_2^2 \geq 1 + 2v_1v_2$$

$$\geq 1 + 2v_1(1+v_1)$$

$$\geq 1 + 2v_1 + 2v_1^2$$

$$\geq 2(\frac{1}{2} + v_1)^2 + \frac{1}{2}$$

\therefore we have to have $v_1 = -\frac{1}{2} \Rightarrow v_2 \geq \frac{1}{2}$

\therefore to minimize $\|\bar{V}^{(2)}\|^2 \Rightarrow v_2 = \frac{1}{2}$

$$\therefore \bar{v}_2 = [1 - \frac{1}{2}, \frac{1}{2}]$$

$$\therefore Y_{22} = \bar{U}^{(2)} \cdot \bar{V}^{(2)} = [1, 0] \cdot [-\frac{1}{2}, \frac{1}{2}] = -\frac{1}{2}$$

$$\therefore \hat{Y}_{22} = -1 < 0 \text{ . we can predict it.}$$

APPENDIX

1.2 (b)

```
In [1]: import numpy as py

In [2]: def calc_z(p_a, p_b, p_c, sunny):
    x_a = ((p_a) ** sunny) * (1 - p_a) ** (10 - sunny)
    x_b = ((p_b) ** sunny) * (1 - p_b) ** (10 - sunny)
    x_c = ((p_c) ** sunny) * (1 - p_c) ** (10 - sunny)
    sum_abc = x_a + x_b + x_c
    return (x_a/sum_abc, x_b/sum_abc, x_c/sum_abc)

In [14]: calc_z(0.75,0.25,0.55,8)
Out[14]: (0.7859387940650213, 0.0010781053416529785, 0.21298310059332584)

In [15]: calc_z(0.75,0.25,0.55,3)
Out[15]: (0.009421150914142334, 0.7631132240455292, 0.2274656250403285)

In [16]: calc_z(0.75,0.25,0.55,5)
Out[16]: (0.1664595841967641, 0.1664595841967641, 0.6670808316064718)

In [17]: calc_z(0.75,0.25,0.55,2)
Out[17]: (0.0012670049499898517, 0.9236466085426018, 0.07508638650740826)

In [18]: calc_z(0.75,0.25,0.55,6)
Out[18]: (0.36446047786379454, 0.04049560865153273, 0.5950439134846727)

In [19]: calc_z(0.75,0.25,0.55,7)
Out[19]: (0.5961224206794469, 0.007359536057770949, 0.39651804326278217)
```

3 (b)

```
for iter in range(0,num_iter):
    """
        E-Step
        In the first step, we find the expected log-likelihood of the data which is equivalent to:
        finding cluster assignments for each point probabilistically
        In this section, you will calculate the values of zk(n,k) for all n and k according to current values
    """
    # TODO

    for i in range(0, N):
        for j in range(0, num_K):
            zk[i, j] = np.log(pk[j]) + calc_logpdf(trainX[i], mu[j], si2)
        zk[i] = zk[i] - logsumexp(zk[i])

    """
        M-step
        Compute the GMM parameters from the expressions which you have in your writeup
    """

    # Estimate new value of pk
    # TODO
    for j in range(0, num_K): pk[j] = np.exp(logsumexp(zk[:,j])) / N

    # Estimate new value for means
    # TODO
    for j in range(0, num_K): mu[j] = np.matmul([np.exp(zk[:,j])], trainX) / np.exp(logsumexp(zk[:,j]))

    # Estimate new value for sigma^2
    # TODO
    transform = np.zeros([N, num_K])
    for i in range(0, N):
        for j in range(0, num_K):
            transform[i, j] = (trainX[i] - mu[j]).dot(trainX[i] - mu[j])
    si2 = np.exp(logsumexp(zk, b = transform)) / (N * D)

    # Computing the expected likelihood of data for the optimal parameters computed
    # TODO
    for i in range(0, N):
        for j in range(0, num_K):
            zk[i, j] = np.log(pk[j]) + calc_logpdf(trainX[i], mu[j], si2)

    log_like = np.sum(logsumexp(zk, 1))

    # Compute the BIC for the current cluster
    # TODO
    BIC = (D + 1) * num_K * np.log(N) - 2 * log_like
```

3 (c) (d)

```
print("We'll try different numbers of clusters with GMM, using multiple runs for each to identify the 'best' results")
trainX = get_data()
num_K = range(2, 9) # List of cluster sizes
BIC_K = np.zeros(len(num_K))
means = {} # Dictionary mapping cluster size to corresponding matrix of means
cluster_proportions = {} # Dictionary mapping cluster size to corresponding mixture proportions vector
z_K = {}
sigma2 = {} # Dictionary mapping cluster size to the learned variance value
for idx in range(len(num_K)):
    # Running
    k = num_K[idx]
    print("%d clusters..." % k)
    bestBIC = float("inf")
    for i in range(1, 11):
        # TODO: Run gmm function 10 times and get the best
        # set of parameters for this particular value of k
        log_like = gmm(trainX, k)[4]
        if log_like < bestBIC: bestBIC = log_like
    BIC_K[idx] = bestBIC

    # TODO: Part d: Make a plot to show BIC as function of clusters K
    plt.plot(num_K, BIC_K)
    plt.xlabel("num_cluster")
    plt.ylabel("BestBic_Value")
    plt.show()
```