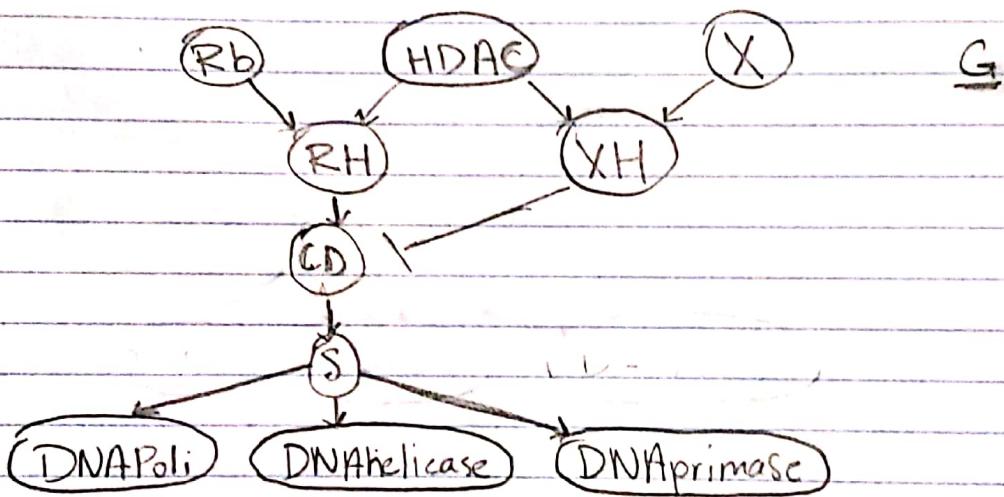


BME 230B - Homework 1

1. a) Nodes we have Rb, HDAC, RH, S, X, XH, Transcription - D-Pol, D-Heli, D-Prim, chromatin Domains CD



- b) RH and XH likely to be hidden b/c both act upon S and depend on HDAC so we could infer which is active based on Rb or X but we can't ever know for sure.

- c) The following independence assumptions are true
 $I(S \text{-phase}, Rb | CD = \text{closed})$,
 $I(HDAC, Rb)$, $I(HDAC, X)$

$$d) P(RH^+ | CD = \text{open}) = \frac{P(CD = \text{open} | RH^+) P(RH^+)}{P(CD = \text{open})}$$

neat page

$$P(A, B) = P(A)P(B)$$

so

both independent

$$\begin{aligned} P(RH+) &= P(RH+ | R+, HDAC+) P(R+, HDAC+) + \\ &\quad P(RH+ | R-, HDAC+) P(R-, HDAC+) + \\ &\quad P(RH+ | R+, HDAC-) P(R+, HDAC-) + \\ &\quad P(RH+ | R-, HDAC-) P(R-, HDAC-) \end{aligned}$$

$$\therefore P(RH-) = 1 - P(RH+) \quad P(A, B) = P(A)P(B)$$

$$\begin{aligned} P(CD=\text{open}) &= P(CD=\text{open} | RH+, XH+) P(RH+) P(XH+) + \\ &\quad P(CD=\text{open} | RH+, XH-) P(RH+) P(XH-) + \\ &\quad P(CD=\text{open} | RH-, XH+) P(RH-) P(XH+) + \\ &\quad P(CD=\text{open} | RH-, XH-) P(RH-) P(XH-) \end{aligned}$$

(needed for $P(CD=\text{open})$)

$$\begin{aligned} P(XH+) &= P(XH+ | HDAC+, X+) P(X+) P(HDAC+) + \\ &\quad P(XH+ | HDAC-, X+) P(X+) P(HDAC-) + \\ &\quad P(XH+ | HDAC+, X-) P(X-) P(HDAC+) + \\ &\quad P(XH+ | HDAC-, X-) P(X-) P(HDAC-) \end{aligned}$$

$$\therefore P(XH-) = 1 - P(XH+)$$

$$\begin{aligned} P(CD=\text{open} | RH+) &= P(CD=\text{open} | RH+, XH+) P(XH+) P(RH+) + \\ &\quad P(CD=\text{open} | RH+, XH-) P(XH-) P(RH+) \end{aligned}$$

Now all components can be plugged in to solve
for bayes rule for $P(RH+ | CD=\text{open})$.

2.a) BN of position independent model (a zero markov chain)

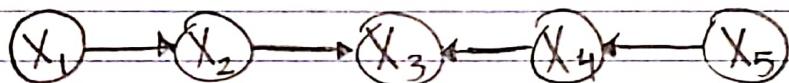


b/c no dependence from one another -

$$\Theta_{(x_i | p_a(i))} = \Theta_{x_i} \text{ for all } i \text{ in vector } \Theta$$

in the PSSM. Θ_i need 4 (for i in 1 to 5)

b) BN of position dependent model



Parameters needed:

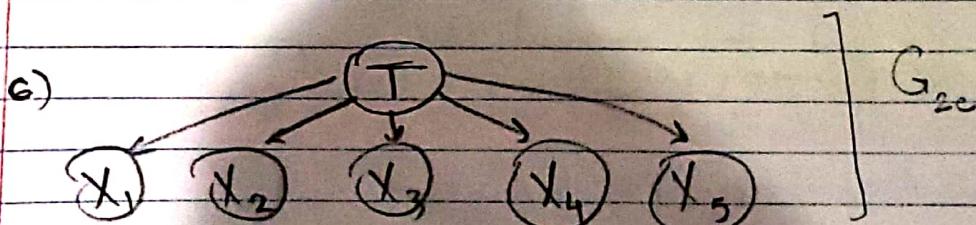
Θ_1 need 4 , 1 for each nuc in "ATGC"

Θ_5 needs 4 , 1 for each nuc in "ATGC"

Θ_2 needs 16 $\rightarrow \sum_{i \in ATGC} \sum_{j \in ATGC} P(\Theta_{2i}, \Theta_{1j})$

Θ_4 needs 16 $\rightarrow \sum_{i \in ATGC} \sum_{j \in ATGC} P(\Theta_{4i}, \Theta_{5j})$

Θ_3 needs 64 $\sum_{i \in ATGC} P(\Theta_{3i} | \Theta_{4j}, \Theta_{2k}) P(\Theta_{4j}) P(\Theta_{2k})$



each Θ_i in Θ needs 8 parameters

A, T, G, or C given TF1 and A, T, G, C given TF2

d) $P(X_1, X_2, \dots, X_5 | \Theta, G_{2c})$ find

$$\Theta = [\Theta_1, \Theta_2, \dots, \Theta_5]$$

G_{2c} = graph made in 2c

$X = X_1, X_2, \dots, X_5$ (the sequence)

$$P(X | \Theta, G_{2c}) = \prod_{i=1}^{n=5} P(X_i | \Theta_i)$$

where

$$P(X_i | \Theta_i) = P(X_i | \Theta_i, T = \text{TF1}) P(T = \text{TF1}) + P(X_i | \Theta_i, T = \text{TF2}) P(T = \text{TF2})$$

e) I would get all the sequences
and find

$$P(T = \text{TF1} | X_i)$$

where X_i is the binding site sequence
and X is the set of sequences.

for each we would find

$P(X_i | T = \text{TF1})$ the likelihood

$P(X_i)$ the normalizer

$P(T = \text{TF1})$ which is the prior

and use naive bayes.

from there I would find

$$P(T = \text{TF1} | X) = \frac{1}{n} \sum_{i=1}^n P(T = \text{TF1} | X_i)$$

where n is number of sequences in X .
and since T is binary so

$$P(T = \text{TF2} | X) = 1 - P(T = \text{TF1} | X)$$

which ever has higher probability is most likely the transcription factor used.

4 a) $\log \left[\frac{P_2(Y | \text{Model}-X)}{P_2(Y | \text{Model}-\text{Null})} \right] =$

$$\log [P_2(Y | \text{Model}-X)] - \log [P_2(Y | \text{Model}-\text{Null})]$$

log-likelihood for Model A = 29.90

$$-34.5305 - (-67.4307) = 29.9002$$

log-likelihood for Model B = 10.86

$$-56.5723 - (-67.4307) = 10.8584$$

b) on PDF

c) on PDF

d) Model A better explains the data
it has a better distribution of
data with a shape like a
gaussian distribution.

$$3 \text{ a) } P(\text{hap}^+) = 0.5 \quad P(X^+) = 0.5$$

$$P(\text{vpsT}^+) = P(\text{vpsT}^+ | \text{hap}^+) (\text{hap}^+) + P(\text{vpsT}^+ | \text{hap}^-)$$

$$P(\text{vpsT}^+) = (0.1)(0.5) + (0.9)(0.5)$$

$$P(\text{vpsT}^+) = 0.05 + 0.45 = 0.5$$

$$P(\text{hapR}^+) = 0.5 \quad P(\text{hapR}^-) = 0.5$$

$$P(\text{vpsT}^-) = 0.5 \quad P(\text{vpsT}^-) = 0.5$$

$$\begin{aligned} P(\text{vpsR}^+) &= P(\text{vpsR}^+ | \text{vpsT}^+, X^+) P(\text{vpsT}^+, X^+) + \\ &\quad P(\text{vpsR}^+ | \text{vpsT}^-, X^+) P(\text{vpsT}^-, X^+) + \\ &\quad P(\text{vpsR}^+ | \text{vpsT}^+, X^-) P(\text{vpsT}^+, X^-) + \\ &\quad P(\text{vpsR}^+ | \text{vpsT}^-, X^-) P(\text{vpsT}^-, X^-) \end{aligned}$$

$$\begin{aligned} P(\text{vpsR}^+) &= (0.9)(0.5 \cdot 0.5) + (0.9)(0.5 \cdot 0.5) + \\ &\quad (0.9)(0.5 \cdot 0.5) + (0.1)(0.5 \cdot 0.5) \end{aligned}$$

$$P(\text{vpsR}^+) = 0.7$$

$$P(\text{vpsR}^-) = 0.3$$

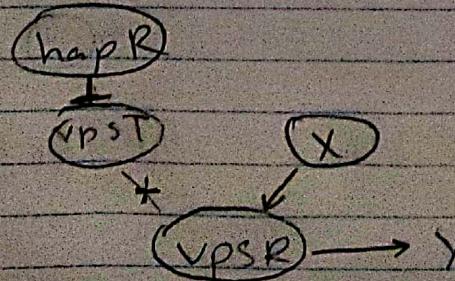
$$\begin{aligned} P(Y^+) &= P(Y^+ | \text{vpsR}^+) P(\text{vpsR}^+) + \\ &\quad P(Y^+ | \text{vpsR}^-) P(\text{vpsR}^-) \end{aligned}$$

$$P(Y^+) = (0.9)(0.7) + (0.1)(0.3)$$

$$P(Y^+) = 0.63 + 0.03 = 0.66$$

$$P(Y^+) = 0.66$$

b)



\ominus includes

- vpsT- (knocked out)
- $P(\text{vpsR}^+)$, still depends on if X^+ or not.
- Y^+ is still wild type prob.

$$c) P(vps^+, X^+ | Y^+, G_n^*, \Theta_n^*) =$$

$$\frac{P(Y^+, G_n^*, \Theta_n^*)}{P(Y^+, G_n^*, \Theta_n^*)} vps^+, X^+ P(vps^+, X^+) \\ \hookrightarrow P(Y^+)$$

$$P(vpsR^+, X^+) = P(vpsR^+ | X^+) P(X^+) \\ \text{b/c } vpsT \text{ is knocked out}$$

$$P(vpsR^+ | X^+) = P(vpsR^+ | X^+, vpsT^-)$$

$$P(vpsR^+ | X^+) = 0.9$$

$$P(X^+) = 0.5$$

$$P(vpsR^+, X^+) = 0.45$$

given G_n^*, Θ_n^*

$$P(Y^+) = P(Y^+ | vpsR^+) P(vpsR^+) + P(Y^+ | vpsR^-) P(vpsR^-)$$

$$P(vpsR^+) = P(vpsR^+ | X^+) P(X^+) + P(vpsR^+ | X^-) P(X^-)$$

$$P(vpsR^+) = (0.9)(0.5) + (0.1)(0.5)$$

$$P(vpsR^+) = .45 + .05 = 0.5$$

$$\therefore P(vpsR^-) = 0.5$$

$$P(Y^+ | vpsR^+) = 0.9$$

$$P(Y^+ | vpsR^-) = 0.1$$

$$P(Y^+) = (0.9)(0.5) + (0.1)(0.5)$$

$$\underline{P(Y^+) = 0.5}$$

$$P(Y^+ | vpsR^+, X^+) = 0.9 \times 0.9 \times 0.5 = 0.405$$

Hence

$$P(vpsR+, Y_1+ | G_A, \Theta_A) = \frac{405}{\cdot 5} \cdot .45 \\ = \boxed{1.3645}$$

3d) Let D be data

find $P(D | G_A, \Theta_A)$

$$P(D | G_A, \Theta_A) = \left(\prod_{i=1}^2 P(Y_i+ | G_A, \Theta_A) \right) \cdot P_{vpsT-}(Y_1- | G_A, \Theta_A) \\ \cdot P_{vpsT-}(Y_2+ | G_A, \Theta_A)$$

$$P(Y_1+ | G_A, \Theta_A) = P(Y_1+ | vspR+) P(vspR+) + \\ P(Y_1+ | vspR-) P(vspR-)$$

$$P(vspR+) = 0.7 \text{ and } P(vspR-) = 0.3 \\ (\text{both calculated in 3A})$$

$$P(Y_1+ | vspR+) = 0.9 \quad P(Y_1+ | vspR-) = 0.1 \\ (\text{from Model A})$$

$$P(Y_1+ | G_A, \Theta_A) = (0.9)(0.7) + (0.1) \cdot (0.3) \\ P(Y_1+ | G_A, \Theta_A) = 0.66$$

$$P_{vpsT-}(Y_1+ | G_A, \Theta_A) = P_{vpsT-}(Y_1+ | vpsR+) P_{vpsT-}(vpsR+) + \\ P_{vpsT-}(Y_1+ | vpsR-) P_{vpsT-}(vpsR-) \\ \downarrow$$

$$\text{b/c } P_{vpsT-}(Y_1- | G_A, \Theta_A) = 1 - P_{vpsT-}(Y_1+ | G_A, \Theta_A)$$

$$P_{vpsT-}(vpsR+) = 0.50 \quad P_{vpsT-}(vpsR-) = 0.5 \\ (\text{from 3c})$$

$$P_{\text{vspt-}}(Y_+ | \text{vspR+}) = 0.9 \quad \text{from model A}$$

$$P_{\text{vspt-}}(Y_+ | \text{vspR-}) = 0.1$$

$$P_{\text{vspt-}}(Y_+ | G_A, \Theta_A) = (0.9 \cdot 0.5) + (0.1 \cdot 0.5)$$

$$P_{\text{vspt-}}(Y_+ | G_A, \Theta_A) = 0.5$$

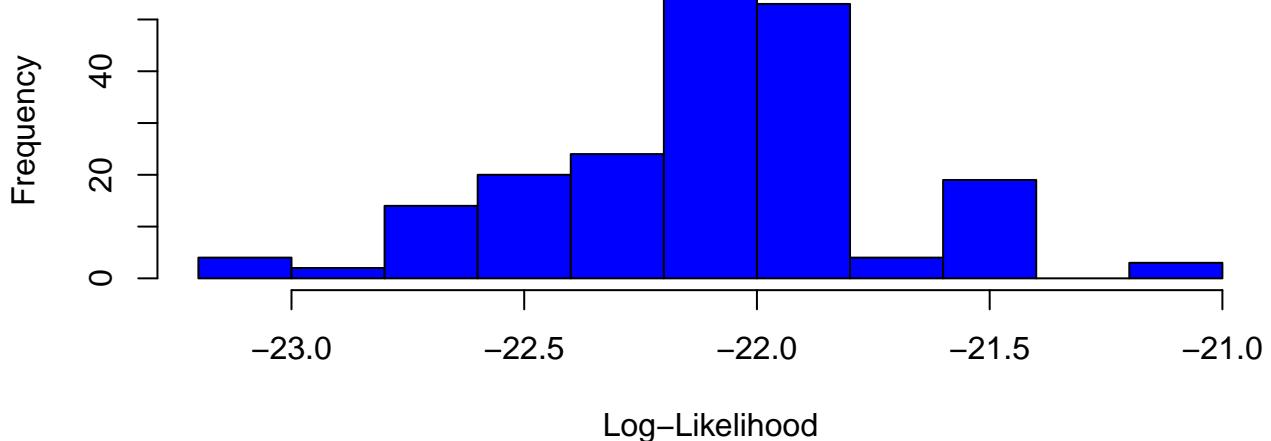
$$P_{\text{vspt-}}(Y_- | G_A, \Theta_A) = 0.5$$

$$\text{b/c same model } P_{\text{vspt-}}(Y_+ | G_A, \Theta_A) = 0.5$$

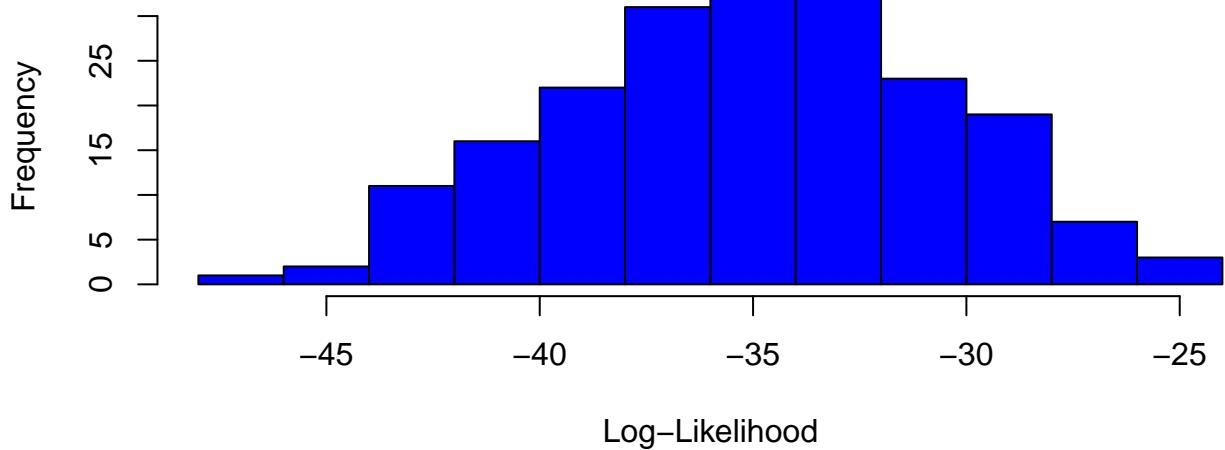
$$P(D | G_A, \Theta_A) = \left[\prod_{i=1}^2 (0.66) \right] \cdot 0.5 \cdot 0.5$$

$$P(D | G_A, \Theta_A) = 0.1089$$

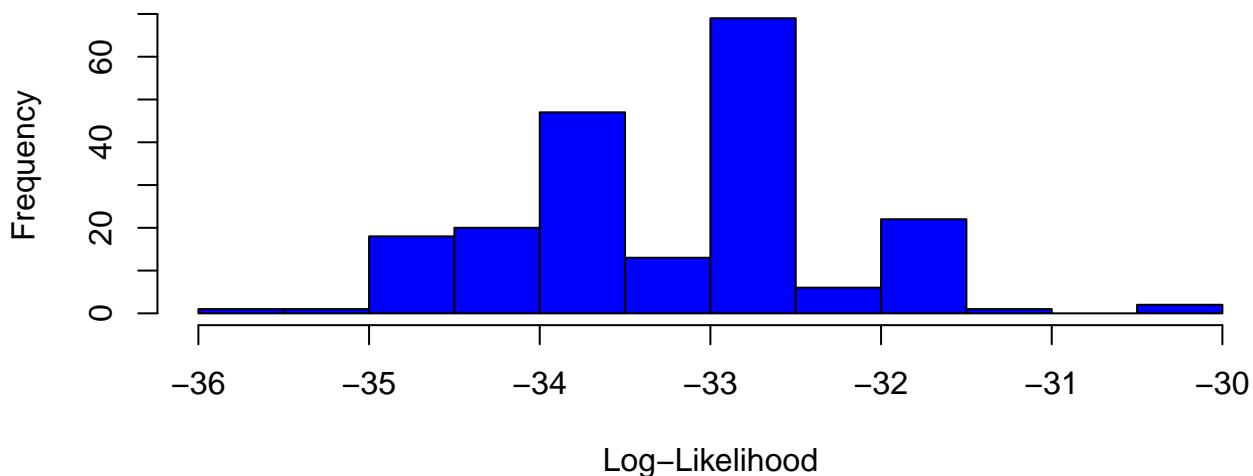
Log-Likelihood from 10 Genes in y.txt – Model A



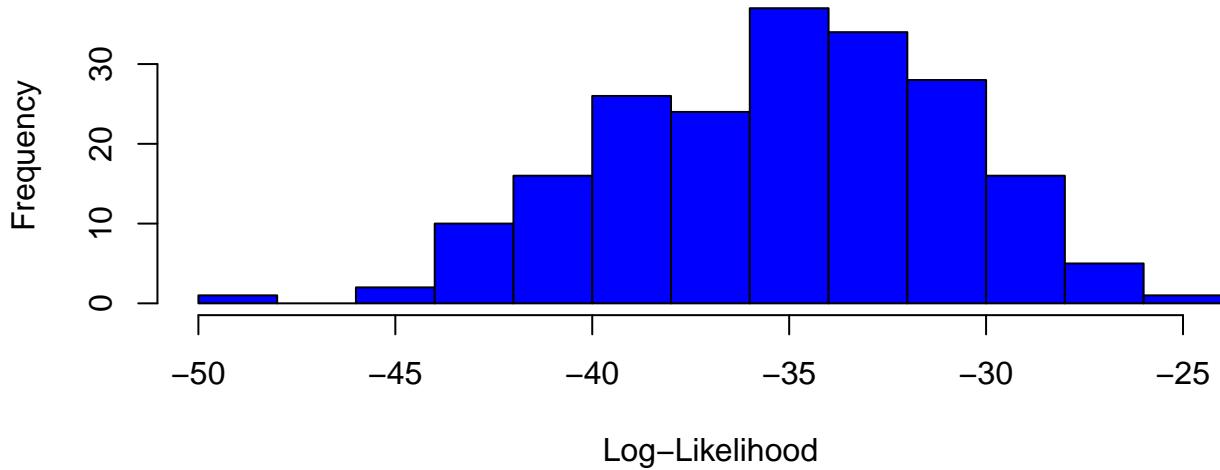
Log-Likelihood from 10 Random Genes – Model A



Log-Likelihood from 10 Genes in y.txt – Model B



Log-Likelihood from 10 Random Genes – Model B



```

# -----
# HW1 Problem 4b
#
# Script by: Mohammad Abdulqader
#
# this shell script is used to solve problem 4b in
# the HW1 assignment for BME 230B
# -----


let a=0

# create these files
echo log_likelihood > experiment_1.txt
echo log_likelihood > experiment_2.txt
tail -3850 data.txt | awk -F '\t' '{print $1}' > y2.txt

while [ $a -lt 200 ]
do

    # shuf to get 10 random Y values.
    shuf -n 10 y.txt > sample_Y.txt
    shuf -n 10 y2.txt > sample_Y2.txt

    # now run pathy.jar for first experiment
    java -jar pathy.jar sample_Y.txt data.txt modela.bif >> experiment_1.txt

    # now run pathy.jar for second experiment
    java -jar pathy.jar sample_Y2.txt data.txt modela.bif >> experiment_2.txt

    # increment a (number of iterations)
    let a=a+1

done

# remove excess files
rm sample_Y.txt ReadMe.txt YgenesIndScore.txt y2.txt sample_Y2.txt

# done

```

```

# -----
# HW1 Problem 4c
#
# Script by: Mohammad Abdulqader
#
# this shell script is used to solve problem 4c in
# the HW1 assignment for BME 230B
# -----


let a=0

# create these files
echo log_likelihood > experiment_3.txt
echo log_likelihood > experiment_4.txt
tail -3850 data.txt | awk -F '\t' '{print $1}' > y2.txt

while [ $a -lt 200 ]
do

    # shuf to get 10 random Y values.
    shuf -n 10 y.txt > sample_Y.txt
    shuf -n 10 y2.txt > sample_Y2.txt

    # now run pathy.jar for first experiment
    java -jar pathy.jar sample_Y.txt data.txt modelb.bif >> experiment_3.txt

    # now run pathy.jar for second experiment
    java -jar pathy.jar sample_Y2.txt data.txt modelb.bif >> experiment_4.txt

    # increment a (number of iterations)
    let a=a+1

done

# remove excess files
rm sample_Y.txt ReadMe.txt YgenesIndScore.txt y2.txt sample_Y2.txt

# done

```

```

# -----
# HW1 Problem 4
#
# Script by: Mohammad Abdulqader
#
# this R script is used to solve problem 4b and 4c in
# the HW1 assignment for BME 230B
# -----


# read data
data1 <-read.csv('experiment_1.txt')
data2 <-read.csv('experiment_2.txt')
data3 <-read.csv('experiment_3.txt')
data4 <-read.csv('experiment_4.txt')


# Model A
pdf('Problem_4b.pdf')
par(mfrow = c(2, 1))

hist(data1$log_likelihood, main='Log-Likelihood from 10 Genes in y.txt - Model A',
      col = 'blue', xlab = 'Log-Likelihood')
hist(data2$log_likelihood, main='Log-Likelihood from 10 Random Genes - Model A',
      col = 'blue', xlab = 'Log-Likelihood')

dev.off()

# Model B
pdf('Problem_4c.pdf')
par(mfrow = c(2, 1))

hist(data3$log_likelihood, main='Log-Likelihood from 10 Genes in y.txt - Model B',
      col = 'blue', xlab = 'Log-Likelihood')
hist(data4$log_likelihood, main='Log-Likelihood from 10 Random Genes - Model B',
      col = 'blue', xlab = 'Log-Likelihood')

dev.off()

par(mfrow = c(1, 1))

```