

# Estimation

Sofia Olhede



October 6, 2020

## 1 Estimation

## 2 Maximum Likelihood Estimation

# Proof of Rao–Blackwell Theorem

- To commence the proof we note that since  $T$  is assumed to be sufficient for  $\theta$ .
- To refer back to Rao-Blackwellising we note that  $\mathbb{E}\{\hat{\theta}|T = t\} = h(t)$  is independent of  $\theta$ .
- Thus  $\hat{\theta}^* = \mathbb{E}\{\hat{\theta}|T\}$  is a well-defined statistic (it is calculated from  $Y$  but is not a function of  $\theta$ ).
- We note that

$$\mathbb{E}\{\hat{\theta}^*\} = \mathbb{E}\{\mathbb{E}\{\hat{\theta}|T\}\} = \mathbb{E}\hat{\theta} = \theta.$$

- By the law of total variance we have due to the non-negativity of variances

$$\text{Var}\{\hat{\theta}\} = \text{Var}\{\mathbb{E}\{\hat{\theta}|T\}\} + \mathbb{E}\text{Var}\{\hat{\theta}|T\} \geq \text{Var}\{\mathbb{E}\{\hat{\theta}|T\}\} = \text{Var}\{\hat{\theta}^*\}.$$

# Proof of Rao–Blackwell Theorem

- We can directly note that

$$\text{Var}\{\hat{\theta}|T\} = \mathbb{E}\{(\hat{\theta} - \mathbb{E}\{\hat{\theta}|T\})^2|T\} = \mathbb{E}\{(\hat{\theta} - \hat{\theta}^*)^2|T\}.$$

?  
o

- Thus it follows that  $E \text{Var}\{\hat{\theta}|T\} = \mathbb{E}\{(\hat{\theta} - \hat{\theta}^*)^2\} > 0$  unless if  $\Pr(\hat{\theta} = \hat{\theta}^*) = 1$ .
- Now assume that  $\hat{\theta}$  is an unbiased estimator of  $\theta$  and that  $T$  and  $S$  are  $\theta$ -sufficient.
- We can also ask ourselves how do  $\text{Var}\{\mathbb{E}(\hat{\theta}|T)\}$  and  $\text{Var}\{\mathbb{E}(\hat{\theta}|S)\}$  relative?
- Our intuition would suggest that whichever of  $T$  and  $S$  carries the least irrelevant information should do better.
- Formally, if  $T = h(S)$  then  $\hat{\theta}_T^*$  should dominate  $\hat{\theta}_S^*$ .

# Variance preference for sufficient statistics

- Proposition: for  $\hat{\theta}$  an unbiased estimator of  $\theta$  and  $T$  and  $S$  sufficient statistics, define


$$\hat{\theta}_T^* = \mathbb{E}\{\hat{\theta} | T\}, \quad \& \quad \hat{\theta}_S^* = \mathbb{E}\{\hat{\theta} | S\}.$$

Then

$$T = h(S) \implies \text{Var}\{\hat{\theta}_T^*\} \leq \text{Var}\{\hat{\theta}_S^*\}.$$

- We can deduce from this result that the best way of Rao–Blackwellising is conditioning on the minimally sufficient statistic.

# Proof

- We start from the tower property of conditional expectation. If  $Y = f(X)$  then we see

$$\mathbb{E}\{Z|Y\} = \mathbb{E}\{\mathbb{E}\{Z|X\}|Y\}.$$

- Since  $T = f(S)$  we can deduce that

$$\begin{aligned}\widehat{\theta}_T^* &= \mathbb{E}\{\widehat{\theta}|T\} \\ &= \mathbb{E}\{\mathbb{E}\{\widehat{\theta}|S\}|T\} \\ &= \mathbb{E}\{\widehat{\theta}_S^*|T\}. \end{aligned} \tag{1}$$

We can now conclude the desired result using the Rao–Blackwell theorem.

- We can now judge the quality of any given estimator. For certain cases we even know how good performance we can hope for.

# ML Estimation

- How do we then come up with (good) estimators?
- We do want a method that works in general to come up with estimators.
- We even want methods that give good estimators. We shall do so using the method of maximum likelihood.
- This estimator was proposed by Ron Fisher in 1921. A more detailed exposition was provided in “On the mathematical foundations of theoretical statistics.” Philos. Trans. Roy. Soc. London Ser. A 222, 309–368.
- How do we motivate the principle of maximum likelihood? Fisher:  
*“On the bases that the purpose of the statistical reduction of data is to obtain statistics which shall contain as much as possible, ideally the whole, of the relevant information contained in the sample . . . ”.*  
Let us see how that is done in practice.

## ML Estimation II

- Statistics is about “inverse probability” (more to follow).
- Likelihood from a probability perspective: Given a parameter  $\theta \in \Theta$  then for any  $(y_1, \dots, y_n)^T \in \mathcal{Y}^n$  then we can evaluate

$$(y_1, \dots, y_n) \xrightarrow{\text{"likelihood given" }} \Pr_{\theta}(Y_1 = y_1, \dots, Y_n = y_n),$$

given a parameter  $\theta$   
 we can define a function  
 that go from  $y_1$  to  $y_n$ ,  
 and give the probability  
 of having that value

namely how the probability varies as a function of the sample.

- Likelihood from a statistics perspective: Given a sample  $(y_1, \dots, y_n)^T \in \mathcal{Y}^n$  then for any  $\theta \in \Theta$  we can calculate

$$\theta \mapsto \Pr_{\theta}(Y_1 = y_1, \dots, Y_n = y_n).$$

given a sample  $y_1 \dots y_n$ ,  
 likelihood gives us the  
 probability of the  
 observations, given  $\theta$

- How the probability varies as a function of  $\theta$  or the model.
- The maximum likelihood principle is to select a value of  $\theta$  that makes the observations most likely.

# ML Estimation III

- Definition: Let  $(Y_1, \dots, Y_n)$  be a sample of random variables with joint density/frequency  $f(y_1, \dots, y_n; \theta)$  where  $\theta \in \mathbb{R}^p$ . The likelihood of  $\theta$  is defined as

$$\rightarrow L(\theta) = f(Y_1, \dots, Y_n; \theta). \quad \begin{matrix} \text{joint density/frequency of the sample, but} \\ \text{viewed as a function of } \theta \end{matrix}$$

- If  $(Y_1, \dots, Y_n)^T$  has i.i.d. entries, each with density/frequency  $f(y_j; \theta)$  then

$$\rightarrow L(\theta) = \prod_{j=1}^n f(y_j; \theta). \quad \begin{matrix} \text{because they're independent:} \\ f(A, \delta | \theta) = f(A | \theta) \cdot f(\delta | \theta) \end{matrix}$$

We get the following estimation method.

- Definition: (Maximum Likelihood Estimation). In the same context, a maximum likelihood estimator (MLE) of  $\hat{\theta}$ ; is an estimator such that

*estimated < real*

$$L(\theta) \leq L(\hat{\theta}), \quad \forall \theta \in \Theta.$$

# ML Estimation IV

- When there exists a unique maximum, we speak of the MLE  
 $\hat{\theta} = \arg_{\theta \in \Theta} \max L(\theta)$ .
- The likelihood is a random function.
- It is the joint density/frequency of the sample, but viewed as a function of  $\theta$ .
- It is NOT the probability of  $\theta$ .
- $L(\theta)$  is the answer to the question how does the joint density / probability of the sample vary as we vary  $\theta$ .
- In the discrete case it is exactly “the probability of observing our sample” as a function of  $\theta$ .
- If a sufficient statistic  $T$  exists for  $\theta$  then Fisher–Neyman factorisation implies

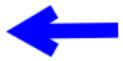
*proportional as a function of  $\theta$*

$$\rightarrow L(\theta) = g(T(Y); \theta) h(Y) \propto g(T(Y); \theta).$$

*↳ max likelihood method will only maximize this*

## ML Estimation V

a bijection, bijective function, one-to-one correspondence, or invertible function, is a function between the elements of two sets, where each element of one set is paired with exactly one element of the other set, and each element of the other set is paired with exactly one element of the first set.

- Any MLE depends on data only through a sufficient statistic.
- Since the sufficient statistic was arbitrary, if a minimally sufficient statistic exists, the MLE will have used an estimator that has achieved the maximal sufficient reduction of the data. 
- MLE's are also equivariant. Sometimes this is also known as invariance. If  $g : \Theta \mapsto \Theta'$  is a bijection, and if  $\hat{\theta}$  is the MLE of  $\theta$ , then  $g(\hat{\theta})$  is the MLE of  $g(\theta)$ .
- When the likelihood is differentiable in  $\theta$ , its maximum must solve 
$$\nabla_{\theta} L(\theta)|_{\theta=\hat{\theta}} = 0.$$
  *as always*  
$$\frac{d \text{post}}{d \theta} = 0$$
- We must verify it is a maximum (don't want the minimum likelihood solution!)

## ML Estimation VI

Hessian is a square matrix of second-order partial derivatives of a scalar-valued function, or scalar field. It describes the local curvature of a function of many variables.

- If the likelihood is twice differentiable in  $\theta$ , we can verify this by checking

$$-\nabla_{\theta}^2 L(\theta)|_{\theta=\hat{\theta}} \succ 0. \quad \leftarrow$$

- The negative of the Hessian is positive definite. In one dimension, this reduces to the standard second derivative criterion.
- Solving  $\nabla_{\theta} L(\theta)|_{\theta=\hat{\theta}} = 0$  when  $Y_i$  are independent corresponds to solving an equation after using the Leibnitz rule.
- Instead of maximising  $L(\theta)$  we take an increasing function of  $L(\theta)$  and maximise that log-likelihood
- We chose to use  $\log(x)$  as that function which is strictly increasing.
- When the  $Y_i$  are independent then  $\ell(\theta) = \log L(\theta)$  is a sum rather than a product

$$\rightarrow \ell(\theta) = \log \left( \prod_{j=1}^n f(y_j; \theta) \right) = \sum_{j=1}^n \log(f(y_j; \theta)).$$

# ML Estimation VI

- Of course, under twice differentiability, verification of a maximum can be checked again by whether or not

$$\nabla_{\theta} L(\theta)|_{\theta=\hat{\theta}} = 0, \quad -\nabla_{\theta}^2 L(\theta)|_{\theta=\hat{\theta}} \succ 0.$$



- Example (MLE for Bernoulli trials). Let  $Y_1, \dots, Y_n \stackrel{iid}{\sim} \text{Bern}(p)$ . The likelihood is

$$L(p) = \prod_{j=1}^n f(Y_j; p) = \prod_{j=1}^n p^{Y_j} (1-p)^{1-Y_j} = p^{\sum_j Y_j} (1-p)^{n-\sum_j Y_j}.$$

- This gives a loglikelihood (or a loglihood) of

$$\ell(p) = \log(p^{\sum_j Y_j} (1-p)^{n-\sum_j Y_j}) = \sum_j Y_j \log(p) + (n - \sum_j Y_j) \log(1-p).$$

## ML Estimation VII

- This is twice differentiable. We can therefore calculate

$$\frac{d}{dp} \ell(p) = \sum_j Y_j \frac{1}{p} - (n - \sum_j Y_j) \frac{1}{1-p}. \quad (2)$$

Setting this equal to 0 and solving yields  $\hat{p} = \frac{1}{n} \sum_j Y_j = \bar{Y}$ .

- We now calculate the second derivative which yields

$$\frac{d^2}{dp^2} \ell(p) = - \sum_j Y_j \frac{1}{p^2} - (n - \sum_j Y_j) \frac{1}{(1-p)^2} < 0. \quad (3)$$

- Thus it follows that

$$\hat{p} = \bar{Y},$$

the mean  
 negative value means it's  
 a maximum likelihood  
 estimator

is a unique MLE for the Bernoulli success probability.

*Because  $\hat{p}$  can take only 1 value?*

# ML Estimation VIII

- Example (MLE for Exponential trials). Let  $Y_1, \dots, Y_n \stackrel{iid}{\sim} \text{Expl}(\lambda)$ .
- The likelihood for this data takes the form

$$L(\lambda) = \prod_{j=1}^n f(Y_j; \lambda) = \prod_{j=1}^n \lambda e^{-\lambda Y_j} = \lambda^n e^{-\lambda \sum_j Y_j}.$$

- This gives a loglikelihood (or a loglihood) of

$$\ell(\lambda) = \log(L(\lambda)) = n \log \lambda - \lambda n \bar{Y}.$$

- Differentiating we find

$$\frac{d}{d\lambda} \ell(\lambda) = n\lambda^{-1} - n\bar{Y}.$$

$$\begin{aligned} \frac{d}{d\lambda} \ell(\lambda) &= 0 \\ n\lambda^{-1} &= n\bar{Y} \Rightarrow \lambda = \bar{Y}^{-1} \end{aligned}$$

Thus we take  $\hat{\lambda} = \bar{Y}^{-1}$ .

## ML Estimation IX

- Differentiating yet one more time we get

$$\frac{d^2}{d\lambda^2} \ell(\lambda) = -n\lambda^{-2} < 0.$$

- Thus we determine that

$$\hat{\lambda} = \bar{Y}^{-1},$$

is a unique MLE.

# ML Estimation X

- Example (MLE for Gaussian trials). Let  $Y_1, \dots, Y_n \stackrel{iid}{\sim} N(\mu, \sigma^2)$ .
- This process has two parameters,  $\mu \in \mathbb{R}$  and  $\sigma^2 \in \mathbb{R}^+$ .
- The likelihood is then

$$\begin{aligned} L(\mu, \sigma^2) &= \prod_{j=1}^n f(Y_j; \mu, \sigma^2) = \prod_{j=1}^n \left\{ \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{1}{2\sigma^2}(Y_j - \mu)^2\right\} \right\} \\ &= \frac{1}{(2\pi\sigma^2)^{n/2}} \exp\left\{-\frac{1}{2\sigma^2} \sum_{j=1}^n (Y_j - \mu)^2\right\}. \end{aligned}$$

- Taking logarithms we get

$$\ell(\mu, \sigma^2) = -(n/2) \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{j=1}^n (Y_j - \mu)^2.$$

We can now seek to maximise this quantity.

## ML Estimation XI

- Taking partial derivatives we get

$$\frac{\partial}{\partial \mu} \ell(\mu, \sigma^2) = 2 \frac{1}{2\sigma^2} \sum_{j=1}^n (Y_j - \mu)$$

$$\frac{\partial}{\partial \sigma^2} \ell(\mu, \sigma^2) = -(n/2) \frac{1}{\sigma^2} + \frac{1}{2\sigma^4} \sum_{j=1}^n (Y_j - \mu)^2.$$

*derivating with respect to  $\sigma^2$  not  $\sigma$*

- This leads to a set of equations that we can solve for  $\hat{\mu}$  and  $\widehat{\sigma^2}$ :

$$\begin{aligned} \frac{\partial \ell}{\partial \mu} &= 0 \Leftrightarrow \frac{1}{2\sigma^2} \sum_{j=1}^n (Y_j - \mu) = 0 \quad \leftarrow \\ &\approx n\mu - \sum_j Y_j \Rightarrow \mu = \frac{\sum_j Y_j}{n} = \bar{Y} \end{aligned}$$

- Now we need to verify that this corresponds to a maximum, by calculating second and mixed partial derivatives.

## ML Estimation XII

- Note that second partial derivatives

$$\frac{\partial^2}{\partial \mu^2} \ell(\mu, \sigma^2) = -\frac{n}{\sigma^2}, \quad \frac{\partial^2}{\partial (\sigma^2)^2} \ell(\mu, \sigma^2) = \frac{(n/2)}{\sigma^4} - \frac{2}{2\sigma^6} \sum_{j=1}^n (Y_j - \mu)^2$$

*Mixed partials sum deriving*

$$\frac{\partial^2}{\partial \mu \partial \sigma} \ell(\mu, \sigma^2) = -2 \frac{1}{\sigma^4} \sum_{j=1}^n (Y_j - \mu).$$

*The order doesn't matter*

We now evaluate the derivatives and get

$$\frac{\partial^2}{\partial \mu^2} \ell(\hat{\mu}, \widehat{\sigma^2}) = -\frac{n}{\widehat{\sigma^2}}, \quad \frac{\partial^2}{\partial (\sigma^2)^2} \ell(\hat{\mu}, \widehat{\sigma^2}) = \frac{(n/2)}{\widehat{\sigma^4}} - \frac{2n}{2\widehat{\sigma^4}} = -\frac{(n/2)}{\widehat{\sigma^4}}$$

$$\frac{\partial^2}{\partial \mu \partial \sigma} \ell(\hat{\mu}, \widehat{\sigma^2}) = 0.$$

*why?*

*how?*

Thus the Hessian is diagonal. We can also note that the negative of the matrix is positive definite and diagonal- we have a maximum.

# ML Estimation XIII

- Example (MLE for Poisson observations). Let  $Y_1, \dots, Y_n \stackrel{iid}{\sim} \text{Poisson}(\mu)$ .
- This process has one parameter,  $\mu \in \mathbb{R}^+$ .
- The likelihood is then

$$L(\mu) = \prod_{j=1}^n f(Y_j; \mu) = \prod_{j=1}^n \left\{ \frac{e^{-\mu} \mu^{Y_j}}{Y_j!} \right\} = e^{-n\mu} \mu^{\sum_{j=1}^n Y_j} \frac{1}{\prod_j Y_j!}.$$

Remember the definition of the factorial sign

$n! = \Gamma(n+1) = n(n-1)(n-2)\dots 1$ . Note that  $0! = 1$ ,  $1! = 1$ ,  $2! = 2$ ,  $3! = 6$  etc. The loglikelihood is

log-likelihood

$$\ell(\mu) = -n\mu + \sum_{j=1}^n Y_j \log(\mu) - \log \prod_j Y_j!$$

$$\frac{\partial}{\partial \mu} \ell(\mu) = -n + \sum_{j=1}^n Y_j / \mu.$$

## ML Estimation XIV

- This is zero if  $\mu = \bar{Y}$ . We check the second derivative:  

$$\frac{\partial^2}{\partial \mu^2} \ell(\mu) = - \sum_{j=1}^n Y_j / \mu^2.$$
- Example (MLE for Uniform Distribution – a non-differentiable case).  
 This example is useful as it shows that differentiation does not always work.
- The uniform distribution takes the form of  

$$f(Y_j; \theta) = \theta^{-1} I(Y_j \in (0, \theta))$$
 here. The likelihood takes the form of

remind:  
 $I(Y_j < \theta)$

$$L(\theta) = \prod_j \theta^{-1} I(Y_j \in (0, \theta)) = \theta^{-n} I(\max_j Y_j \in (0, \theta)) I(\min_j Y_j \in (0, \theta)).$$

We cannot differentiate this sensibly in theta as  $\frac{\partial}{\partial \theta} L(\theta)$  has a discontinuity. We see that  $\theta^{-n}$  has derivative  $-n\theta^{-n+1}$ , which is negative and so is decreasing in  $\theta$ . As  $\theta > Y_i$  for all  $i$ , that means we must take  $\hat{\theta} = \max_j Y_j$ .

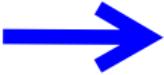
## ML Estimation XV

- The next property we shall cover is the equi-variance or invariance of MLEs. If  $g(\theta)$  is a bijection, recall that if we are attempting to estimate  $\tau = g(\theta)$  then if we form the likelihood

$\tau$  given by a bijection  $\downarrow$  on  $\theta$

$$L(\theta) = \prod_{j=1}^n f(Y_j; \theta),$$

- We can use the chain rule to determine



$$\frac{d}{d\tau} L(\tau) = \frac{d}{d\theta} L(\theta) \frac{d\theta}{d\tau},$$

and as  $g(\theta)$  is a bijection, we know that  $\frac{d\theta}{d\tau} \neq 0$ . Thus we can easily map between zeros.

# ML Estimation XVI

Invariance and equi-variance allows us to create new estimators that are also MLE estimator

- Let  $Y_1, \dots, Y_n \stackrel{iid}{\sim} N(\mu, \sigma^2)$ . Suppose we are interested in estimating  $\tau = \Pr\{Y_1 \leq y\}$ . Assume that  $y \in \mathbb{R}$  and fix  $\sigma \in \mathbb{R}^+$ .

- We note that

$$\tau(\mu) = \Pr\{Y_1 \leq y\} = \Pr\left\{\frac{Y_1 - \mu}{\sigma} \leq \frac{y - \mu}{\sigma}\right\} = \Phi\left\{\frac{y - \mu}{\sigma}\right\}.$$

remember: Pr is  
 Gaussian prob

$\hookrightarrow$  we don't know  $\mu$  but we know  $\sigma^2$  and  $y$

- We note that

to estimate  $\mu$

we derive w.r.t. the new estimator  $\hat{\tau}$  instead  $\tau(\mu)$

$\hookrightarrow$  we convert it into a standard unit Gaussian:  $N(\mu, \sigma^2) \rightarrow N(0, 1)$

$\hookrightarrow$  means Normal CDF

$$\frac{d\tau}{d\mu} = \phi\left\{\frac{y - \mu}{\sigma}\right\} \left(-\frac{1}{\sigma}\right) < 0 \quad \forall \mu.$$

Normal PDF

$\hookrightarrow$  positive to negative : function is strictly decreasing

- This implies that we have a bijective map, and so  $\hat{\tau}$  is  $\tau(\hat{\mu})$ .

slide 18: estimate of  $\hat{\mu}$  for a Gaussian is  $\bar{Y}$  is a bijective

# ML Estimation–Exponential Family

- Assume that  $Y_1, \dots, Y_n$  are all drawn from the distribution  $f$  which is specified by

$$f(y) = \exp\{\phi T(y) - \gamma(\phi) + S(y)\}, \quad y \in \mathcal{Y}.$$

- We set  $\phi$  as the natural parameter.
- Assume also that  $\phi = \eta(\theta)$ , where  $\theta \in \Theta$  is the usual parameter.
- Assume that  $\eta : \Theta \rightarrow \Phi$  is a differentiable bijection. Thus  
 $\gamma(\phi) = \gamma(\eta(\theta)) = d(\theta)$ .  $d = \gamma \circ \eta$ . *is any transformation  $\delta$  on the natural param.  $\phi$  matches a transf. of transf. ( $\delta \circ \eta$ ) on usual param  $\Theta$*
- Thus the density can be rewritten:

  $f(y) = \exp\{\phi T(y) - \gamma(\phi) + S(y)\} = \exp\{\eta(\theta)T(y) - d(\theta) + S(y)\}.$

- Equivariance implies that if  $\hat{\theta}$  is an MLE for  $\theta$  then so is  $\eta(\hat{\theta})$  for  $\phi = \eta(\theta)$ . If  $\hat{\phi}$  is an MLE of  $\phi$  then  $\eta^{-1}(\hat{\phi})$  is an MLE of  $\eta^{-1}(\phi)$ .
- also can compute MLE in either natural and usual param*

# ML Estimation–Exponential Family

*→ in the estimator converges in probability to the true value*

- **Consistency** of MLE in  $\theta \in \mathbb{R}$ . Let  $Y_1, \dots, Y_n \sim f(y; \theta_0)$  where  $f \in C^1$  wrt to  $\theta$ . Assume that for all  $n$  there is an unique MLE  $\hat{\theta}_n$ . We will show  $\hat{\theta}_n \xrightarrow{P} \theta$ .
- Define

$$\Xi_n(u) = \frac{1}{n} \sum_{i=1}^n \left[ \frac{\partial}{\partial u} \log \left( \frac{f(Y_i; u)}{f(Y_i; \theta)} \right) \right], \quad \Xi(u) = \mathbb{E} \left[ \frac{\partial}{\partial u} \log \left( \frac{f(Y_i; u)}{f(Y_i; \theta)} \right) \right],$$

- so that  $\Xi_n(\hat{\theta}_n) = 0$  by uniqueness of MLEs;
  - $\Xi(\theta) = 0$  uniquely, assuming regularity allowing interchange of  $\mathbb{E}()$  and  $\frac{\partial}{\partial u}$ .
  - The law of large numbers implies that  $\Xi_n(u) \xrightarrow{P} \Xi(u)$ .
  - This implies  $\hat{\theta}_n$  converges to  $\theta$ . Consistency is equivalent to this statement.
- Law of Large Numbers: the average of the results obtained from a large number of trials should be close to the expected value and will tend to become closer to the expected value as more trials are performed. Weak law of large numbers (Khinchin's law): the sample average converges in probability towards the expected value

# ML Estimation–Exponential Family

- Consistency of MLE in  $\mathbb{R}^k$  for exponential families.
- Consider  $Y_1, \dots, Y_n \sim f(y; \phi)$  from a  $k$  exponential family:

$$f(y) = \exp\left\{\sum_{j=1}^k \phi_j T_j(y) - \gamma(\phi) + S(y)\right\}.$$

- The likelihood and log-likelihoods are given by

$$L(\phi) = \exp\{\phi^T \tau - n\gamma(\phi)\}, \quad \ell(\phi) = \phi^T \tau - n\gamma(\phi).$$

where  $\tau = (\tau_1, \dots, \tau_k)^T$  and  $\tau_j(\mathbf{y}) = \sum_i T_j(y_i)$ .

- If it exists the MLE  $\hat{\phi}$  must satisfy

$$\nabla_\phi \ell(\hat{\phi}_n) = 0 \Rightarrow \nabla_\phi \gamma(\hat{\phi}_n) = n^{-1} \tau.$$

Furthermore, existence of the MLE guarantees uniqueness by strict concavity.

# ML Estimation–Exponential Family

- Provided it exists, the MLE of the natural parameter in a  $k$ -parameter natural exponential family with open parameter space  $\Phi$  is consistent.
  - Assuming we can get consistency, we can focus on understanding the sampling distribution of the MLE.
  - For simplicity, assume  $X_1, \dots, X_n$  are iid with density/frequency)  $f(x; \theta)$  for  $\theta \in \Theta$ . Write
  - Let  $\ell(x_i; \theta) = \log f(x_i; \theta)$ .
  - Let  $\ell'(x_i; \theta)$ ,  $\ell''(x_i; \theta)$  and  $\ell'''(x_i; \theta)$  denote the partial derivatives wrt  $\theta$ .
  - We need some regularity conditions:
  - A1:  $\Theta$  is an open subset of  $\mathbb{R}$ .
  - A2: The support of  $f$ ,  $\text{supp } f$  is independent of  $\theta$ .
  - A3:  $f$  is thrice continuously differentiable w.r.t.  $\theta$  for all the support of  $f$ .
- Important:  
to do: review!*
- ?

# ML Estimation–Exponential Family

- A4:  $\mathbb{E}(\ell'(X_i; \theta)) = 0$  for all  $\theta$  and  $\text{Var}\{\ell'(X_i; \theta)\} = \mathcal{I}_1(\theta)$ .
- A5:  $-\mathbb{E}(\ell''(X_i; \theta)) = \mathcal{J}_1(\theta)$ .
- A6:  $\exists M(x) > 0$  and  $\delta > 0$  such that  $\mathbb{E} M(x) < \infty$  and

?  
o

$$|\theta - \theta_0| < \delta \Rightarrow |\ell'''(x; \theta)| < M(x).$$