

Notes

MA4107: Statistical Inference

Sabarno Saha

ss22ms037@iiserkol.ac.in

Date: 03 September 2025

1. Data and Models

Definition 1.1 (Data): Let G generate a vector in \mathbb{R}^n according to some model.

$$G : \Omega \longrightarrow \mathbb{R}^n \quad (1.1)$$

Then \tilde{X} is called a realization of the data $G(\omega)$ for some $\omega \in \Omega$

\tilde{X} comes from some distribution or model F . This course is all about parametric models.

Definition 1.2 (Parametric Model): Let F be some parameter model. Then

$$F \in \mathfrak{F} = \{F_\theta : \theta \in \Theta \subset \mathbb{R}^k \text{ for some } k \in \mathbb{N}\} \quad (1.2)$$

Θ is called the parameter space.

We define g to be the map that maps a set of parameters to a model. g is onto. For every F_θ , there always exists θ , such that

$$\theta \xrightarrow{g} F_\theta \quad (1.3)$$

g also needs to be injective, otherwise we run into the problem of identifiability. The problem of identifiability arises when we cannot infer the parameters from the model F_θ .

Unless explicitly mentioned otherwise, the data is represented by \tilde{X} . For all models in this course, we assume the following,

1. The parametrization is bijective.
2. The models $F_\theta \in \mathfrak{F}$ are all either discrete or continuous, not a mixture of both.

2. Statistic

Definition 2.1 (Statistic): Let the data be $\tilde{X} \in \mathbb{R}^n$. A statistic $T = T(\tilde{X})$ is a measurable function of data and data only.

Definition 2.2 (Ancillary Statistic): Suppose the distribution of the test statistic $T(\tilde{X})$, \mathfrak{T} is independent of the parameter vector θ . Then $T(\tilde{X})$ is called an Ancillary Statistic.

Theorem 2.1:

Let $f(x)$ be a pdf. Let μ and $\sigma > 0$ be any constants. Then

$$g(x|\mu, \sigma) = \frac{1}{\sigma} f\left(\frac{x - \mu}{\sigma}\right) \quad (2.1)$$

is a valid pdf.

The proof is in Casella and burger theorem 3.5.1. One just needs to show non-negativity and normalization.

Definition 2.3 (Location Family of distributions):

Let a family of distributions be,

$$\mathfrak{F} = \{f_\theta : \theta \in \Theta \text{ where } f_\theta(x) = g(x - \theta) \text{ for some known function } g \text{ on } \mathbb{R}^n\} \quad (2.2)$$

Then \mathfrak{F} is a location family of distributions with the standard pdf $f(x)$ and the location parameter θ for the family.

Generally we will talk about these families on \mathbb{R} rather than \mathbb{R}^n .

Definition 2.4 (Scale family of distributions):

Let a family of distributions be,

$$\mathfrak{F} = \left\{ f_\theta : \theta \in \Theta \text{ where } f_\theta(x) = \frac{1}{\theta} g\left(\frac{x}{\theta}\right) \text{ for some known function } g \text{ on } \mathbb{R}^n \right\} \quad (2.3)$$

Then \mathfrak{F} is a scale family of distributions with the standard pdf $f(x)$ and the scale parameter θ for the family.

Definition 2.5 (Location-Scale family of distribution):

Let a family of distributions be,

$$\mathfrak{F} = \left\{ f_{\mu, \sigma} : (\mu, \sigma) \in \mathbb{R} \times \mathbb{R}^+ \text{ where } f_{\mu, \sigma}(x) = \frac{1}{\sigma} g\left(\frac{x - \mu}{\sigma}\right) \right. \\ \left. \text{for some known function } g \text{ on } \mathbb{R}^n \right\} \quad (2.4)$$

Then \mathfrak{F} is a location family of distributions with the standard pdf $f(x)$ and the location parameter μ and scale parameter σ for the family.

Now comes to the choice of test statistics. Note that ancillary statistics are not useful in inferring parameters from data. So we shall not use them. The test statistic just compresses the data, from \tilde{X} to $T(\tilde{X})$. We want to see when this data compression is lossless.

We can define the level sets of the test statistic $T(X)$. Let D be the set on which the data can lie in. Let us define the relation,

$$x \sim_T y \text{ if and only if } T(x) = T(y) \quad (2.5)$$

We can easily see that this is an equivalence relation. Then D can be partitioned into sets D_t s.t.,

$$D_t = \{X \in D : T(X) = t\} \quad (2.6)$$

These sets D_t are called level sets of the test Statistic T . If we are in some level set of the test statistic, the inferred parameters should not change if we stay in the same level set.

Definition 2.6 (Sufficient Statistic):

The test statistic $T = T(\tilde{X})$ is said to be sufficient iff the distribution of $(\tilde{X}|T = t)$ is free of θ for all values of t .

Theorem 2.2 (Fisher Neyman Factorization theorem):

Suppose \tilde{X} is an iid sample from $f_\theta(\cdot)$, which might be either a pmf or a pdf. A statistic $T = T(\tilde{X})$ is sufficient for θ if and only if,

$$f_\theta(\tilde{X}) = g(T(\tilde{X}), \theta) h(x) \quad (2.7)$$

where $\theta \in \Theta \subset \mathbb{R}^k$, and the functions $g(\cdot)$ and $h(\cdot)$ are non negative functions and $h(\cdot)$ is independent of θ .

In principle we can get multiple sufficient statistics, with finer and finer partitions of D . We try to find the statistic which is sufficient and has the coarsest partitions in D .

Definition 2.7 (Minimal Sufficient Statistic):

A statistic T is said to be minimal sufficient for θ iff it is sufficient for θ and for any other sufficient statistic $S \exists$ a function g such that,

$$T = g(S) \quad (2.8)$$

Lemma 2.3: If T_1 and T_2 are minimal sufficient statistics, there exist injective functions g_1 and g_2 such that,

$$T_1 = g_2(T_2) \quad T_2 = g_1(T_1) \quad (2.9)$$

INSERT PROOF.

Theorem 2.4 (Characterization of minimal sufficiency):

Let \tilde{X} be a joint pmf/pdf $f_\theta(x)$ and $T = T(\tilde{X})$ be a statistic. Suppose the following property holds,

$$\frac{f_\theta(x)}{f_\theta(y)} \text{ is free of } \theta \text{ iff } T(x) = T(y) \quad \forall x, y \text{ s.t. } f_\theta(y) \neq 0 \quad (2.10)$$

Then T is minimal sufficient for θ .

Proof: Assume for simplicity that $f_\theta(x) > 0 \quad \forall x \in \mathbb{R}^n$ and $\forall \theta \in \Theta$. We need to show two things,

1. T is sufficient for θ
2. For any other sufficient statistic S , \exists a function g such that $T = g(S)$.

1. Sufficiency of T :

Let $\mathfrak{T} = \{T(y) : y \in \mathbb{R}^n\}$ be the image of T . Let $\{A_t : t \in \mathfrak{T}\}$ be the level sets of T . We can pick an appropriate y^* such that,

$$f_\theta(x) = f_\theta(y^*) \frac{f_\theta(x)}{f_\theta(y^*)} = g(T(x), \theta) h(x) \quad (2.11)$$

We choose y^* such that $T(y^*) = T(x)$, i.e. we choose y^* from the level set $A_{T(x)}$. Then by Fisher Neyman factorization theorem [Theorem 2.2](#), T is sufficient for θ . ■

Note a sufficient statistic can also contain garbage information. A sufficient statistic can be paired with an ancillary statistic and the combined statistic would still be sufficient.

Definition 2.8 (Complete Statistic):

Let T be a statistic. Let $\{g_\theta(T = t) : \theta \in \Theta\}$ be a family of pdf/pmf of T . The statistic T is **complete** if given any measurable function h the following holds for all $\theta \in \Theta$,

$$\mathbb{E}_\theta[h(T)] = 0 \implies P_\theta(h(T) = 0) = 1 \quad (2.12)$$

Unfortunately we have no easy characterization of complete statistic, like [Theorem 2.4](#) for minimal sufficiency.

Theorem 2.5 (Basu's Theorem):

A complete sufficient statistic is independent of any ancillary statistic.

Theorem 2.6 (Lehmann Scheffe):

Let X have a joint pdf/pmf. If T is a complete sufficient statistic for θ , then T is minimal sufficient.

The reverse implication does not hold and an example shall be provided below.

3. Families of distribution

Definition 3.1 (Exponential family of distribution):

Let X have a joint pdf/pmf $f_\theta(\cdot)$ with $\theta \in \Theta \subset \mathbb{R}^p$. We say that the $f_\theta(\cdot)$ belongs to the k -parameter exponential family if $f_\theta(\cdot)$ admits the functional form,

$$f_\theta(X) = \exp \left[\left\{ \sum_{j=1}^k c_j(\theta) T_j(X) \right\} - d(\theta) \right] S(X) \quad (3.1)$$

for all $x \in \mathfrak{X}, \forall \theta \in \Theta$. \mathfrak{X} is the space of all values taken by X .

Equivalently the last restriction can also be rewritten as that the support of f ,

$$\text{supp}(f) = \{x : f_\theta(x) > 0\} \quad (3.2)$$

is free of θ .

c_j 's are called Natural or canonical parameters

T_j 's are called Natural or canonical statistics

The expression (3.1) needs to have some restrictions placed on it, otherwise we can do simple algebraic manipulations such that the upper limit in the sum k increases, which causes a pdf $f_\theta(\cdot)$ to belong to multiple parameter exponential families.

The expression is said to be minimal in the sense that the expression cannot be reduced further without breaking the functional form (3.1).

1. $c_j(\theta)$ must explicitly depend on theta. If some c_m is constant and free of theta, the corresponding $T_m(X)$ can be absorbed into $S(X)$.
2. If $i \neq j$, $c_i \neq c_j$. This means that no two c_j 's can be the same. If $c_i = c_j$ for some $i \neq j$, we can define a combined statistic $\tilde{T} = T_i + T_j$.

Note that $\theta \in \Theta \subset \mathbb{R}^p$. To avoid the problem of identifiability, $k \geq p$. If $k < p$ we will lose information and not have an injective map. Generally $k = p$ for most cases. Later we will see an explicit example where $k > p$.

Result 1 (Sufficiency of Canonical statistic): The statistic $T = (T_1, \dots, T_k)$, where $T_j \forall j \in \{1, \dots, k\}$ are the canonical statistics defined in [Definition 3.1](#), is sufficient for θ .

Proof: The proof of this is obvious. Using [Theorem 2.2](#) and (3.1), we can see that

$$\begin{aligned} f_{\theta}(\tilde{X}) &= \exp \left[\left\{ \sum_{j=1}^k c_j(\theta) T_j(\tilde{X}) \right\} - d(\theta) \right] S(\tilde{X}) \\ &= g(T(\tilde{X}), \theta) h(\tilde{X}) \end{aligned} \quad (3.3)$$

Note that the pdf $f_{\theta}(x)$ and the function $g(T(\tilde{X}), \theta)$ are non-negative. So $h(\tilde{X}) = S(\tilde{X})$ must also be non negative. Thus using [Theorem 2.2](#), $T = (T_1, \dots, T_k)$ is a sufficient statistic for θ . ■

Theorem 3.1 (Completeness of Canonical Statistics for the Exponential family):

Consider a k -parameter exponential family. Define the natural parameter space to be ,

$$C := \{(c_1(\theta), \dots, c_k(\theta)) : \theta \in \Theta\} \quad (3.4)$$

where Θ is the parameter space. If C contains an open set in \mathbb{R}^k , then the statistic $T(\tilde{X})$ is complete. Hence T is also minimal sufficient.

TRY TO SHOW.

We cannot drop the open set restriction on the parameter space due to the following counter example. Moreover this counter example also serves as an example for the $k > p$ case outlined above. The canonical statistic vector here is also minimal sufficient but not complete.

Example (Curved Normal distribution):

Let $\tilde{X} = (X_1, \dots, X_n) \stackrel{iid}{\sim} \mathcal{N}(\theta, \theta^2), \theta > 0$. The pdf is then given by,

$$f_{\theta}(\tilde{X}) = \frac{e^{-\frac{n}{2\theta^2}}}{(2\pi)^{\frac{n}{2}}} \exp \left[-\frac{1}{2\theta^2} \sum_{i=1}^n x_i^2 + \frac{1}{\theta} \sum_{i=1}^n x_i - n \ln(\theta) \right] \quad (3.5)$$

The canonical parameters and canonical statistics are given by,

$$\begin{aligned} c_1(\theta) &= -\frac{1}{2\theta^2} & T_1(\tilde{X}) &= \sum_{i=1}^n x_i^2 \\ c_2(\theta) &= \frac{1}{\theta} & T_2(\tilde{X}) &= \sum_{i=1}^n x_i \end{aligned} \quad (3.6)$$

The natural parameter space is then defined as,

$$C = \left\{ \left(-\frac{1}{2\theta^2}, \frac{1}{\theta} \right) : \theta > 0 \right\} \quad (3.7)$$

This forms a graph in \mathbb{R}^2 and thus this does not contain an open set in \mathbb{R}^2 . We will now show that

$$T(\tilde{X}) = (T_1(\tilde{X}), T_2(\tilde{X})) = \left(\sum_{i=1}^n x_i^2, \sum_{i=1}^n x_i \right) \quad (3.8)$$

is not complete. To do this we explicitly construct a function,

$$h(T) = h(t_1, t_2) = \frac{t_2^2}{n + n^2} - \frac{t_1}{2n} \quad (3.9)$$

Note that in general $h(T) \neq 0$. Then we get ,

$$\begin{aligned} \mathbb{E}_\theta(h(T)) &= \frac{1}{n^2 + n} \mathbb{E}_\theta \left(\left(\sum_{i=1}^n x_i \right)^2 \right) - \frac{1}{2n} \mathbb{E}_\theta \left(\sum_{i=1}^n x_i^2 \right) \\ &= \left[\frac{1}{n^2 + n} - \frac{1}{2n} \right] \mathbb{E}_\theta \left(\sum_{i=1}^n x_i^2 \right) + \frac{1}{n^2 + n} \mathbb{E}_\theta \left(\sum_{i,j;i \neq j} x_i x_j \right) \\ &= \theta^2 \left[\frac{1 - n}{1 + n} \right] + \frac{n(n-1)}{n^2 + n} \theta^2 = 0 \end{aligned} \quad (3.10)$$

Thus T is not complete. Also we can very easily show using the functional form of the curved normal distribution, that T satisfies [Theorem 2.4](#). Thus T is a minimal statistic. This serves as an example why the reverse implication in [Theorem 2.6](#) does not hold.

This also serves as an example where $k > p$. The dimension of the parameter space $p = 1$ and this belongs to a $k = 2$ parameter exponential family.

4. Optimal Estimation

We now need to quantify the “error of approximation” to select among potential candidates for statistics.

Definition 4.1 (Mean Squared Error):

Let T be a test statistic. Let $\theta \in \Theta$ be a parameter. Then the mean squared error, function of the parameter θ , is given by,

$$\text{MSE}_\theta(T) = \mathbb{E}_\theta[(T - \theta)^2] \quad (4.1)$$

Theorem 4.1 (Bias Variance Decomposition): The Mean squared error can be decomposed into,

$$\text{MSE}_\theta(T) = \mathbb{E}_\theta[(T - \theta)^2] = \left[\underbrace{\mathbb{E}_\theta(T) - \theta}_{\text{Bias}_\theta(T)} \right]^2 + \text{Var}_\theta(T) \quad (4.2)$$

This is called the Bias-Variance decomposition.

Proof: This is trivial as,

$$\mathbb{E}_\theta[(T - \theta)^2] = \text{Var}_\theta(T - \theta) + [\mathbb{E}_\theta(T) - \mathbb{E}_\theta(\theta)]^2 = \text{Var}_\theta(T) + [\mathbb{E}_\theta(T) - \theta]^2 \quad (4.3)$$

■

To find an optimal statistic, we need to find one such that the MSE is minimized uniformly in θ over the parameter space Θ . Unfortunately this is not possible due to one simple fact. Just fix some $\theta_a \in \Theta$. Then define the statistic to be $T = \theta_a$. This is a valid statistic. Any statistic chosen will be worse than this in some neighbourhood of θ_a . (Is MSE a continuous function of θ ?)

We have some possible alternatives to error quantifiers, and look for one dimensional summaries of the curve $\theta \mapsto \text{MSE}_\theta$

We can have the

1. Bayes Approach :

$$R_1(T) = \int_{\Theta} \text{MSE}_\theta(T) \omega(\theta) d\theta \quad (4.4)$$

$\omega(\theta)$ is called the prior.

2. Minimax Approach:

$$R_2(T) = \sup_{\theta \in \Theta} \text{MSE}_\theta(T) \quad (4.5)$$

This quantifies the worst possible scenario.

We will be heading in a different direction and only consider the case of parameters with $\text{Bias}_\theta(T) = 0 \forall \theta \in \Theta$.

Definition 4.2 (Unbiased Estimator): The statistic T is said to unbiased for θ if

$$\mathbb{E}_\theta[T] = \theta \forall \theta \in \Theta \quad (4.6)$$

Our aim is to find T^* among the class of unbiased estimators such that the mean squared error is minimized, which means the following property holds,

$$\underbrace{\text{MSE}_\theta(T^*)}_{\text{Var}_\theta(T^*)} \leq \underbrace{\text{MSE}_\theta(T)}_{\text{Var}_\theta(T)} \forall \theta \in \Theta, \forall T \text{ unbiased for } \theta \quad (4.7)$$

If such a T^* exists, it is called the Uniform Minimum Variance Unbiased Estimator (UMVUE).

Some comments on unbiased estimators,

1. An unbiased estimator need not exist.
2. For a given class of estimators, the MSE error need not be the lowest for an unbiased estimator.

Unbiased estimators might not very good for certain purposes. In general this approach is taken so that the MSE can be minimized over all values of θ in Θ , provided we restrict ourselves to unbiased estimators.

Theorem 4.2 (Rao-Blackwell theorem):

Let X have a joint pdf/pmf $f_\theta(\cdot)$ with $\theta \in \Theta \subset \mathbb{R}^p$. Let $T = T(\tilde{X})$ be a sufficient statistic for θ . Let $S = S(\tilde{X})$ be a statistic such that

1. $\mathbb{E}_\theta[S] = \theta \quad \forall \theta \in \Theta$
2. $\text{Var}_\theta(S) < \infty \quad \forall \theta \in \Theta$

Define the new statistic,

$$S^* = \mathbb{E}_\theta[S \mid T] \quad (4.8)$$

This is called the Rao-Blackwellization of S with respect to T .

Then the following holds,

1. S^* is unbiased for θ
2. $\text{Var}_\theta(S^*) \leq \text{Var}_\theta(S) \quad \forall \theta \in \Theta$

with equality if and only if $P_\theta(S^* = S) = 1 \quad \forall \theta \in \Theta$.

Proof: Note that S^* is a statistic as it is a conditional expectation on the sufficient statistic T . To show unbiasedness,

$$\mathbb{E}_\theta[S^*] = \mathbb{E}_\theta[\mathbb{E}_\theta[S \mid T]] = \mathbb{E}_\theta[S] = \theta \quad (4.9)$$

using the law of total expectation.

To show the variance reduction, we use a formula for the conditional variance.

Let X, Y be random variables, with $\mathbb{E}[Y^2] < \infty$. Then,

$$\text{Var}(Y) = \mathbb{E}[\text{Var}(Y \mid X)] + \text{Var}(\mathbb{E}[Y \mid X]) \geq \text{Var}(\mathbb{E}[Y \mid X]) \quad (4.10)$$

Since $\text{Var}(Y \mid X) \geq 0$. (WILL WRITE EQUALITY PROOF LATER) ■

The above theorem leaves us with two questions,

- What unbiased estimator S should we start with to minimize MSE ?
- What sufficient statistic T should we use to minimize MSE ?

Theorem 4.3 (Minimal Statistic for unbiased estimation):

Let S be an unbiased estimator for θ and let T_1 and T_2 be two sufficient statistics for θ . Define the Rao-Blackwellizations,

$$\begin{aligned} S_1 &= \mathbb{E}_\theta[S \mid T_1] \\ S_2 &= \mathbb{E}_\theta[S \mid T_2] \end{aligned} \quad (4.11)$$

If $T_1 = h(T_2)$, for some function h , then

$$\text{Var}_\theta(S_1) \leq \text{Var}_\theta(S_2) \quad \forall \theta \in \Theta \quad (4.12)$$

Thus if a minimal sufficient statistic T exists, the Rao-Blackwellization with respect to T is the best among all Rao-Blackwellizations with respect to any other sufficient statistic.

Proof: To show this we use the tower property of conditional expectation. If $Y = f(X)$, then

$$\mathbb{E}(Z|Y) = \mathbb{E}(\mathbb{E}(Z|X)|Y) \quad (4.13)$$

Since $T_1 = h(T_2)$,

$$\begin{aligned} S_1^* &= \mathbb{E}(S|T_1) \\ &= \mathbb{E}(\mathbb{E}(S|T_2)|T_1) \\ &= \mathbb{E}(S_2^*|T_1) \end{aligned} \quad (4.14)$$

Using [Theorem 4.2](#),

$$\text{Var}_\theta(S_1^*) = \text{Var}_\theta(\mathbb{E}(S_2^*|T_1)) \leq \text{Var}_\theta(S_2^*) \quad \forall \theta \in \Theta \quad (4.15)$$

■

Theorem 4.4 (Choice of Unbiased estimator): Assume that T is complete sufficient for θ . Let S_1 and S_2 be two unbiased estimators for θ . Define the Rao-Blackwellizations,

$$\begin{aligned} S_1^* &= \mathbb{E}_\theta[S_1 | T] \\ S_2^* &= \mathbb{E}_\theta[S_2 | T] \end{aligned} \quad (4.16)$$

Then

$$P_\theta(S_1^* = S_2^*) = 1 \quad \forall \theta \in \Theta \quad (4.17)$$

Proof: Note that both S_1^* and S_2^* are unbiased for θ . Thus

$$\mathbb{E}_\theta[S_1^* - S_2^*] = 0 \quad \forall \theta \in \Theta \quad (4.18)$$

We have achieved first order ancillarity. S_1^* and S_2^* are functions of the complete sufficient statistic T . Thus using the definition of completeness,

$$P_\theta(S_1^* - S_2^* = 0) = 1 \quad \forall \theta \in \Theta \quad (4.19)$$

■

This theorem shows that the choice of unbiased estimator does not matter if we have a complete sufficient statistic. The Rao-Blackwellization will always be the same.

Theorem 4.5 (Lehmann-Scheffe for UMVUE): Let T be a complete sufficient statistic for θ . Let S be any unbiased estimator for θ , such that $\text{Var}_\theta(S) < \infty \forall \theta \in \Theta$. Then the Rao-Blackwellization,

$$S^* = \mathbb{E}_\theta[S \mid T] \quad (4.20)$$

Then,

1. $\text{Var}_\theta(S^*) \leq \text{Var}_\theta(U) \forall \theta \in \Theta, \forall U$ unbiased for θ
2. If for some unbiased estimator U , $\text{Var}_\theta(S^*) = \text{Var}_\theta(U) \forall \theta \in \Theta$, then

$$P_\theta(U = S^*) = 1 \forall \theta \in \Theta \quad (4.21)$$

Thus S^* is the unique UMVUE for θ .

This theorem is the culmination of the previous theorems in this section. If a complete sufficient statistic exists, we can find the UMVUE for θ by Rao-Blackwellizing any unbiased estimator with respect to the complete sufficient statistic.

Remark:

1. If a complete sufficient statistic T exists, the UMVUE is a function of T .
2. If the UMVUE exists, it is unique.
3. This theorem can be used to check whether unbiased estimators exist or not. If we can find a complete sufficient statistic T and show that no function of T is unbiased for θ , then no unbiased estimator for θ exists. This means that the UMVUE exists iff an unbiased estimator exists.

In general, we have two approaches to find the UMVUE for θ . First find a complete sufficient statistic T for θ .

1. Find some unbiased estimator S for θ and Rao-Blackwellize it with respect to T .
2. Find a function of T which is unbiased for θ . We also need check that the function has finite variance.