

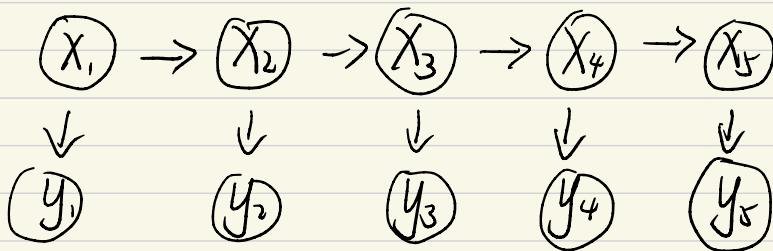
## Problem 1 (Hidden Markov Model)

assum1  $P(X_t | X_{t-1}, \dots, X_1) = P(X_t | X_{t-1}) \quad t=2, \dots, T$

assum2  $P(Y_t | X_1, \dots, X_t, \text{other } y's) = P(Y_t | X_t) \quad t=1, \dots, T$

homogenous: when  $P(X_{t+1}=j | X_t=i)$  and  $P(Y_t | X_t=k)$   
do not change as a function of  $t$

①



$$P(X_1, \dots, X_T, Y_1, \dots, Y_T) = \prod_{t=1}^T P(Y_t | \text{pac}(Y_t)) P(X_t | \text{pac}(X_t))$$

$$\stackrel{\text{by assum1}}{=} \left( \prod_{t=1}^T P(Y_t | X_t) \right) \cdot \left( \prod_{t=2}^T P(X_t | X_{t-1}) \right) \cdot P(X_1)$$

②

let  $P(Y_t | X_t=k) \sim N(\mu_k, \Sigma_k)$  where  $\Sigma_k$  is  $d$  by  $d$   
 $\mu_k$  is  $1$  by  $d$

(a) if cov matrix  $\Sigma_k$  is diagonal

(\ ) needs  $d$  parameters for each  $k$

$\mu_k$  needs  $d$  parameter for each  $k$

since each  $X_{t+1}$  has  $K$  value and  $X_t$  has  $K$  values and  
thus  $P(X_t | X_{t-1})$  requires  $\underbrace{\dots}_{t \in \{2, \dots, T\}}$  each row sum to 1  
 $K \cdot$  degree of freedom  $K-1$

$P(X_t)$  has  $K-1$  degree of freedom, thus  $K-1$  parameters

thus needs

$$Kd + kd + K(K-1) + K-1$$

$$= 2Kd + K^2 - K + K - 1$$

$$= 2Kd + K^2 - 1$$

(b) if  $\Sigma_k$  is full, since it's symmetric  
 $\binom{1+1+d}{2}$  has  $1+2+\dots+d = \frac{(1+d)d}{2}$  degree of freedom

thus total needs

$$K \left( \frac{1+d)d}{2} + kd + K(K-1) + k - 1 \right)$$

$$= kd \left( \frac{3+d}{2} \right) + K^2 - 1$$

③

$$P(X_5=k | y_1^*, \dots, y_5^*) = \frac{P(X_5=k, y_1^*, \dots, y_5^*)}{P(y_1^*, \dots, y_5^*)}$$

$$\text{joint prob} = P(X_5=k, y_1^*, \dots, y_5^*)$$

$$\begin{aligned} &= P(y_5^* | X_5=k, y_1^*, \dots, y_4^*) \cdot P(X_5=k, y_1^*, \dots, y_4^*) \\ &= P(y_5^* | X_5=k) P(X_5=k, y_1^*, \dots, y_4^*) \\ &= P(y_5^* | X_5=k) \sum_{j=1}^k P(X_5=k, X_4=j, y_1^*, \dots, y_4^*) \\ &= P(y_5^* | X_5=k) \sum_{j=1}^k P(X_5=k | X_4=j) P(X_4=j, y_1^*, \dots, y_4^*) \\ &= P(y_5^* | X_5=k) \sum_{j=1}^k P(X_5=k | X_4=j) P(Y_4^* | X_4=j) P(X_4=j, y_1^*, y_2^*, y_3^*) \end{aligned}$$

from hint 2

$$\begin{aligned} P(X_3=i, y_1^*, y_2^*, y_3^*) &= P(y_3^* | X_3=i) P(X_3=i, y_1^*, y_2^*) \\ &= P(y_3^* | X_3=i) \sum_{j=1}^k P(X_3=i, X_2=j, y_1^*, y_2^*) \\ &= P(y_3^* | X_3=i) \sum_{j=1}^k P(X_3=i | X_2=j) P(X_2=j, y_1^*, y_2^*) \end{aligned}$$

In general

$$P(X_t=i, y_1^*, \dots, y_t^*) = P(y_t^* | X_t=i) \sum_{j=1}^k P(X_t=i | X_{t-1}=j) \cdot P(X_{t-1}=j, y_1^*, y_2^*)$$

$$\text{let } g_{t,i} = P(X_t=i, y_1^*, \dots, y_t^*)$$

then  
systematically

$$\left\{ \begin{array}{l} g_{t,i} = P(Y_t^* | X_t=i) \sum_{j=1}^k P(X_t=j | X_{t-1}=j) g_{t-1,j} \\ g_{1,i} = P(X_1=i, Y_1^*) \end{array} \right.$$

by this recursion start from  $g_{1,i}$  to  $g_{t,i}$  we can get

$$P(X_5=k, Y_1^*, \dots, Y_5^*) = g_{t,k}$$

for magmal  $P(Y_1^*, \dots, Y_5^*) = \sum_{k=1}^K P(X_5=k, Y_1^*, \dots, Y_5^*) = \sum_{k=1}^K g_{t,k}$

thus  $P(X_5=k | Y_1^*, \dots, Y_5^*) = \frac{g_{t,k}}{\sum_{k=1}^K g_{t,k}}$  by bayes formula

### Problem 3

(Maximum likelihood for the Multinomial Model)

1.  $L(\theta) = P(D|\theta) = \prod_{i=1}^N P(X_i|\theta)$   
 $= \prod_{i=1}^N \theta_{X_i}^{r_i}$  where  $N = \sum_{i=1}^V r_i$ ,  $X_i$  is the  $i$ th trial  
 $= \prod_{k=1}^V \theta_k^{r_k}$

2.  $\ell(\theta) = \log L(\theta) = \log \prod_{k=1}^V \theta_k^{r_k}$   
 $= \sum_{k=1}^V r_k \log \theta_k$

adding lagrange operator

$$\ell(\theta) = \sum_{k=1}^V r_k \log \theta_k + \lambda (1 - \sum_{k=1}^V \theta_k)$$

$$\frac{d\ell(\theta)}{d\theta_k} = \frac{r_k}{\theta_k} - \lambda \stackrel{\text{want}}{=} 0 \quad \text{for } k \in \{1, \dots, V\}$$

$$\theta_k = \frac{r_k}{\lambda}$$

consider  $\sum_{k=1}^V \theta_k = \sum_{k=1}^V \frac{r_k}{\lambda}$

then  $\lambda = \sum_{k=1}^V r_k$

$$\text{thus } \hat{\theta}_k^{\text{ml}} = \frac{r_k}{n} = \frac{r_k}{\sum_{i=1}^n r_i}$$

### Problem 4

(Maximum likelihood for Geometric Model)

$$P(X=k) = (1-\theta)^{k-1} \theta, \quad k=1, 2, \dots \quad 0 < \theta < 1$$

①

$$\begin{aligned} L(\theta) &= P(D|\theta) = \prod_{i=1}^n P(X_i=k_i|\theta) \\ &= \prod_{i=1}^n (1-\theta)^{k_i-1} \theta \\ &= \theta^n (1-\theta)^{\sum_{i=1}^n k_i - n} \end{aligned}$$

②

$$\begin{aligned} \frac{d \ell(\theta)}{d \theta} &= \frac{d \log(L(\theta))}{d \theta} = (n \theta^{n-1} (1-\theta)^{\sum_{i=1}^n k_i - n}) - (\sum_{i=1}^n k_i - n) (1-\theta)^{\sum_{i=1}^n k_i - n-1} \theta^n \\ &\stackrel{\text{want}}{=} 0 \end{aligned}$$

$$\text{then } n(1-\theta) = (\sum_{i=1}^n k_i - n) \theta$$

$$\sum_{i=1}^n k_i \theta = n$$

$$\hat{\theta}^{\text{ml}} = \frac{n}{\sum_{i=1}^n k_i}$$

### Problem 5

(Maximum likelihood for Poisson Model)

$$D = \{X_1, \dots, X_n\}, \quad X_i \in \{0, 1, 2, \dots\}$$

$$P(X=k) = \frac{e^{-\lambda} \lambda^k}{k!} \quad \lambda > 0 \quad k \in \{0, 1, 2, 3, \dots\}$$

①

$$L(\theta) = P(D|\theta) = \prod_{i=1}^n P(X_i=k_i|\theta)$$

$$\begin{aligned} &= \prod_{i=1}^n \frac{\theta^{k_i} \lambda^{k_i}}{k_i!} \\ &= e^{-n\lambda} \lambda^{\sum_{i=1}^n k_i} \prod_{i=1}^n \frac{1}{k_i!} \end{aligned}$$

②

$$\ell(\theta) = \log L(\theta) = -n\lambda + \sum_{i=1}^n k_i \log(\lambda) - \sum_{i=1}^n \log(k_i!)$$

consider  $\frac{\partial \ell(\theta)}{\partial \lambda} = -n + \frac{\sum_{i=1}^n k_i}{\lambda} \stackrel{\text{want}}{=} 0$

thus  $\hat{\lambda}^{ml} = \frac{\sum_{i=1}^n k_i}{n}$

**Problem 6**

(Maximum likelihood with Measurement Variance per point)

D consists of N scalar measurements  $x_i$ ,  $1 \leq i \leq N$   
 same mean  $\mu$ , different variance  $\sigma_i^2$   $1 \leq i \leq N$  know  $\sigma_i^2$   
 brightness  $\mu$  of a star, measurement  $x_i$  taken at  
 different locations  $i$  on planet with noise  $\sigma_i^2$

$$\textcircled{1} \quad P(x_i | \mu, \sigma_i^2) = \frac{1}{\sqrt{2\pi\sigma_i^2}} e^{-\frac{1}{2\sigma_i^2}(x_i - \mu)^2}$$

$$\begin{aligned} L(\theta) &= P(D|\theta) = \prod_{i=1}^n P(x_i = k_i | \theta) \\ &= \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma_i^2}} e^{-\frac{1}{2\sigma_i^2}(x_i - \mu)^2} \end{aligned}$$

$$l(\theta) = \log L(\theta) = \sum_{i=1}^n \left( -\frac{1}{2} \log(2\pi\sigma_i^2) - \frac{1}{2\sigma_i^2}(x_i - \mu)^2 \right)$$

$$\textcircled{2} \quad \frac{d l(\theta)}{d \mu} = \sum_{i=1}^n \frac{1}{\sigma_i^2} (x_i - \mu) \stackrel{\text{want}}{=} 0$$

thus  $\sum_{i=1}^n \frac{1}{\sigma_i^2} \mu = \sum_{i=1}^n \frac{x_i}{\sigma_i^2}$

$$\hat{\mu}^{ml} = \sum_{i=1}^n \frac{x_i}{\sigma_i^2} \cdot \left( \sum_{i=1}^n \frac{1}{\sigma_i^2} \right)^{-1}$$

$$\textcircled{3} \quad \hat{\mu}^{ml} = \sum_{i=1}^n \frac{1}{\sigma_i^2 \left( \sum_{i=1}^n \frac{1}{\sigma_i^2} \right)} x_i$$

let  $w_i = \frac{1}{\sigma_i^2 \left( \sum_{i=1}^n \frac{1}{\sigma_i^2} \right)}$

let  $\tau_i = 1/\sigma_i^2$  be precision

thus weights for  $X_i$  are its precision normalized by the sum of all precision for  $X_k$   $k=\{1, \dots, n\}$

## Problem 8

moment matching is trying to find parameters such as the moments of the model curve with its estimated parameters) match the empirical moments in the observed data

uniformly distributed  $p(x) = \frac{1}{b-a}$

$$(1) E[X] = \int_a^b x p(x) dx = \frac{1}{b-a} \left. \frac{x^2}{2} \right|_a^b = \frac{1}{b-a} \frac{b^2 - a^2}{2} =$$

$$= \frac{ab}{2}$$

$$E[X^2] = \int_a^b x^2 p(x) dx = \frac{1}{b-a} \left. \frac{x^3}{3} \right|_a^b = \frac{1}{b-a} \frac{b^3 - a^3}{3}$$

$$= \frac{b^2 + ab + a^2}{3}$$

want:  $\frac{ab}{2} = E[X] = \frac{1}{N} \sum_{i=1}^N x_i$

$$\frac{b^2 + ab + a^2}{3} = E[X^2] = \frac{1}{N} \sum_{i=1}^N x_i^2$$

$$ab = (2E[X])^2 - 3E[X^2] = 4(E[X])^2 - 3E[X^2]$$

$$ab = 2E[X] \Rightarrow b = 2E[X] - a$$

$$a^2 - 2E[X]a + [4(E[X])^2 - 3E[X^2]] = 0$$

$$a = \frac{2E[X] \pm \sqrt{-12(E[X])^2 + 12E[X^2]}}{2}$$

$$= E[X] \pm \sqrt{3} \sqrt{E[X^2] - E[X]^2}$$

$$b = E[X] \mp \sqrt{3} \sqrt{E[X^2] - E[X]^2}$$

$b > 0$

$$\text{thus } a^* = \frac{1}{N} \sum_{i=1}^N x_i - \sqrt{3} \sqrt{\frac{1}{N} \sum_{i=1}^N x_i^2 - \left( \frac{1}{N} \sum_{i=1}^N x_i \right)^2}$$

$$b^* = \frac{1}{N} \sum_{i=1}^N x_i + \sqrt{3} \sqrt{\frac{1}{N} \sum_{i=1}^N x_i^2 - \left( \frac{1}{N} \sum_{i=1}^N x_i \right)^2}$$

$$\textcircled{2} \quad L(\theta) = PCD(\theta) = \prod_{i=1}^n P(X_i | a, b)$$

$$= \prod_{i=1}^n \frac{1}{b-a}$$

$$= \left( \frac{1}{b-a} \right)^n$$

$\max(L(\theta))$  is equal to minimize  $b-a$

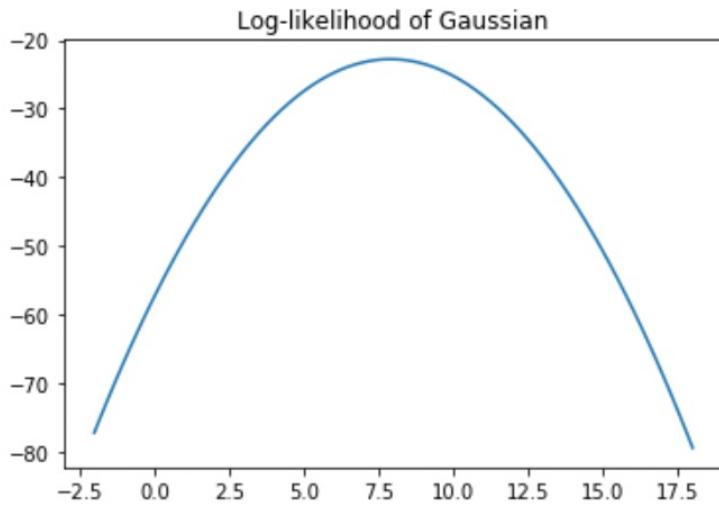
$$\text{WLOG let } \begin{aligned} b^* &= \max\{x_i\} \quad i \in \{1, \dots, N\} \\ a^* &= \min\{x_i\} \quad i \in \{1, \dots, N\} \end{aligned}$$

\textcircled{3}

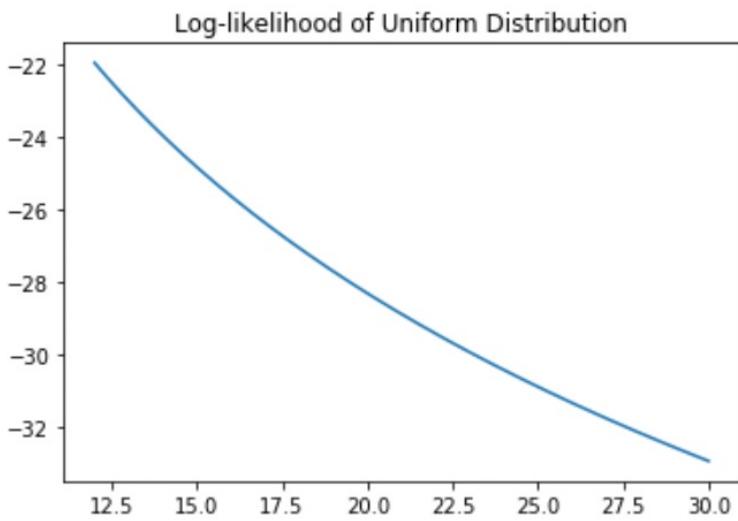
the maximum likelihood function estimator is highly depends on data's maximum and minimum however if we add more data, there may have some number less than  $a$  or larger than  $b$  which  $P(\theta)$  will give probability 0 for prediction than

method of moments gives probability range which  $a^*, b^*$  that  $a^* \leq \min(x)$ ,  $b^* \geq \max(x)$  which means that by thus  $P(\theta_{\text{moments}})$  may predict some number that even not exist in the future data observation

## Problem 2

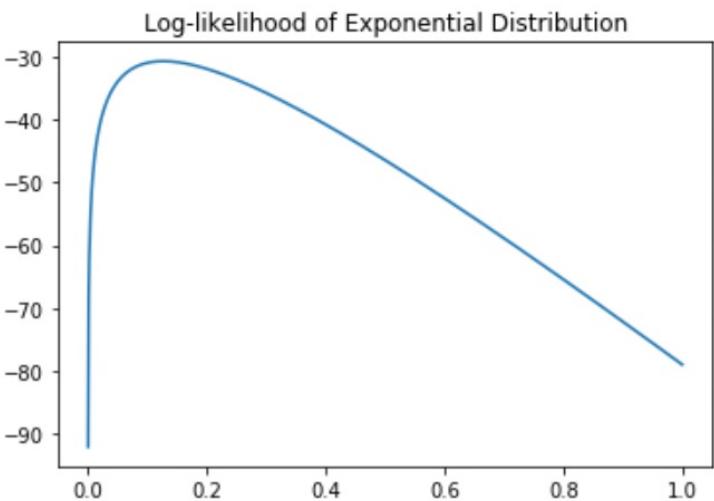


[1)



(2)

since data has maximum  
12, if  $b < 12$  then  
likelihood = 0  
log likelihood doesn't exist



## Problem 7

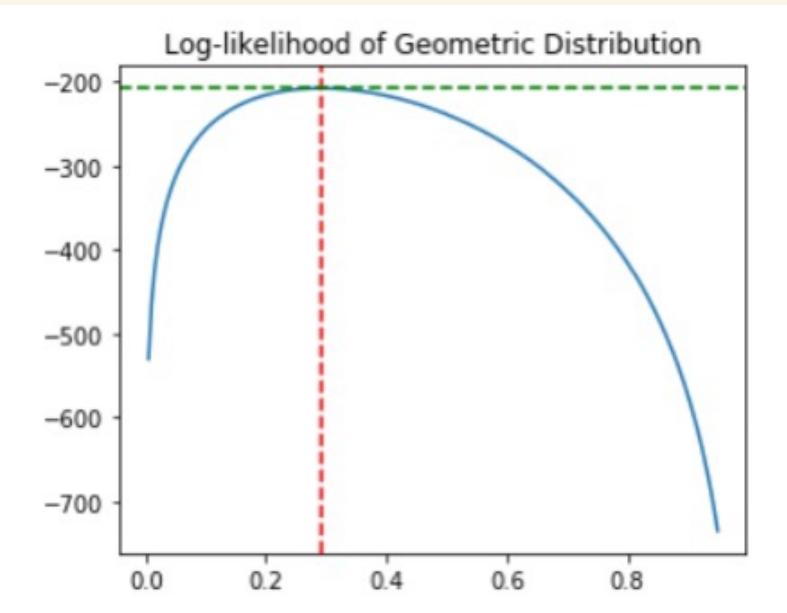
1.

Geometric Model  
Poisson model

$$P(X=k) = (1-\theta)^{k-1} \theta \quad k=0, 1, 2, \dots$$

$$P(X=k) = \frac{e^{-\lambda} \lambda^k}{k!} \quad \lambda > 0 \quad k \in \{0, 1, 2, 3, \dots\}$$

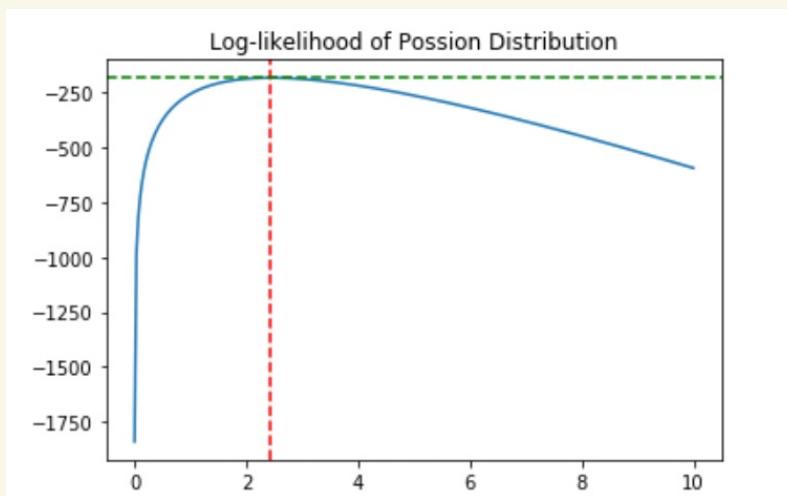
2.



Geometric Distribution

$$\hat{\theta}^{ml} = 0.29069767$$

$$\ell(\theta)_{\max} = -207.35$$

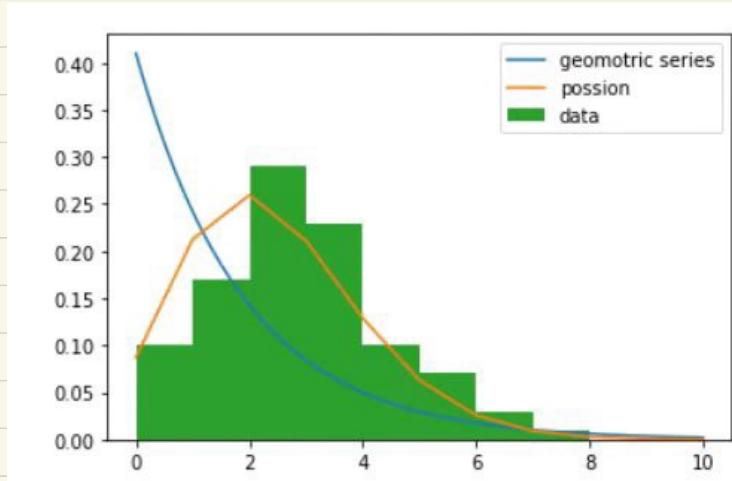


Poisson Distribution

$$\hat{\lambda}^{ml} = 2.44$$

$$\ell(\theta)_{\max} = -181.22$$

3.



4. Poisson is better to fit data both by observation and has smaller maximum likelihood between two types.

5. No. since model with two parameters has degree of freedom 2.

We can use extra-sum-of-squares F test compare or use Information theory approach Akaike's Criterion (AIC) to help compare them.

[graphpad.com/guides/prism/7/curve-fitting/reg-approaches-to-comparing-models.htm?fcid=8printWindow](http://graphpad.com/guides/prism/7/curve-fitting/reg-approaches-to-comparing-models.htm?fcid=8printWindow)