

Estimating the sampling distribution of the maximum likelihood estimator of the parameters of a series system given a sample of masked system failures

Alexander Towell
atowell@siue.edu

Abstract

First, the *parametric family* of an abstract *series system* is derived. Then, the *asymptotic sampling distribution* (*multivariate normal*) of a *maximum likelihood estimator* of the true parameter is derived conditioned on the *information* given by a random sample of *masked system failure times*. The asymptotic sampling distributions of statistics that are functions of the parameter follow as a result. Finally, these results are applied to series systems in which component lifetimes are *exponentially* distributed, which lead to closed formulas on a set of *minimally sufficient statistics*.

Contents

Contents	1
1 Introduction	2
2 Probabilistic model	2
2.1 Series system	3
2.2 Probability distribution functions and parametric families	4
2.3 α -masked component failure	5
3 Derivation of parametric probability distribution functions	8
3.1 System lifetime	8
3.2 Component cause of system failure	10
3.3 α -masked component failures	12
4 Likelihood and Fisher information	15
5 Sampling distribution of parameter estimators	20
5.1 Maximum likelihood estimator for α -masked system failure times each consisting of w components	21
5.2 Maximum likelihood estimator	24
Bibliography	27

Appendix	27
A Alternative proof of equation 3.6	27
B Alternative proof of equation ??	28
C Proof of corollary 3.7.5	30
D Numerical solutions to the MLE	31
E General joint distribution	32
F Sampling distribution of functions of parameters	33
G Bootstrap of sample covariance	33

1 Introduction

Empirical modeling involves a blending of substantive subject matter and statistical information. In our case, the substantive information includes assumptions like a series system and the α -masked candidate set models described later. Such information suggests the relevant variables and simplifies learning from data (in our case, maximum likelihood estimation).

In contrast, statistical information stems from the chance regularities in data, which are statistical model attempts to “adequately” capture.

Discuss latent or hidden variables.

Draw a graph in prob. model section. What do we know, how we know it, etc.

We consider series systems consisting of some fixed number of components with uncertain lifetimes. We desire a mathematical model of the system that may be used to predict the failure time of the system and its component cause. However, since the system lifetime and the component cause of failure is uncertain, we are interested in probabilistically modeling the system lifetime, e.g., there is a 75% chance that component 3 will cause a system failure in the next 3 years.

We assume the system belongs to some parametric family but whose true parameter index is unknown. Furthermore, we also assume there is a sample of *masked system failure times* that carry information about the true parameter index.

Using the information in the sample, we employ the frequentist approach of repeated experiments which induces a sampling distribution on the maximum likelihood point estimator of the fixed but unknown true parameter index where *a priori* any value in the parameter space is equally likely. The point estimator converges to a multivariate normal sampling distribution with a mean given by the true parameter index and a variance that is inversely proportional to sample size.

2 Probabilistic model

A probabilistic model describes how observable data about the system is generated. In what follows, we describe our model and observable data.

A *set* is an unordered collection of distinct elements. The set of all subsets of a set is called the *powerset*. The powerset of \mathcal{X} is denoted by $2^{\mathcal{X}}$. An interval in $2^{\mathbb{R}}$ is a set of all the real numbers lying between two indicated real numbers, e.g.,

$$(a, b] := \{x \in \mathbb{R} \mid x \geq a \wedge x \leq b\} \in 2^{\mathbb{R}}. \quad (2.1)$$

Definition 2.1. A probability space (Ω, \mathcal{F}, P) , where $P: 2^\Omega \mapsto [0, 1]$ is the probability set function, Ω is the sample space, and \mathcal{F} is the set of measurable events (subsets of Ω).

A random variable is a measurable function defined on a probability space.

Definition 2.2. Given a probability space (Ω, \mathcal{F}, P) , a random variable X is a function of type $\Omega \mapsto \text{codom}(X)$.

The range (or image) of X is also called the space of X .

For a random variable X , we use $X = x$ to denote the subset containing all elements $e \in \Omega$ that X maps to the value of x .

2.1 Series system

The system consists of m components. We make the following assumption about the state of the components.

Assumption 1. The j^{th} component is initially in a non-failed state and permanently enters a failed state at some uncertain time $t_j > 0$ for $j = 1, \dots, m$.

Components have uncertain lifetimes as given by the following definition.

Definition 2.3. The lifetime of the j^{th} component is given by the continuous random variable $T_j \in (0, \infty)$ for $j = 1, \dots, m$.

We make the following assumption about the joint distribution of the component lifetimes.

Assumption 2. The random variables T_1, \dots, T_m are mutually independent.

The system is a series system.

Assumption 3. The system is in a non-failed state if all m components are in a non-failed state, otherwise the system is in a failed state.

A system in which at least k of the components must be in a non-failed state for the system to be in a non-failed system is denoted a k -out-of- m system, thus a series system is an m -out-of- m system.

The uncertain lifetime of the system is given by the following definition.

Theorem 2.1. The lifetime of the system is given by the random variable

$$S = \min(T_1, \dots, T_m) \tag{2.2}$$

with a sample space $(0, \infty)$, where T_1, \dots, T_m are the component lifetimes.

Proof. A series system fails whenever any component fails, therefore its lifetime is equal to the lifetime of the component with the minimum lifetime. \square

2.2 Probability distribution functions and parametric families

As random variables, it is not certain which values S, T_1, \dots, T_m will realize. However, their respective probability distribution functions quantify the uncertainty. The cumulative distribution function is given by the following definition.

Definition 2.4. *Let X be some random variable. The probability that $X \leq x$ is given by its cumulative distribution function, denoted by*

$$F_X(x) = P[X \leq x]. \quad (2.3)$$

The probability density function is given by the following definition.

Definition 2.5. *Let X be some continuous random variable. The relative likelihood that $X = x$ is given by its probability density function, denoted by*

$$f_X(x) = \frac{d}{dt} F_X(t). \quad (2.4)$$

By the Second Fundamental Theorem of Calculus, the probability that a continuous random variable $X \in (a, b]$, $a \leq b$, is given by

$$P[a < X \leq b] = F_X(b) - F_X(a) = \int_a^b f_X(s) ds. \quad (2.5)$$

By these definitions, the relative likelihood that the system will fail at time t is $f_S(t)$ and probability the j^{th} component will before time t is $F_{T_j}(t)$.

Notation. *We simplify the notation and let F_j denote F_{T_j} and f_j denote f_{T_j} . Column vectors are denoted by boldface lower case letters, e.g., \mathbf{x} , and row vectors are denoted by the transpose of column vectors, e.g., \mathbf{x}^T . Matrices are denoted by boldface upper case greek letters, e.g., \mathbf{A} . The j^{th} column of \mathbf{A} is denoted by \mathbf{A}_j (or $[\mathbf{A}]_j$) and the (j, k) -th element of \mathbf{A} is denoted by \mathbf{A}_{jk} (or $[\mathbf{A}]_{jk}$).*

We restrict our attention to parametrized families of probability distribution functions. The system lifetime has a probability density function that is a member of the following parametric family.

Definition 2.6. *The system lifetime has a probability density function that is a member of the parametric family of probability density functions denoted by*

$$\mathcal{F}_{\Theta} = \left\{ f_S(\cdot | \Theta) : \Theta \in \mathbb{R}^{q \times m} \right\}, \quad (2.6)$$

where $\Theta_j \in \Omega$ corresponds to the parameter index of the j^{th} component lifetime.

Particular values of Θ index particular probability density functions within the \mathcal{F}_{Θ} -family. The true probability density function of the system lifetime is given by the following definition.

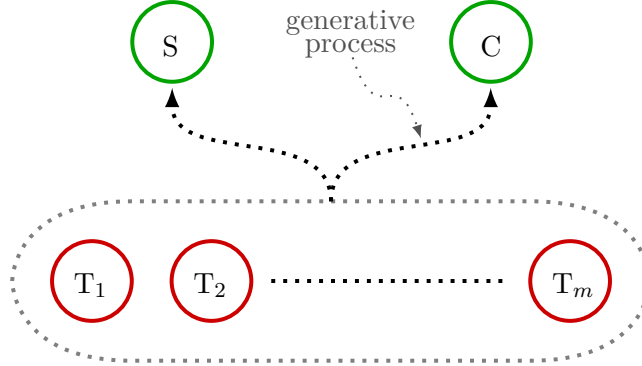
Definition 2.7. *The distribution of the system lifetime has a true parameter index denoted by Θ^* , and we say that*

$$S \sim f_S(\cdot | \Theta^*). \quad (2.7)$$

By these definitions, it follows that

$$T_j \sim f_j(\cdot | \Theta_j^*). \quad (2.8)$$

Figure 1: Graphical model of random variables where $S = \min\{T_1, T_2, \dots, T_m\}$ and masked component failure C is generated by some process that is a function of Θ^* .



2.3 α -masked component failure

By theorem 2.1, if a system failure occurs at time t , then one of the components failed at time t .¹

Definition 2.8. *The discrete random variable indicating the component responsible for a system failure is denoted by K with a sample space $\mathcal{K} = \{1, \dots, m\}$. That is, $K = j$ indicates that $T_j < T_k$ for $k \in \mathcal{K} \setminus \{j\}$.*

We may not be certain about which component caused a system failure, but we may be able to narrow the component cause to some *subset*. We denote these a *masked* component failure as given by the following definition.

Definition 2.9. *Let C denote the random masked component failure with a sample space $2^{\mathcal{K}}$.*

We suppose there is some process that, informed by some degree of knowledge about the components, generates a masked component failure which *plausibly*² contains the failed component as depicted in fig. 1.

A masked component failure represents the notion of a diagnostician who is able to, with varying degrees of success, isolate the failed component to some subset of the system. We consider a particular generative model given by the following definition.

Definition 2.10. *An α -masked component failure C contains the failed component with probability $\alpha \cdot 100\%$, where α is denoted its accuracy. Let A denote the random accuracy of a masked component failure with a sample space $\mathcal{A} \subseteq [0, 1]$. The distribution of a random masked component failure C given $A = \alpha$ is denoted a random α -masked component failure.*

The accuracy of a masked component failure could theoretically convey information about the system. As a simplification, we make the following assumption.

¹Only a single component can cause a system failure since two or more continuous random variables cannot realize the same value.

²The random masked component failures carry information about the true parameter index Θ^* by being dependent in some way on the component lifetimes.

Assumption 4. *The random accuracy A of masked component failures is independent of the random system lifetime S and the random component cause K and therefore does not carry information about the true parameter index Θ^* .*

A random masked component failure may have an uncertain cardinality. We model this uncertainty with the random variable given by the following definition.

Definition 2.11. *Let $W = |C|$ denote the random cardinality of a random masked component failure with a sample space $\{1, \dots, m-1\}$. W may never realize m since that does not isolate a proper subset of the components and it may never realize 0 since that does not estimate the component cause of failure.*

Note that if $W = m$ were to be realized, contrary to the above definition, then a consistent point estimator for the lifetime distribution S is possible but no consistent estimator of the component lifetimes T_1, \dots, T_m is possible. In ??, we provide a point estimator in such a situation that is inconsistent but minimizes the *mean squared error*.

The cardinality of a masked component failure could theoretically convey information about the system. As a simplification, we make the following assumption.

Assumption 5. *The random cardinality W of masked component failures is independent of the random system lifetime S and the random component cause K and therefore does not carry information about the true parameter index Θ^* .*

We make the following simplifying assumption.

Assumption 6. *The random masked component failure C conditioned K is independent of S .*

The α -masked component model is somewhat naive since C given K is conditionally independent of S . By assumption 6, an α -masked component failure carries information about Θ^* through its dependency on K .

The indicator function and the k -subset functions are needed for subsequent material and are given by the following definitions.

Definition 2.12 (Indicator function). *Given a subset \mathcal{A} of a set \mathcal{B} , the indicator function $\mathbb{1}_{\mathcal{A}}: \mathcal{B} \mapsto \{0, 1\}$ is equal to 1 if $x \in \mathcal{A}$ and 0 otherwise.*

Definition 2.13. *The subsets of \mathcal{A} that are of cardinality k is defined as*

$$[\mathcal{A}]^k := \{ \mathcal{X} \in \mathcal{A} \mid |\mathcal{X}| = k \}, \quad (2.9)$$

which has a cardinality

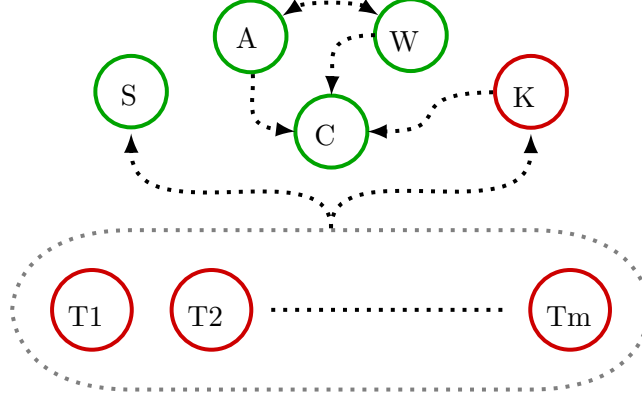
$$\binom{|\mathcal{A}|}{k}. \quad (2.10)$$

The distribution of candidate sets conditioned on a specific cardinality is given by the following definition.

Definition 2.14. *The random candidate set C conditioned on $W = w$ has a sample space given by*

$$\mathcal{W}_w = [2^K]^w. \quad (2.11)$$

Figure 2: Graphical model of random variables where $S = \min\{T_1, T_2, \dots, T_m\}$, C is a random α -masked candidate set of cardinality w , K is the random component failure, W is the cardinality of random candidate sets, and A is the random α .



By assumptions 4 and 5, the distributions of W and A is not dependent on the true parameter index Θ^* . Any inferences about the parametric model will be *independent* of W and A except in the sense that conditioning on larger α -masked component failures or less accurate accuracy α generates estimators with greater variance. This, of course, is quite important, but the fact that a particular W or A is realized says nothing by itself about Θ^* .

In general, a sample of masked system failures may have variable sized candidate sets. In ??, we derive the sampling distribution of the maximum likelihood estimator of this kind of sample using the inverse-variance weighted mean.

The *primary* data point is the *masked system failure* as given by the following definition.

Definition 2.15. An α -masked system failure consists of a system failure time t , a corresponding set of failed candidate components C , and an expected accuracy α of the C , denoted by $\langle t, C, \alpha \rangle$.

Assumption 7. A random α -masked system failure is a jointly distributed random system lifetime S , random failed candidate set C , and random α -accuracy A .

Assumption 8. A random sample of n masked system failures, denoted by

$$\mathbf{M}_n = (\langle S_1, C_1, A_1 \rangle, \dots, \langle S_n, C_n, A_n \rangle) , \quad (2.12)$$

consists of n independent and identically jointly distributed pairs of system failure times and component sets.

Assumption 9. The only information we have about the system and its true parameter index Θ^* is given by observing a random sample of n masked component failures,

$$(\langle S_1 = t_1, C_1 = C_1, A_1 = \alpha_1 \rangle, \dots, \langle S_n = t_n, C_n = C_n, A_n = \alpha_n \rangle) . \quad (2.13)$$

Given the assumed model, the object of statistical interest is the (unknown) true parameter index Θ^* . Provided an estimate of Θ^* , any characteristic that is a function of Θ^* can be estimated

as a result, thus the primary objective is to use the *information* in a sample of n masked system failures to estimate Θ^* with some quantifiable uncertainty.

3 Derivation of parametric probability distribution functions

Chapter 2 described the probabilistic model. In what follows, we derive the parametric functions that are entailed by the model.

3.1 System lifetime

The complement of the cumulative distribution function is the reliability function and is given by the following definition.

Definition 3.1. *Let X be some random variable. The probability that $X > x$ is given by its reliability function, denoted by*

$$R_X(x) = P[X > x] = 1 - F_X(x) . \quad (3.1)$$

The reliability function of the system lifetime is given by the following theorem.

Theorem 3.1. *The system lifetime has a reliability function given by*

$$R_S(t|\Theta^*) = \prod_{j=1}^m R_j(t|\Theta_j^*) , \quad (3.2)$$

where R_j is the reliability function of the j^{th} component.

Proof. By definition 3.1, the reliability function is given by

$$R_S(t|\Theta^*) = P[S > t] . \quad (a)$$

By eq. (2.2), $S = \min(T_1, \dots, T_m)$. Performing this substitution yields

$$R_S(t|\Theta^*) = P[\min(T_1, \dots, T_m) > t] . \quad (b)$$

For the minimum to be larger than t , every component must be larger than t , leading to the equivalent equation

$$R_S(t|\Theta^*) = P[T_1 > t, \dots, T_m > t] . \quad (c)$$

By assumption 2, the component lifetimes are independent, thus by the axioms of probability

$$R_S(t|\Theta^*) = P[T_1 > t] \cdots P[T_m > t] . \quad (d)$$

By ??, the probability that $T_j > t$ is equivalent to the reliability function of the j^{th} component lifetime evaluated at t . Performing these substitutions results in

$$R_S(t|\Theta^*) = \prod_{j=1}^m R_j(t|\Theta_j^*) . \quad (e)$$

□

The probability density function of the system's lifetime is given by the following theorem.

Theorem 3.2. *The lifetime of the m -out-of- m system has a probability density function given by*

$$f_S(t|\Theta^*) = \sum_{j=1}^m \left(f_j(t|\Theta_j^*) \prod_{\substack{k=1 \\ k \neq j}}^m R_k(t|\Theta_k^*) \right), \quad (3.3)$$

where R_k and f_k are respectively the reliability and probability density functions of the k^{th} component.

Proof. By ??,

$$f_S(t|\Theta^*) = -\frac{d}{dt} R_S(t|\Theta^*). \quad (a)$$

By theorem 3.1, this is equivalent to

$$f_S(t|\Theta^*) = -\frac{d}{dt} \prod_{j=1}^m R_j(t|\Theta_j^*). \quad (b)$$

By the product rule, this is equivalent to

$$\begin{aligned} f_S(t) &= \frac{d}{dt} R_1(t|\Theta_1^*) \prod_{j=2}^m R_j(t|\Theta_j^*) \\ &\quad - R_1(t|\Theta_1^*) \frac{d}{dt} \prod_{j=2}^m R_j(t|\Theta_j^*). \end{aligned} \quad (c)$$

By ??, we can substitute $-\frac{d}{dt} R_1(t|\Theta_1^*)$ with $f_1(t|\Theta_1^*)$, resulting in

$$\begin{aligned} f_S(t) &= f_1(t) \prod_{j=2}^m R_j(t|\Theta_j^*) \\ &\quad - R_1(t|\Theta_1^*) \frac{d}{dt} \prod_{j=2}^m R_j(t|\Theta_j^*). \end{aligned} \quad (d)$$

Applying the product rule again results in

$$\begin{aligned} f_S(t|\Theta^*) &= f_1(t|\Theta_1^*) \prod_{j=2}^m R_j(t|\Theta_j^*) \\ &\quad + f_2(t|\Theta_2^*) \prod_{\substack{j=1 \\ j \neq 2}}^m R_j(t|\Theta_j^*) \\ &\quad - R_1(t|\Theta_1^*) R_2(t|\Theta_2^*) \frac{d}{dt} \prod_{j=3}^m R_j(t|\Theta_j^*). \end{aligned} \quad (e)$$

We see a pattern emerge. Applying the product rule $m-1$ times results in

$$\begin{aligned} f_S(t|\Theta^*) &= \sum_{j=1}^{m-1} f_j(t|\Theta_j^*) \prod_{\substack{k=1 \\ k \neq j}}^m R_k(t|\Theta_k^*) \\ &\quad - \prod_{j=1}^{m-1} R_j(t|\Theta_j^*) \frac{d}{dt} R_m(t|\Theta_m^*) \\ &= \sum_{j=1}^m f_j(t|\Theta_j^*) \prod_{\substack{k=1 \\ k \neq j}}^m R_k(t|\Theta_k^*). \end{aligned} \quad (f)$$

□

The failure rate function simplifies some of the subsequent material and is given by the following definition.

Definition 3.2. *A random variable X with a probability density function f_X and reliability function R_X has a failure rate function given by*

$$h_X(t) = \frac{f_X(t)}{R_X(t)}. \quad (3.4)$$

The failure rate function of the system lifetime is given by the following theorem.

Theorem 3.3. *The system lifetime has a failure rate function given by*

$$h_S(t|\Theta^*) = \sum_{j=1}^m h_j(t|\Theta_j^*) \quad (3.5)$$

where h_j is the failure rate function of the j^{th} component lifetime.

Proof. By eq. (3.4), the failure rate function for the system is given by

$$h_S(t|\Theta^*) = \frac{f_S(t|\Theta^*)}{R_S(t|\Theta^*)}. \quad (a)$$

Plugging in these parametric functions results in

$$h_S(t|\Theta^*) = \frac{\sum_{j=1}^m \left(f_j(t|\Theta_j^*) \prod_{\substack{k=1 \\ k \neq j}}^m R_k(t|\Theta_k^*) \right)}{\prod_{j=1}^m R_j(t|\Theta_j^*)}, \quad (b)$$

which can be simplified to

$$h_S(t|\Theta^*) = \sum_{j=1}^m f_j(t|\Theta_j^*) / R_j(t|\Theta_j^*) \quad (c)$$

$$= \sum_{j=1}^m h_j(t|\Theta_j^*). \quad (d)$$

□

3.2 Component cause of system failure

According to definition 2.8, the component responsible for a system failure is given by the discrete random variable K . The joint likelihood that component k is the cause of a system failure occurring at time t is given by the following theorem.

Theorem 3.4. *The joint probability density function of random variable K and random system lifetime S is given by*

$$f_{K,S}(k, t|\Theta^*) = f_k(t|\Theta^*) \prod_{\substack{j=1 \\ j \neq k}}^m R_j(t|\Theta_j^*), \quad (3.6)$$

where f_k and R_k are respectively the probability density and reliability functions of the k^{th} component.

Proof. See appendix A for a detailed proof. However, a quick justification follows. The likelihood that component k is the cause of a system failure at time t is equivalent to the likelihood that component k fails at time t and the other components do not fail before time t . According to assumption 2, the random variables T_1, \dots, T_m are mutually independent, therefore the likelihood is the product of the probability density function $f_k(t|\Theta_k^*)$ and the reliability functions $R_j(t|\Theta_j^*)$ for $j \in \mathcal{K} \setminus \{k\}$. \square

By definition 2.8, if a system failure occurs at time t , then each component has a particular probability of being the cause as given by the following theorem.

Theorem 3.5. *The conditional probability mass function $p_{K|S}$ maps each component $k \in \mathcal{K}$ to the probability that component k is the cause of the system failure at time t and is given by*

$$p_{K|S}(k|t, \Theta^*) = \frac{h_k(t|\Theta_k^*)}{h_S(t|\Theta^*)}, \quad (3.7)$$

where h_j and h_S are respectively the failure rate functions of the j^{th} component and the system.

Proof. By the axioms of probability,

$$p_{K|S}(k|t, \Theta^*) = \frac{f_{K,S}(k, t|\Theta^*)}{f_S(t|\Theta^*)}, \quad (a)$$

thus we wish to show that

$$\frac{h_k(t|\Theta_k^*)}{h_S(t|\Theta^*)} = \frac{f_{K,S}(k, t|\Theta^*)}{f_S(t|\Theta^*)}. \quad (b)$$

By theorem 3.3, the system's failure rate function is given by

$$h_S(t|\Theta^*) = \sum_{j=1}^m h_j(t|\Theta_j^*). \quad (3.5 \text{ revisited})$$

Performing this substitution results in

$$\frac{h_k(t|\Theta_k^*)}{h_S(t|\Theta^*)} = \frac{h_k(t|\Theta_k^*)}{\sum_{j=1}^m h_j(t|\Theta_j^*)}. \quad (c)$$

By definition 3.2, the failure rate function of the p^{th} component is given by

$$h_p(t|\Theta_p^*) = \frac{f_p(t|\Theta_p^*)}{R_p(t|\Theta_p^*)}. \quad (d)$$

Performing this substitution results in

$$\frac{h_k(t|\Theta_k^*)}{h_S(t|\Theta^*)} = \frac{f_k(t|\Theta_k^*)}{R_k(t|\Theta_k^*)} \cdot \left[\underbrace{\frac{f_1(t|\Theta_1^*)}{R_1(t|\Theta_1^*)} + \dots + \frac{f_m(t|\Theta_m^*)}{R_m(t|\Theta_m^*)}}_A \right]^{-1}. \quad (e)$$

To combine the fractions labeled A in the above equation, we make their respective denominators the same, resulting in

$$\frac{h_k(t|\Theta_k^*)}{h_S(t|\Theta^*)} = \underbrace{\frac{f_k(t|\Theta_k^*)}{R_k(t|\Theta_k^*)}}_B \cdot \frac{\overbrace{\prod_{j=1}^m R_j(t|\Theta_j^*)}^C}{\sum_{j=1}^m \left[\prod_{\substack{p=1 \\ p \neq j}}^m R_p(t|\Theta_p^*) \right] f_j(t|\Theta_j^*)}. \quad (f)$$

Dividing the part of the above equation labeled C in the numerator by the part labeled B in the denominator results in

$$\frac{h_k(t|\Theta_k^*)}{h_S(t|\Theta^*)} = \frac{f_k(t|\Theta_k^*) \prod_{\substack{j=1 \\ j \neq k}}^m R_j(t|\Theta_j^*)}{\sum_{j=1}^m f_j(t|\Theta_j^*) \prod_{\substack{p=1 \\ p \neq j}}^m R_p(t|\Theta_p^*)}. \quad (g)$$

By eq. (3.3), the denominator in the above equation is the density $f_S(t|\Theta^*)$ and, by eq. (3.6), the numerator is the joint density $f_{K,S}(k, t|\Theta^*)$. Performing these substitutions results in

$$\frac{h_k(t|\Theta_k^*)}{h_S(t|\Theta^*)} = \frac{f_{K,S}(k, t|\Theta^*)}{f_S(t|\Theta^*)}. \quad (h)$$

□

Corollary 3.5.1. *The joint probability density function of random variable K and random system lifetime S is given by*

$$f_{K,S}(k, t|\Theta^*) = h_k(t|\Theta_k^*) R_S(t|\Theta^*). \quad (3.8)$$

The proof of eq. (3.8) immediately follows from eqs. (3.6) and (3.7).

3.3 α -masked component failures

The series system consists of m components with random lifetimes T_1, \dots, T_m . These random variables are *latent* (unobservable) and thus, for instance, we cannot know with certainty from a sample of masked system failures which component failed.

The α -candidate sets have a conditional probability distribution given by the following theorem.

Theorem 3.6. *The random α -candidate set C conditioned on W , K , and A has a probability distribution given by*

$$p_{C|K,W,A}(C|k, w, \alpha) = \frac{\alpha}{\binom{m-1}{w-1}} \mathbb{1}_C(k) + \frac{1-\alpha}{\binom{m-1}{w}} \mathbb{1}_{\bar{C}}(k) \quad (3.9)$$

with a sample space \mathcal{W}_w .

Proof. Suppose we condition on $K = k$, i.e., the k^{th} -component failed. By ??, a random α -candidate set C contains k with probability α and does *not* contain k with probability $1 - \alpha$.

We partition the sample space $\mathcal{W}_w := [\mathcal{K}]^w$ into the set of sets containing k and its complement. There are ${}^{m-1}C_{w-1}$ candidate sets that contain k since we may place k in the “first” position and of the remaining $w - 1$ positions uniformly choosing any of the remaining $m - 1$ components. There are ${}^{m-1}C_w$ candidate sets that do not contain k since, by removing k from the set of m choices, uniformly choosing w of the remaining $m - 1$ components.

By definition, the outcomes in the partition containing k occur with probability α and we assign uniform probability $(m^{-1}C_{w-1})^{-1}$ to any such outcome. By the product rule the joint probability that a random α -candidate set contains k is

$$\frac{\alpha}{\binom{m-1}{w-1}}. \quad (\text{a})$$

Similarly, the outcomes in the partition not containing k occur with probability $1 - \alpha$ and we assign uniform probability $(m^{-1}C_w)^{-1}$ to any such outcome. By the product rule the joint probability that a random α -candidate set does not contain k is

$$\frac{1 - \alpha}{\binom{m-1}{w}}. \quad (\text{b})$$

□

Remark. The α -candidate set model is somewhat naive since if $\mathbb{1}_{C_1}(k) = \mathbb{1}_{C_2}(k)$ then

$$p_{C|K,S,W,A}(\mathcal{C}_1|k, t_1, w, \alpha) = p_{C|K,S,W,A}(\mathcal{C}_2|k, t_2, w, \alpha). \quad (3.10)$$

That is, α -masked component failures that include the failed component are equally likely to be chosen and C given K is conditionally independent of the system failure time S . A more informative generative model of masked component failures may preferentially include components that are more likely to fail at the given system failure time. \triangle

The random variable K is *latent*. In a masked system failure sample, the only observable is the candidate set, the system failure time, and the expected accuracy α of the α -masked candidate set. The joint probability density function of C and S conditioned on W and A is given by the following theorem.

Theorem 3.7. The conditional joint probability density function of random variable C and random system lifetime S given $W = w$ is given by

$$f_{C,S|W,A}(\mathcal{C}, t|w, \alpha, \Theta^*) = \sum_{k=1}^m p_{C|K,W,A}(\mathcal{C}|k, w, \alpha) f_{K,S}(k, t|\Theta^*) \quad (3.11)$$

with a sample space $\mathcal{W}_w \times (0, \infty)$. Equation (3.11) may be rewritten as

$$f_{C,S|W,A}(\mathcal{C}, t|w, \alpha, \Theta^*) = \frac{R_S(t|\Theta^*)}{\binom{m-1}{w}} \left((1 - \alpha) h_S(t|\Theta^*) - \left(1 - \alpha \frac{m}{w}\right) \sum_{k \in \mathcal{C}} h_k(t|\Theta_k^*) \right). \quad (3.12)$$

Proof. By the laws of probability, the joint distribution of C , K and S conditioned on W and A may be written as

$$f_{C,K,S|W,A}(\mathcal{C}, k, t|w, \alpha, \Theta^*) = f_{C|K,S,W,A}(\mathcal{C}|k, t, w, \alpha, \Theta^*) f_{K,S|W,A}(k, t|w, \alpha, \Theta^*). \quad (\text{a})$$

We observe that C conditioned on K , W , and A is independent of S and the true parameter index Θ^* and the joint distribution of S and K is independent of W and A . Thus, we simplify the joint probability density function to

$$f_{C,K,S|W,A}(\mathcal{C}, k, t|w, \alpha, \Theta^*) = p_{C|K,W,A}(\mathcal{C}|k, w, \alpha) f_{K,S}(k, t|\Theta^*). \quad (\text{b})$$

The marginal joint distribution of C and S given W is calculated by summing the joint probability distribution over K which has a sample space $\{1, \dots, m\}$. \square

The *optimal* case is provided by the 1-masked candidate set as given by the following corollary.

Corollary 3.7.1. *The joint probability density of C and S conditioned on $W = w$ and $\alpha = 1$ is given by*

$$f_{C,S|W,A}(\mathcal{C}, t|w, \alpha = 1, \Theta^*) = \frac{R_S(t|\Theta^*)}{\binom{m-1}{w-1}} \sum_{k \in \mathcal{C}} h_k(t|\Theta^*) \quad (3.13)$$

with a sample space $\mathcal{W}_w \times (0, \infty)$.

Suppose the candidate sets are chosen *randomly*, which is equivalent to the $\frac{w}{m}$ -masked candidate set model.

Corollary 3.7.2. *The joint probability density of C and S conditioned on $W = w$ and $\alpha = \frac{w}{m}$ is given by*

$$f_{C,S|W,A}(\mathcal{C}, t|w, \alpha = w/m, \Theta^*) = \frac{f_S(t|\Theta^*)}{\binom{m}{w}} \quad (3.14)$$

with a sample space $\mathcal{W}_w \times (0, \infty)$.

Proof. Intuitively, since random candidate sets of w components (out of m) are randomly chosen, any of the $\binom{m}{w}$ subsets of size w are equally likely and thus the probability mass of each candidate set is just $\binom{m}{w}^{-1}$. Taking the product of the density for the system lifetime S and $\binom{m}{w}^{-1}$ yields the result. However, a more rigorous proof is given by substituting α with $\frac{w}{m}$ in eq. (3.12) and simplifying, yielding the intermediate result

$$f_{C,S|W,A}(\mathcal{C}, t|w, \alpha, \Theta^*) = \frac{R_S(t|\Theta^*)}{\binom{m-1}{w}} \frac{m-w}{m} h_S(t|\Theta^*) \quad (a)$$

$$= \frac{R_S(t|\Theta^*)}{\binom{m}{w}} h_S(t|\Theta^*). \quad (b)$$

By ??, $f_S(t|\Theta^*) = R_S(t|\Theta^*) h_S(t|\Theta^*)$. Making this substitution yields the result. \square

As intuition suggests, candidate sets that are randomly chosen carry no information about the true parameter index. Therefore, the set of informative α -masked candidate set models for estimating Θ^* is indexed by the interval $\frac{w}{m} < \alpha \leq 1$.

Remark. *Conversely, suppose we have an adversarial diagnostician who generates α -masked candidate sets where $\alpha < \frac{w}{m}$. Then, the random candidate sets serve as misinformation. If it is known or estimated that the candidate sets are misinforming, then we may estimate the answer to a related but different question, "What is Θ^* not likely to be?" While not nearly as informative, the only case in which no information about Θ^* is conveyed is when candidate sets are randomly selected, i.e., $\alpha = \frac{w}{m}$. \triangle*

Corollary 3.7.3. *The conditional probability mass function $p_{C|S,W,A}$ maps each candidate set \mathcal{C} to the probability that the given candidate set contains the component responsible for the system failure at time t . By the axioms of probability,*

$$p_{C|S,W,A}(\mathcal{C}|t, w, \alpha, \Theta^*) = \frac{(1 - \alpha) h_S(t|\Theta^*) - (1 - \alpha \frac{m}{w}) \sum_{j \in \mathcal{C}} h_j(t|\Theta_j^*)}{\binom{m-1}{w-1} h_S(t|\Theta^*)} \mathbb{1}_{\mathcal{W}_w}(\mathcal{C}), \quad (3.15)$$

where h_S and h_j are the failure rate functions of the system and the j^{th} component respectively.

Proof. Suppose it is given that $W = w$ and $A = \alpha$. By assumption 5 and ?? and the axioms of probability, we may rewrite the joint distribution of C and S as

$$f_{C,S|W,A}(\mathcal{C}, t|w, \alpha, \Theta^*) = p_{C|S,W,A}(\mathcal{C}|t, w, \alpha, \Theta^*) f_S(t|\Theta^*), \quad (a)$$

which may be rearranged to

$$p_{C|S,W,A}(\mathcal{C}|t, w, \alpha, \Theta^*) = \frac{f_{C,S|W,A}(\mathcal{C}, t|w, \alpha, \Theta^*)}{f_S(t|\Theta^*)}, \quad (b)$$

where $f_{C,S|W,A}$ is the joint probability density function given by corollary 3.7.1. Making this substitution and further simplifying yields the result. \square

By assumption 5, the random lifetime S and the cardinality of the random candidate set $W = |C|$ are independent, thus by the axioms of probability the joint density of C , S , and W is given by

$$f_{C,S,W|A}(\mathcal{C}, t, w|\alpha, \Theta^*) = p_{C|S,W,A}(\mathcal{C}|t, w, \alpha, \Theta^*) p_{W|A}(w|\alpha) f_S(t|\Theta^*). \quad (3.16)$$

Corollary 3.7.4. When $\alpha = 1$, $p_{C|S,W,A}$ simplifies to

$$p_{C|S,W,A}(\mathcal{C}|t, w, \alpha = 1, \Theta^*) = \frac{\mathbb{1}_{\mathcal{W}_w}(\mathcal{C}) \sum_{j \in \mathcal{C}} h_j(t|\Theta_j^*)}{\binom{m-1}{w-1} h_S(t|\Theta^*)}, \quad (3.17)$$

where h_S and h_j are the failure rate functions of the system and the j^{th} component respectively.

Corollary 3.7.5. The random component failure K given S is conditionally independent of C , W , and A , i.e.,

$$p_{K|C,S,W,A}(k|\mathcal{C}, t, w, \alpha, \Theta^*) = p_{K|S}(k|t, \Theta^*). \quad (3.18)$$

See appendix C for a proof.

4 Likelihood and Fisher information

Given a random variable X , the *information content* of observing a particular realization $X = x$ is defined as

$$I_X(x) := -\log p_X(x). \quad (4.1)$$

Intuitively, the more likely the outcome x , the less information content of observing x , e.g., if x is certain then observing x conveys no information.

The *expected value* of the information content is given by the following definition.

Definition 4.1. The entropy of discrete random X is a measure of uncertainty defined as

$$H(X) := E[I_X(X)] = E[-\log p_X(X)]. \quad (4.2)$$

The following example may provide some intuition about the basis of the entropy metric.

Example 1 Assuming the only information diagnosticians have about a series system is provided by the distribution of K , the lower-bound on average number of “yes” or “no” questions needed to determine the component cause of a failed system is given by $H(K)$ (where the logarithm is base-2). In the worst-case, K is uniformly distributed and $H(K) = \log_2 m$.

The expectation $E[-\log f_{X(x)}]$ of a *continuous* random variable X is called *differential* entropy. Since *negative* values are possible in this case, the previous example no longer applies but the metric is still useful.

Suppose we are interested in the entropy of X given $Y = y$,

$$H(X|Y = y) := E[-\log f_{X|Y}(X|y)]. \quad (4.3)$$

More information is never expected to increase the uncertainty, i.e., $H(X|Y = y) \leq H(X)$, and if X and Y are statistically dependent then $H(X|Y = y) < H(X)$.

If the diagnosticians described in example 1 are provided with the additional information that the system failed at time t , then the lower-bound on the average number of questions required to isolate the component cause is given by $H(K|S = t) \leq H(K)$.

If we *average* over all of the values Y may realize, we get the conditional entropy

$$H(X|Y) := E[-\log f_{X|Y}(X|Y)]. \quad (4.4)$$

The joint entropy of X and Y is defined as

$$H(X, Y) := E[-\log f_{X,Y}(X, Y)] \quad (4.5)$$

which may also be computed by

$$H(X, Y) = H(X|Y) + H(Y). \quad (4.6)$$

If X and Y are *statistically independent*, then $H(X|Y) = H(X)$ and $H(Y|X) = H(Y)$.

When we make an *observation* of a random variable, the entropy is a measure of the *expected information content* provided by the observation. The expected information provided by a sample of n masked system failures is given by the following theorem.

Theorem 4.1. *Since S is statistically independent of W and A ,*

$$H(C, S|W = w, A = \alpha) = H(C|S, W = w, A = \alpha) + H(S). \quad (4.7)$$

Given some distribution S , the *maximum* entropy is when the failed component is in the candidate set by random chance, $A = \frac{w}{m}$, as given by

$$H(C, S|W = w, A = w/m) = H(S) + \log_2 \binom{m}{w} \quad (4.8)$$

and the *minimum* entropy is when the failed component is *certainly* in the candidate set, $A = 1$, as given by

$$H(C, S|W = w, A = 1) = H(S) + H(C|S, W = w, A = 1) \quad (4.9)$$

Proof. The ... □

The rate of change of the entropy with respect to α is given by the following theorem.

Theorem 4.2. *test*

In what follows, we frequently prefer to work with vectors rather than matrices.

Definition 4.2. Let $\text{vec}: \mathbb{R}^{m \times q} \mapsto \mathbb{R}^{m \cdot q}$ denote the linear transformation of a matrix of m rows and q columns into a column vector of $m \cdot q$ elements, where consecutive matrix rows are placed side by side. Additionally, the j^{th} component of a vector \mathbf{a} is denoted by $[\mathbf{a}]_j$ or a_j .

Notation. By definition 4.2, Θ^* has an equivalent representation denoted by $\theta^* \equiv \text{vec}(\Theta^*)$ and the j^{th} component of θ^* is denoted by θ_j^* . We use matrix and vector representations interchangeably.

The likelihood of observing a particular realization $\mathbf{M}_n = \langle \mathcal{C}_1, t_1 \rangle, \dots, \langle \mathcal{C}_n, t_n \rangle$ conditioned on $W = w$ and $A = \alpha$ is given by the following definition.

Definition 4.3. A random sample of n masked system failure times conditioned on candidate sets of cardinality w has a joint probability density given by

$$f_{\mathbf{M}_n|W,A}(\langle \mathcal{C}_1, t_1 \rangle, \dots, \langle \mathcal{C}_n, t_n \rangle \mid w, \alpha, \Theta^*) = \prod_{i=1}^n f_S(t_i | \Theta^*) \cdot f_{C|S,W,A}(C_i \mid t_i, w, \alpha, \Theta^*), \quad (4.10)$$

where f_S is the marginal density given by theorem 3.2 and $f_{C|S,W,A}$ is the conditional probability mass given by corollary 3.7.4.

A random variable is a measurable function from a probability space to a measurable space of values that the variable can take on. Such a space of values is called a random variable.

The subsequent material makes a few assumptions, denoted *regularity conditions*, given by the following.

Assumption 10. The following regularity conditions hold for the \mathcal{F}_Θ -family of the series system:

1. θ^* is interior to the parameter support Ω .
2. The log-likelihood ℓ is continuous, thrice differentiable, and bounded.
3. The true parameter value θ^* is identified, i.e.,

$$\theta^* = \arg \max_{\theta \in \Omega} E_{\theta^*} [\ln f_{C,S|W,A}(C, S \mid w, \alpha, \theta)]. \quad (4.11)$$

By definition 4.3, the probability density that $\mathbf{M}_n = \langle \mathcal{C}_1, t_1, \alpha \rangle, \dots, \langle \mathcal{C}_n, t_n, \alpha \rangle$ conditioned on $W = w$ and $A = \alpha$ is given by

$$f_{\mathbf{M}_n|W,A}(\langle \mathcal{C}_1, t_1 \rangle, \dots, \langle \mathcal{C}_n, t_n \rangle \mid w, \alpha, \theta^*) = \prod_{i=1}^n f_{C,S|W,A}(C_i, t_i \mid w, \alpha, \theta^*). \quad (4.10 \text{ revisited})$$

If we fix the sample of masked system failures and allow the parameter θ to change, we have the likelihood function.

Definition 4.4. The likelihood function represents the likelihood of parameter index θ with \mathbf{M}_n fixed at $\langle \mathcal{C}_1, t_1 \rangle, \dots, \langle \mathcal{C}_n, t_n \rangle$ is given by

$$\mathcal{L}_n(\theta \mid w, \alpha) := \prod_{i=1}^n f_{C,S|W,A}(C_i, t_i \mid w, \alpha, \theta). \quad (4.12)$$

The likelihood is a measure of the extent to which the given sample provides support for particular values of a parameter in \mathcal{F}_Θ ; its surface captures all the information there is about θ^* from the given sample.

The log-likelihood plays a more central role than the likelihood for reasons that will be made clear later on.

Definition 4.5. *The log-likelihood with respect to θ given a particular $\mathbf{M}_n = \langle \mathcal{C}_1, t_1, \alpha \rangle, \dots, \langle \mathcal{C}_n, t_n, \alpha \rangle$ conditioned on $W = w$ and $A = \alpha$ is given by*

$$\ell_n(\theta|w, \alpha) = \sum_{i=1}^n \ln f_{C,S|W,A}(\mathcal{C}_i, t_i|w, \alpha, \theta). \quad (4.13)$$

In subsequent material, we need the gradient operator given by the following definition.

Definition 4.6 (Gradient operator). *The gradient operator $\nabla: (\mathbb{R}^q \mapsto \mathbb{R}) \mapsto (\mathbb{R}^q \mapsto \mathbb{R}^q)$ is defined as*

$$\nabla g := \left(\frac{\partial g}{\partial x_1}, \dots, \frac{\partial g}{\partial x_q} \right)^\top. \quad (4.14)$$

The score function is given by the following definition.

Definition 4.7. *The Fisher's score function, denoted by $s_n: \mathbb{R}^{m \cdot q} \mapsto \mathbb{R}^{m \cdot q}$, is the gradient of the log-likelihood function given a sample of n masked failure times,*

$$s_n(\theta|w, \alpha) := \nabla \ell_n(\theta|w, \alpha). \quad (4.15)$$

The variance-covariance is given by the following definition.

Definition 4.8. *A random vector $\mathbf{X} \in \mathbb{R}^p$ has a variance-covariance given by*

$$\text{var}[\mathbf{X}] := \mathbb{E}[(\mathbf{X} - \mathbb{E}[\mathbf{X}]) (\mathbf{X} - \mathbb{E}[\mathbf{X}])^\top], \quad (4.16)$$

which is a p -by- p semi-positive definite matrix.

The score function given a realization $\mathbf{M}_n = \langle \mathcal{C}_1, t_1, \alpha \rangle, \dots, \langle \mathcal{C}_n, t_n, \alpha \rangle$ given $W = w$ and $A = \alpha$ is a statistic. When we consider it as a function of a random sample, it is a random vector. When evaluated at the true parameter index θ^* , it has a variance-covariance given by the following definition.

Definition 4.9 (Fisher information matrix). *The score as a function of a random sample \mathbf{M}_n given $W = w$ and $A = \alpha$ has a variance-covariance known as the Fisher information matrix, a $(m \cdot q)$ -by- $(m \cdot q)$ matrix denoted by*

$$\mathcal{I}_n(\theta^*|w, \alpha) := \text{var}_{\theta^*}[s_n(\theta^* | w, \alpha)]. \quad (4.17)$$

The Hessian is given by the following definition.

Definition 4.10. *The Hessian operator $\mathcal{H}: (\mathbb{R}^q \mapsto \mathbb{R}) \mapsto (\mathbb{R}^q \mapsto \mathbb{R}^{q \times q})$ applied to $g: \mathbb{R}^q \mapsto \mathbb{R}$ is a q -by- q matrix where the (j, k) -th element is defined as $\frac{\partial^2 g}{\partial x_j \partial x_k}$.*

This is somewhat related to the *entropy*, except instead of the logarithm of the density we take the logarithm of the gradient of the density, i.e.,

$$\mathcal{I}(\boldsymbol{\theta}^*|w, \alpha) = \text{var}_{\boldsymbol{\theta}^*} \left[\left(\nabla_{\boldsymbol{\theta}} \ln f_{C,S|W,A}(C, S | w, \alpha, \boldsymbol{\theta}) \right) \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}^*} \right]. \quad (4.18)$$

Under the regularity conditions, the information matrix may be computed from the Hessian and is given by the following definition. For clarity, the (i, j) -th element of the information matrix is given by the following definition.

Definition 4.11. *The information matrix of $\boldsymbol{\theta}^*$ given by the joint distribution of C and S conditioned on $W = w$ and $A = \alpha$ is given by the expectation of the negative of the Hessian of the log-likelihood function ℓ evaluated at $\boldsymbol{\theta}^*$. The (i, j) -th element of \mathcal{I} is given by*

$$\mathcal{I}(\boldsymbol{\theta}^*|w, \alpha)_{ij} = - \sum_{C \in \mathcal{W}_w} \int_0^\infty \left[\frac{\partial^2}{\partial \theta_i \partial \theta_j} \ln f_{C,S|W,A}(C, t|w, \alpha, \boldsymbol{\theta}) \right] \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}^*} \times f_{C,S|W,A}(C, t|w, \alpha, \boldsymbol{\theta}^*) dt, \quad (4.19)$$

where $\mathcal{W}_w \times (0, \infty)$ is the sample space.

Observe that

?? may be difficult to compute.

Definition 4.12. *The information matrix is given by the expectation of the negative of the Hessian of the log-likelihood function ℓ evaluated at $\boldsymbol{\theta}^*$. That is,*

$$\mathcal{I}_n(\boldsymbol{\theta}^*|w) = - \mathbb{E}_{\boldsymbol{\theta}^*} \left[\mathcal{H}(\ell_n(\boldsymbol{\theta}|w)) \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}^*} \right], \quad (4.20)$$

where the expectation is taken with respect to $\boldsymbol{\theta}^*$ and the Hessian is taken with respect to $\boldsymbol{\theta}$ and then evaluated at $\boldsymbol{\theta}^*$.

Notation. We denote the Fisher information matrix of $\boldsymbol{\theta}^*$ given $\mathbf{M}(w, 1)$ by

$$\mathcal{I}(\boldsymbol{\theta}^*|w) := \mathcal{I}_1(\boldsymbol{\theta}^*|w). \quad (4.21)$$

The information matrix is a key quantity in large sample theory; it quantifies the *expected* information (as a reduction of uncertainty) $\mathbf{M}(w, n)$ carries about $\boldsymbol{\theta}^*$. For sufficiently large samples, the log-likelihood has a Taylor series approximation given by

$$\ell_n(\boldsymbol{\theta}|w) \approx -\frac{1}{2} (\boldsymbol{\theta} - \boldsymbol{\theta}^*)^\top \mathcal{I}_n(\boldsymbol{\theta}^*|w) (\boldsymbol{\theta} - \boldsymbol{\theta}^*), \quad (4.22)$$

which is maximized at $\boldsymbol{\theta}^*$. Intuitively, the more “peaked” the log-likelihood is expected to be at $\boldsymbol{\theta}^*$, the more information a random sample is expected to carry about $\boldsymbol{\theta}^*$.

The information independent samples have about $\boldsymbol{\theta}^*$ is *additive* as given by the following postulate.

Postulate 4.1. *Algebraically, the Fisher information independent samples have about the true parameter index $\boldsymbol{\theta}^*$ is additive such that*

$$\mathcal{I}_{n_1+n_2}(\boldsymbol{\theta}^*|\cdot) = \mathcal{I}_{n_1}(\boldsymbol{\theta}^*|\cdot) + \mathcal{I}_{n_2}(\boldsymbol{\theta}^*|\cdot). \quad (4.23)$$

By postulate 4.1, the Fisher information of a random sample of n masked system failure times drawn from \mathbf{M}_n given $W = w$ and $A = \alpha$ is given by

$$\mathcal{I}_n(\boldsymbol{\theta}^*|w, \alpha) = n \mathcal{I}(\boldsymbol{\theta}^*|w, \alpha). \quad (4.24)$$

A *masked system failure time* in which the cardinality of the candidate set is distributed according to W , as described by ??, has a Fisher information matrix given by the following theorem.

Theorem 4.3. *The Fisher information matrix of true parameter index $\boldsymbol{\theta}^*$ with respect to a realization of jointly distributed random system lifetime S and candidate set C is given by*

$$\mathcal{I}(\boldsymbol{\theta}^*) = \sum_{\langle w, \alpha \rangle \in \mathcal{W} \times \mathcal{A}} p_{W,A}(w, \alpha) \mathcal{I}(\boldsymbol{\theta}^*|w, \alpha), \quad (4.25)$$

where $p_{W,A}$ is the joint distribution of W and A with a support $\mathcal{W} \times \mathcal{A}$.

Proof. Suppose we have a random sample of n masked system failure times. By ??, the expected number of sample points that have w α -masked candidate sets is given by

$$q := n \cdot p_{W,A}(w, \alpha). \quad (a)$$

Thus, by the additive property of *Fisher information* given by postulate 4.1, the information matrix of $\boldsymbol{\theta}^*$ conditioned on q such observations is given by

$$\mathcal{I}_q(\boldsymbol{\theta}^*|w, \alpha) = n \cdot p_{W,A}(w, \alpha) \mathcal{I}(\boldsymbol{\theta}^*|w, \alpha). \quad (b)$$

When we sum over the entire support of W and A , we arrive at the *information matrix* of $\boldsymbol{\theta}^*$ conditioned on n masked candidate sets as given by

$$\mathcal{I}_n(\boldsymbol{\theta}^*) = \sum_{\langle w, \alpha \rangle \in \mathcal{W} \times \mathcal{A}} n \cdot p_{W,A}(w, \alpha) \mathcal{I}(\boldsymbol{\theta}^*|w, \alpha). \quad (c)$$

When we set n to 1, we get the result

$$\mathcal{I}(\boldsymbol{\theta}^*) = \sum_{\langle w, \alpha \rangle \in \mathcal{W} \times \mathcal{A}} p_{W,A}(w, \alpha) \mathcal{I}(\boldsymbol{\theta}^*|w, \alpha). \quad (d)$$

□

5 Sampling distribution of parameter estimators

By definitions 2.6 and 2.7, the random system lifetime S has a probability density function that is a member of the \mathcal{F}_Θ -family with a true parameter index $\boldsymbol{\theta}^*$. Let some estimator of $\boldsymbol{\theta}^*$ be a function of the information in a random sample of n masked system failure times. The estimator is a function of a random sample, thus it has *sampling distribution*, a random vector denoted by \mathbf{Y}_n . That is,

$$\mathbf{Y}_n = \psi(\mathbf{M}_n). \quad (5.1)$$

A particular realization of the estimator is denoted by $\bar{\boldsymbol{\theta}}_n$, i.e., $\mathbf{Y}_n = \bar{\boldsymbol{\theta}}_n$. All else being equal, we prefer estimators which vary only *slightly* from sample to sample with a *central tendency* around $\boldsymbol{\theta}^*$. That is, we prefer unbiased estimators in which each component has small variance.

The sampling distribution of $\bar{\boldsymbol{\theta}}_n$ has a variance-covariance given by $\text{Var}(\mathbf{Y}_n)$ (see eq. (4.16)) and a bias given by the following definition.

Definition 5.1. The bias of a point estimator \mathbf{Y}_n is given by

$$\text{bias}(\mathbf{Y}_n) = \mathbb{E}[\mathbf{Y}_n] - \boldsymbol{\theta}^*. \quad (5.2)$$

We are often faced with a trade-off between bias and variance. A measure of estimator accuracy that is a function of both the bias and the variance is the mean squared error as given by the following definition.

Definition 5.2. The mean squared error (MSE) of the sampling distribution of $\bar{\boldsymbol{\theta}}_n$ is given by

$$\text{MSE}(\mathbf{Y}_n) = \mathbb{E}[(\mathbf{Y}_n - \boldsymbol{\theta}^*)^\top (\mathbf{Y}_n - \boldsymbol{\theta}^*)] \quad (5.3)$$

An equivalent way to compute the mean squared error is given by the following postulate.

Postulate 5.1. The mean squared error of the sampling distribution of $\bar{\boldsymbol{\theta}}_n$ as given by definition 5.2 is equivalent to

$$\text{MSE}(\mathbf{Y}_n) = \text{tr}(\text{Var}(\mathbf{Y}_n)) + \text{bias}^2(\mathbf{Y}_n), \quad (5.4)$$

where $\text{tr}(\mathbf{A})$ computes the sum of the diagonal elements of square matrix \mathbf{A} .

Fisher information reduces the uncertainty about $\boldsymbol{\theta}^*$. The *minimum-variance unbiased estimator* (UMVUE) has a lower-bound given by the inverse of the *Fisher information matrix*, denoted the *Cramér-Rao lower-bound*. The minimum variance obtainable from a random sample of n masked system failure times drawn from \mathbf{M}_n is given by

$$\mathbf{CRLB}_n = \frac{1}{n} \mathcal{I}^{-1}(\boldsymbol{\theta}^*). \quad (5.5)$$

By eq. (5.4), if $\bar{\boldsymbol{\theta}}_n$ is an unbiased estimator of $\boldsymbol{\theta}^*$, then the mean squared error is given by

$$\text{MSE}(\mathbf{Y}_n) = \text{tr}(\text{Var}(\mathbf{Y}_n)). \quad (5.6)$$

By eqs. (5.5) and (5.6), the mean squared error of any unbiased point estimator of $\boldsymbol{\theta}^*$ in which the only *a priori* information is given by a random sample \mathbf{M}_n has a lower-bound given by the trace of \mathbf{CRLB}_n ,

$$\text{MSE}(\mathbf{Y}_n) \geq \text{tr}(\mathbf{CRLB}_n). \quad (5.7)$$

5.1 Maximum likelihood estimator for α -masked system failure times each consisting of w components

Suppose the only information about the true parameter index $\boldsymbol{\theta}^*$ is given by a sample of n α -masked system failure times in which only those system failures with w components with accuracy α are observed. The maximum likelihood estimator based on this conditional sample is given in the following definition.

Definition 5.3. A parameter index $\hat{\boldsymbol{\theta}}(n|w)$ that maximizes the likelihood of observing some realization of \mathbf{M}_n given $W = w$ and A (see assumption 8) is a maximum likelihood estimator of $\boldsymbol{\theta}^*$. That is,

$$\hat{\boldsymbol{\theta}}(n|w) = \arg \max_{\boldsymbol{\theta} \in \boldsymbol{\Omega}} \mathcal{L}_n(\boldsymbol{\theta}|w, \alpha), \quad (5.8)$$

where $\boldsymbol{\Omega}$ is the feasible parameter space of the parametric family.

As further justification of assumption 7, the maximum likelihood estimator maximizes the likelihood of generating the given \mathbf{m}_n such that for each system failure time a component in the corresponding candidate set is the cause.

The accuracy of the maximum likelihood estimator is a function of the agreement between the model and reality. In our model, the primary point of disagreement is given by the assumption of a series system under a parametric model of independent component lifetimes, i.e., \mathcal{F}_Θ . However, if the model is a reasonable abstraction of the objects of interest,

The distribution of candidate sets, i.e, C, disagrees with the way the actual candidate sets are generated. Even if this is the case, the logic that components that are more likely to have caused a system failure *with respect to the parameters* (other ways in which a component might fail are outside the scope of the parametric model)

The accuracy of the maximum likelihood estimator is a function of the agreement between the model and reality. In our model, the primary point of disagreement is given by the assumption of a series system under a parametric model of independent component lifetimes, i.e., \mathcal{F}_Θ .

Another point of disagreement is the assumption that candidates are generated from the joint probability densities of the component lifetimes T_1, \dots, T_m and random candidate sets C. However, recall the assumption is that, under the assumed model $S = \min(T_1, \dots, T_m)$, components that are more likely to have failed at observable time $S = t$ are more likely to be in the observable candidate set.

Note that *any other* distribution of candidate sets generates a smaller likelihood given the parametric model. We claim this is the *optimal* candidate set since it provides the most information about θ . If the parametric family \mathcal{F}_Θ is a very accurate model of the reality, then whether the failed component is in the candidate set or not is irrelevant. What matters more is that the most likely components are in the candidate set, as that will provide the most information about θ .

If this is not a very accurate model, the maximum likelihood estimator may not be a consistent estimator even if . Rather, the maximum likelihood estimator will generate an estimate that makes the provided sample the most likely to be generated by the assumed model, i.e., the component lifetimes will not be accurately modeled since they must be modeled in such a way as to make the biased sample more likely.

The logarithm is a monotonically increasing function, thus the parameter value $\hat{\theta}(n|w)$ that maximizes \mathcal{L} also maximizes the log-likelihood ℓ (see definition 4.5), i.e.,

$$\hat{\theta}(n|w) = \arg \max_{\theta \in \Omega} \ell(\theta|w, \mathbf{m}_n). \quad (5.9)$$

According to Bickel [1], if the parameter support Ω of the parametric family is open, the log-likelihood ℓ is differentiable with respect to θ , and $\hat{\theta}(n|w)$ exists, then $\hat{\theta}(n|w)$ must be a stationary point of ℓ as given by

$$s(\hat{\theta}(n|w)|w, \mathbf{m}_n) = \mathbf{0}, \quad (5.10)$$

where s is the score function given by eq. (4.15). In cases where there are no closed-form solutions to eq. (5.10), iterative methods may be used as described in appendix D.

An estimator of θ^* given a sample $\mathbf{M}(w, n) = \mathbf{m}_n$ is the maximum likelihood estimator $\hat{\theta}(n|w)$ described in definition 5.3. Since $\hat{\theta}(n|w)$ is a function of the random sample $\mathbf{M}(w, n)$, it has a *sampling distribution*.

Definition 5.4. The sampling distribution of $\hat{\theta}(n|w)$ is a random vector denoted by $\mathbf{Y}(n|w) \in \mathbb{R}^{m \cdot q}$.

The generative model for the maximum likelihood estimator, `generate_mle`($w|\cdot$), is described by algorithm 1. This generative model depends on `generate_msft`, the generative model for the *masked system failure time* as described by ???. Note that in algorithm 1, line 7 may be approximated with `find_mle`, a function that numerically solves the stationary points of the maximum likelihood equations as described by algorithm 3.

Algorithm 1: Generative model of the maximum likelihood estimator conditioned on w candidates

Result: a realization of a maximum likelihood estimate from the sampling distribution of $\hat{\theta}(n|w)$.

Input:

- θ^* , the true parameter index.
- w , the cardinality of the candidate set.
- n , the number of masked system failure times.

Output:

$\hat{\theta}(n|w)$, a realization of $\mathbf{Y}(n|w)$.

```

1 Model generate_mle( $n, w, \theta^*$ )
2    $\mathbf{m}_n \leftarrow \emptyset$ 
3   for  $i \leftarrow 1$  to  $n$  do
4      $\mathbf{M} \leftarrow \text{generate\_msft}(w, \theta^*)$ 
5      $\mathbf{m}_n \leftarrow \mathbf{m}_n \cup \{\mathbf{M}\}$ 
6   end
7    $\hat{\theta}(n|w) \leftarrow \arg \max_{\theta \in \Omega} \ell(\theta|w, \mathbf{m}_n)$ 
8   return  $\hat{\theta}(n|w)$ 

```

The asymptotic sampling distribution of $\hat{\theta}(n|w)$ is a consistent estimator of θ^* .

Postulate 5.2. *As $n \rightarrow \infty$, the sampling distribution of $\hat{\theta}(n|w)$ converges in probability to θ^* , written*

$$\mathbf{Y}(n|w) \xrightarrow{P} \theta^*, \quad (5.11)$$

since

$$\lim_{n \rightarrow \infty} \mathbb{P}[\text{MSE}(\mathbf{Y}(n|w)) < \epsilon] = 1 \quad (5.12)$$

for every $\epsilon > 0$.

By eqs. (4.16) and (5.11), the variance-covariance of the asymptotic sampling distribution of $\hat{\theta}(n|w)$ is given by

$$\text{Var}[\mathbf{Y}(n|w)] = \mathbb{E}[(\mathbf{Y}(n|w) - \theta^*)(\mathbf{Y}(n|w) - \theta^*)^\top]. \quad (5.13)$$

Postulate 5.3. *The variance-covariance of the asymptotic sampling distribution of $\hat{\theta}(n|w)$ obtains the Cramér-Rao lower-bound for point estimators that are strictly a function of n masked system failure times in which candidate sets are of cardinality w . We denote this variance-covariance by*

$$\Sigma_n(\theta^*|w, \alpha) \equiv \frac{1}{n} \mathcal{I}^{-1}(\theta^*|w, \alpha). \quad (5.14)$$

The asymptotic sampling distribution of $\hat{\theta}(n|w)$ is normally distributed.

Postulate 5.4. As $n \rightarrow \infty$, the sampling distribution of $\hat{\boldsymbol{\theta}}(n|w)$ converges in distribution to a multivariate normal distribution with a mean $\boldsymbol{\theta}^*$ and a variance-covariance $\Sigma_n(\boldsymbol{\theta}^*|w, \alpha)$, written

$$\mathbf{Y}(n|w) \xrightarrow{d} \text{MVN}(\boldsymbol{\theta}^*, \Sigma_n(\boldsymbol{\theta}^*|w, \alpha)) . \quad (5.15)$$

The maximum likelihood estimator $\hat{\boldsymbol{\theta}}(n|w)$ is an asymptotically *efficient* estimator since it obtains the Cramér-Rao lower-bound as given by ???. Thus, for a sufficiently large sample of masked system failure times, the sampling distribution of $\hat{\boldsymbol{\theta}}(n|w)$ varies only *slightly* from sample to sample with a *central tendency* around $\boldsymbol{\theta}^*$.

By eqs. (5.5), (5.6) and (5.11) and ??, the asymptotic sampling distribution of $\hat{\boldsymbol{\theta}}(n|w)$ has the minimum mean squared error of any unbiased estimator,

$$\text{MSE}(\mathbf{Y}(n|w)) = \frac{1}{n} \text{tr} \left(\mathcal{I}^{-1}(\boldsymbol{\theta}^*|w, \alpha) \right) = \text{tr}(\text{CRLB}_{n,w}) . \quad (5.16)$$

Remark. Look up Slutsky's theorem to justify the next paragraph. △

To estimate the sampling distribution of $\hat{\boldsymbol{\theta}}(n|w)$, we may assume the sample size is sufficiently large such that the asymptotic distribution becomes a reasonable approximation. Since $\mathbf{Y}(n|w) \xrightarrow{P} \boldsymbol{\theta}^*$ and $\mathbf{Y}(n|w) \xrightarrow{d} \text{MVN}(\boldsymbol{\theta}^*, \Sigma_n(\boldsymbol{\theta}^*|w, \alpha))$, it follows that

$$\mathbf{Y}(n|w) \xrightarrow{d} \text{MVN}(\hat{\boldsymbol{\theta}}(n|w), \Sigma_n(\hat{\boldsymbol{\theta}}(n|w)|w, \alpha)) . \quad (5.17)$$

TODO: talk about $\mathbf{A} = \alpha$. We don't really want to think about modeling it, only saying that it is given in a sample. Say, for instance, the masked system failure time also has a label for how accurate they determine the candidate set to be, i.e., α -candidate set model where α can vary over the sample. Group by w and α , compute the information matrix, then combine them as described in this section to get the final estimator.

Thus, we can approximate $\boldsymbol{\theta}^*$ and $\mathcal{I}(\boldsymbol{\theta}^*|w, \alpha)$ and obtain the following result. The proof of this theorem is beyond the scope of this paper.

Theorem 5.1. For sufficiently large sample size n , the sampling distribution of $\hat{\boldsymbol{\theta}}(n|w)$ is approximately normally distributed with a mean $\hat{\boldsymbol{\theta}}(n|w)$ and a variance-covariance matrix $\Sigma_n(\hat{\boldsymbol{\theta}}(n|w)|w, \alpha)$, written

$$\mathbf{Y}(n|w) \overset{\text{approx.}}{\sim} \text{MVN}(\hat{\boldsymbol{\theta}}(n|w), \Sigma_n(\hat{\boldsymbol{\theta}}(n|w)|w, \alpha)) . \quad (5.18)$$

5.2 Maximum likelihood estimator

Consider an i.i.d. sample of r asymptotically unbiased estimates of $\boldsymbol{\theta}^*$ denoted by $\bar{\boldsymbol{\theta}}^{(i)}$ which have sampling distributions with variance-covariances given respectively by $\Sigma^{(i)}$ for $i = 1, \dots, r$. The maximum likelihood estimator of $\boldsymbol{\theta}^*$ given these point estimates is the inverse-variance weighted mean and is given by

$$\hat{\boldsymbol{\theta}} = \left(\sum_{i=1}^r \mathbf{A}_i \right)^{-1} \left(\sum_{i=1}^r \mathbf{A}_i \bar{\boldsymbol{\theta}}^{(i)} \right) , \quad (5.19)$$

where \mathbf{A}_i is the inverse of $\Sigma^{(i)}$.

Suppose that the estimates are given by the maximum likelihood estimator described in ??. The maximum likelihood estimator given these estimates has a sampling distribution given by the following definition.

Definition 5.5. Let \mathbf{M}_n be a random sample of n masked system failure times in which n_i realizations have w_i α_i -masked component failures for $i = 1, \dots, r$. The maximum likelihood estimator given by

$$\hat{\boldsymbol{\theta}} = \left(\sum_{i=1}^r \mathbf{A}_i \right)^{-1} \left(\sum_{i=1}^r \mathbf{A}_i \hat{\boldsymbol{\theta}}^{(i)} \right), \quad (5.20)$$

has a sampling distribution given by

$$\mathbf{Y}_n = \left(\sum_{i=1}^r \mathbf{A}_i \right)^{-1} \left(\sum_{i=1}^r \mathbf{A}_i \mathbf{Y}^{(i)}(n_i | w_i, \alpha_i) \right), \quad (5.21)$$

where $\mathbf{Y}^{(i)}(n_i | w_i, \alpha_i)$ is the sampling distribution of $\hat{\boldsymbol{\theta}}(n_i | w_i, \alpha_i)$ for $i = 1, \dots, r$ and

$$\mathbf{A}_i = n_i \mathcal{I}(\boldsymbol{\theta}^* | w_i, \alpha_i) \quad (5.22)$$

is the information matrix for $\boldsymbol{\theta}^*$ with respect to $\mathbf{M}(w_i, \alpha_i, n_i)$.

Theorem 5.2. The random vector \mathbf{Y}_n is an asymptotically unbiased estimator of $\boldsymbol{\theta}^*$.

Proof. The expectation of \mathbf{Y}_n is given by

$$\mathbb{E}[\mathbf{Y}_n] = \mathbb{E} \left[\left(\sum_{i=1}^r \mathbf{A}_i \right)^{-1} \left(\sum_{i=1}^r \mathbf{A}_i \mathbf{Y}^{(i)}(n_i | w_i, \alpha_i) \right) \right] \quad (a)$$

$$= \left(\sum_{i=1}^r \mathbf{A}_i \right)^{-1} \left(\sum_{i=1}^r \mathbf{A}_i \mathbb{E}[\mathbf{Y}^{(i)}(n_i | w_i, \alpha_i)] \right) \quad (b)$$

$$= \left(\sum_{i=1}^r \mathbf{A}_i \right)^{-1} \left(\sum_{i=1}^r \mathbf{A}_i \right) \boldsymbol{\theta}^* \quad (c)$$

$$= \boldsymbol{\theta}^*. \quad (d)$$

□

Theorem 5.3. The variance-covariance of \mathbf{Y}_n is given by

$$\text{Var}[\mathbf{Y}_n] = \left(\sum_{i=1}^r n_i \mathcal{I}(\boldsymbol{\theta}^* | w_i, \alpha_i) \right)^{-1}. \quad (5.23)$$

Proof. Let

$$\mathbf{B} = \left(\sum_{i=1}^r \mathbf{A}_i \right)^{-1}. \quad (a)$$

The variance-covariance of \mathbf{Y}_n is given by

$$\text{Var} [\mathbf{Y}_n] = \text{Var} \left[\mathbf{B} \left(\sum_{i=1}^r \mathbf{A}_i \mathbf{Y}^{(i)}(n_i | w_i, \alpha_i) \right) \right] \quad (\text{b})$$

$$= \mathbf{B} \text{Var} \left[\sum_{i=1}^r \mathbf{A}_i \mathbf{Y}^{(i)}(n_i | w_i, \alpha_i) \right] \mathbf{B}^\top \quad (\text{c})$$

$$= \mathbf{B} \left(\sum_{i=1}^r \text{Var} [\mathbf{A}_i \mathbf{Y}^{(i)}(n_i | w_i, \alpha_i)] \right) \mathbf{B}^\top \quad (\text{d})$$

$$= \mathbf{B} \left(\sum_{i=1}^r \mathbf{A}_i \text{Var} [\mathbf{Y}^{(i)}(n_i | w_i, \alpha_i)] \mathbf{A}_i^\top \right) \mathbf{B}^\top. \quad (\text{e})$$

By ??, the variance-covariance of $\mathbf{Y}^{(i)}(n_i | w_i, \alpha_i)$ is given by

$$\frac{1}{n_i} \mathcal{I}^{-1}(\boldsymbol{\theta}^* | w_i, \alpha_i). \quad (\text{f})$$

By eq. (5.22), this is equivalent to \mathbf{A}_i^{-1} . Performing this substitution results in

$$\text{Var} [\mathbf{Y}_n] = \mathbf{B} \left(\sum_{i=1}^r \mathbf{A}_i \mathbf{A}_i^{-1} \mathbf{A}_i^\top \right) \mathbf{B}^\top \quad (\text{g})$$

$$= \mathbf{B} \left(\sum_{i=1}^r \mathbf{A}_i^\top \right) \mathbf{B}^\top. \quad (\text{h})$$

The summation is equivalent to $(\mathbf{B}^\top)^{-1}$. Performing this substitution results in

$$\text{Var} [\mathbf{Y}_n] = \mathbf{B} (\mathbf{B}^\top)^{-1} \mathbf{B}^\top \quad (\text{i})$$

$$= \mathbf{B} \quad (\text{j})$$

By eq. (a), \mathbf{B} is given by

$$\left(\sum_{i=1}^r \mathbf{A}_i \right)^{-1} \quad (\text{k})$$

and by eq. (5.22), \mathbf{A}_i is given by

$$n_i \mathcal{I}(\boldsymbol{\theta}^* | w_i, \alpha_i). \quad (\text{l})$$

Performing these substitution results in

$$\text{Var} [\mathbf{Y}_n] = \left(\sum_{i=1}^r n_i \mathcal{I}(\boldsymbol{\theta}^* | w_i, \alpha_i) \right)^{-1}. \quad (\text{m})$$

□

The weighted estimator asymptotically achieves the Cramér-Rao lower-bound as given by eq. (5.5) thus it is the asymptotic UMVUE estimator of $\boldsymbol{\theta}^*$ given a random sample of masked samples failures in which n_i realizations have w_i α_i -masked component failures for $i = 1, \dots, r$.

A linear combination of multivariate normal distributions is a multivariate normal distribution, thus the asymptotic sampling distribution of $\hat{\boldsymbol{\theta}}_{\mathbf{n}}$ is normally distributed.

Postulate 5.5. *As $n \rightarrow \infty$, the sampling distribution of $\hat{\boldsymbol{\theta}}_{\mathbf{n}}$ converges in distribution to a multivariate normal with a mean $\boldsymbol{\theta}^*$ and a variance-covariance given by eq. (5.23), written*

$$\mathbf{Y}_{\mathbf{n}} \xrightarrow{d} \text{MVN} \left(\boldsymbol{\theta}^*, \left(\sum_{i=1}^r n_i \mathcal{I}(\boldsymbol{\theta}^* | w_i) \right)^{-1} \right). \quad (5.24)$$

The generative model for $\mathbf{Y}_{\mathbf{n}}$ is given by algorithm 4.

Algorithm 2: Generative model of maximum likelihood estimator

Input:

$\boldsymbol{\Theta}^*$, the true parameter value of the series system.

Output:

$\hat{\boldsymbol{\theta}}_{\mathbf{n}}$, a realization of $\mathbf{Y}_{\mathbf{n}}$.

```

1 Model generate_mle( $\boldsymbol{\Theta}^*$ )
2   draw accuracy  $\alpha \sim p_A(\cdot)$ 
3   draw  $\alpha$ -masked failure cardinality  $w \sim p_{W|A}(\cdot | \alpha)$ 
4    $\hat{\boldsymbol{\theta}}_{\mathbf{n}} \leftarrow \text{generate\_mle}(\boldsymbol{\Theta}^*, w, \alpha)$ 
5   return  $\hat{\boldsymbol{\theta}}_{\mathbf{n}}$ 

```

Bibliography

- [1] Peter Bickel and Kjell Doksum. *Mathematical Statistics*, volume 1, chapter 2, page 117. Prentice Hall, 2 edition, 2000.

A Alternative proof of equation 3.6

Equation (3.6) on page 10 asserts that

$$f_{K,S}(k, t | \boldsymbol{\Theta}^*) = f_k(t | \boldsymbol{\Theta}_k^*) \prod_{\substack{j=1 \\ j \neq k}}^m R_j(t | \boldsymbol{\Theta}_j^*).$$

Proof. First, note that $f_{K,S}(k, t | \boldsymbol{\Theta}^*)$ is a proper density since

$$\int_0^\infty \sum_{j=1}^m f_{K,S}(j, t | \boldsymbol{\Theta}^*) dt = \int_0^\infty f_S(t | \boldsymbol{\Theta}^*) dt = 1. \quad (\text{a})$$

Consider a 3-out-of-3 system. By assumption 2, T_1 , T_2 , and T_3 are mutually independent. Thus, the joint density that $T_1 = t_1$, $T_2 = t_2$, and $T_3 = t_3$ is given by

$$f_{T_1, T_2, T_3}(t_1, t_2, t_3 | \Theta^*) = \prod_{p=1}^3 f_p(t_p | \Theta_p^*). \quad (b)$$

Component 1 causes a system failure during the interval $(0, t)$, $t > 0$, if component 1 fails during the given interval and components 2 and 3 survive longer than component 1. That is, $0 < T_1 < t$, $T_1 < T_2$, and $T_1 < T_3$.

Let $p(t)$ denote the probability that component 1 causes a system failure during the interval $(0, t)$,

$$p(t) = P[0 < T_1 < t \cap T_1 < T_2 \cap T_1 < T_3]. \quad (c)$$

This probability is given by

$$p(t) = \int_{t_1=0}^{t_1=t} \int_{t_2=t_1}^{t_2=\infty} \int_{t_3=t_1}^{t_3=\infty} \prod_{p=1}^3 f_p(t_p | \Theta_p^*) dt_3 dt_2 dt_1. \quad (d)$$

Performing the integration over t_3 results in

$$p(t) = \int_{t_1=0}^{t_1=t} \int_{t_2=t_1}^{t_2=\infty} \prod_{p=1}^2 f_p(t_p | \Theta_p^*) R_3(t_1 | \Theta_3^*) dt_2 dt_1. \quad (e)$$

Performing the integration over t_2 results in

$$p(t) = \int_{t_1=0}^{t_1=t} f_1(t_1 | \Theta_1^*) \prod_{p=2}^3 R_p(t_1 | \Theta_p^*) dt_1. \quad (f)$$

Taking the derivative of $p(t)$ with respect to t produces the density of probability near t . By the Second Fundamental Theorem of Calculus,

$$\frac{dp}{dt} = f_1(t | \Theta_1^*) \prod_{p=2}^3 R_p(t | \Theta_p^*). \quad (g)$$

The probability density $\frac{dp}{dt}$ is equivalent to the joint density that $K = 1$ and $S = t$ as given by theorem 3.4, i.e., $\frac{dp}{dt} = f_{K,S}(1, t | \Theta^*)$. Generalizing from this completes the proof. \square

B Alternative proof of equation ??

By ??, S and W are independent and the marginal distribution of W is independent of θ^* , thus the likelihood of observing a particular realization, $\mathbf{M}_n = \mathbf{m}_n$, is given by

$$\begin{aligned} f_{\mathbf{M}_n}(\mathbf{m}_n | \theta^*) &= \prod_{i=1}^n f_{C,S|W}(c_i, t_i | |c_i|, \theta^*) p_W(|c_i|) \\ &= \left[\prod_{i=1}^n f_{C,S|W}(c_i, t_i | |c_i|, \theta^*) \right] \left[\prod_{i=1}^n p_W(|c_i|) \right]. \end{aligned} \quad (B.1)$$

If we fix \mathbf{m}_n and allow the parameter $\boldsymbol{\theta}$ to change, we have the likelihood function

$$\mathcal{L}(\boldsymbol{\theta} | \mathbf{m}_n) = c \prod_{i=1}^n f_{C,S|W}(\mathcal{C}_i, \mathbf{t}_i | |\mathcal{C}_i|, \boldsymbol{\theta}^*) , \quad (\text{B.2})$$

where the product over $p_W(\cdot)$ is constant with respect to $\boldsymbol{\theta}$ and has been relabeled as c . The log-likelihood with respect to $\boldsymbol{\theta}$ then is given by

$$\begin{aligned} \ell(\boldsymbol{\theta} | \mathbf{m}_n) &= \ln c + \sum_{i=1}^n \ln f_{C,S|W}(\mathcal{C}_i, \mathbf{t}_i | |\mathcal{C}_i|, \boldsymbol{\theta}) \\ &= c' + \sum_{w=1}^{m-1} \left[\sum_{i \in \mathbb{A}(w)} \ln f_{C,S|W}(\mathcal{C}_i, \mathbf{t}_i | w, \boldsymbol{\theta}) \right] , \end{aligned} \quad (\text{B.3})$$

where $\mathbb{A}(w) = \{i \in \{1, \dots, n\} : |\mathcal{C}_i| = w\}$. Let $\mathbf{m}_{n_w} = \{(\mathcal{C}_i, \mathbf{t}_i) : i \in \mathbb{A}(w)\}$. The part in brackets is the log-likelihood $\ell(\boldsymbol{\theta} | \mathbf{m}_{n_w})$. Performing this substitution yields

$$\ell(\boldsymbol{\theta} | \mathbf{m}_n) = \ln c + \sum_{w=1}^{m-1} \ell(\boldsymbol{\theta} | \mathbf{m}_{n_w}) . \quad (\text{B.4})$$

The score function is given by

$$\mathbf{s}(\boldsymbol{\theta} | \mathbf{m}_n) = \sum_{w=1}^{m-1} \mathbf{s}(\boldsymbol{\theta} | w, \mathbf{m}_{n_w}) , \quad (\text{B.5})$$

The information matrix may be computed by taking the negative of the expectation of the Jacobian of the score over the true joint distribution of S , C , and W , giving

$$\mathcal{I}(\boldsymbol{\theta}^*) = -\mathbb{E}_{\boldsymbol{\theta}^*} \left[\mathcal{J}(\mathbf{s}(\boldsymbol{\theta} | \mathbf{M}_1)) \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}^*} \right] . \quad (\text{B.6})$$

Substituting eq. (B.5) with its definition gives

$$\mathcal{I}(\boldsymbol{\theta}^*) = -\mathbb{E}_{\boldsymbol{\theta}^*} \left[\mathcal{J} \left(\sum_{w=1}^{m-1} \mathbf{s}(\boldsymbol{\theta} | w, \mathbf{M}(w, 1)) \right) \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}^*} \right] . \quad (\text{B.7})$$

The Jacobian is a linear operator so it can be moved inside the summation, giving

$$\mathcal{I}(\boldsymbol{\theta}^*) = -\mathbb{E}_{\boldsymbol{\theta}^*} \left[\sum_{w=1}^{m-1} \mathcal{J}(\mathbf{s}(\boldsymbol{\theta} | w, \mathbf{M}(w, 1))) \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}^*} \right] . \quad (\text{B.8})$$

The Jacobian of \mathbf{s} is the Hessian of ℓ . Performing this substitution gives

$$\mathcal{I}(\boldsymbol{\theta}^*) = -\mathbb{E}_{\boldsymbol{\theta}^*} \left[\sum_{w=1}^{m-1} \mathcal{H}(\ell(\boldsymbol{\theta} | w, \mathbf{M}(w, 1))) \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}^*} \right] . \quad (\text{B.9})$$

The expectation is a linear operator so it may be moved inside the summation, giving

$$\mathcal{I}(\boldsymbol{\theta}^*) = - \sum_{w=1}^{m-1} \mathbb{E}_{\boldsymbol{\theta}^*} \left[\mathcal{H}(\ell(\boldsymbol{\theta} | w, \mathbf{M}(w, 1))) \right]_{\boldsymbol{\theta}=\boldsymbol{\theta}^*}. \quad (\text{B.10})$$

The expectation sums over the probability mass function $W \sim p_W(\cdot)$, giving

$$\mathcal{I}(\boldsymbol{\theta}^*) = - \sum_{w=1}^{m-1} p_W(w) \mathbb{E}_{\boldsymbol{\theta}^*} \left[\mathcal{H}(\ell(\boldsymbol{\theta} | w, \mathbf{M}(w, 1))) \right], \quad (\text{B.11})$$

where the expectation is now over the true marginal joint distribution of C and S given $W = w$. The expectation of the Hessian of ℓ is equivalent to the negative of $\mathcal{I}(\boldsymbol{\theta}^* | w, \mathbf{M}(w, 1))$. Performing this substitution gives

$$\mathcal{I}(\boldsymbol{\theta}^*) = \sum_{w=1}^{m-1} p_W(w) \mathcal{I}(\boldsymbol{\theta}^* | w). \quad (\text{B.12})$$

□

C Proof of corollary 3.7.5

Corollary 3.7.5 on page 15 asserts that K given S is conditionally independent of C , W , and A .

Proof. By the laws of probability,

$$p_{K|C,S,W,A}(k|\mathcal{C}, t, w, \alpha, \boldsymbol{\Theta}^*) = \frac{f_{K,C,S|W,A}(k, \mathcal{C}, t|w, \alpha, \boldsymbol{\Theta}^*)}{f_{C,S|W,A}(\mathcal{C}, t|w, \alpha, \boldsymbol{\Theta}^*)} \quad (\text{a})$$

which may be rewritten as

$$p_{K|C,S,W,A}(k|\mathcal{C}, t, w, \alpha, \boldsymbol{\Theta}^*) = \frac{p_{C|K,S,W,A}(\mathcal{C}|k, t, w, \alpha, \boldsymbol{\Theta}^*) f_{K,S|W,A}(k, t|w, \alpha, \boldsymbol{\Theta}^*)}{f_{C,S|W,A}(\mathcal{C}, t|w, \alpha, \boldsymbol{\Theta}^*)}. \quad (\text{b})$$

By assumptions 5 and 6, we may simplify the above equation to

$$p_{K|C,S,W,A}(k|\mathcal{C}, t, w, \alpha, \boldsymbol{\Theta}^*) = \frac{p_{C|K,W,A}(\mathcal{C}|k, w, \alpha) f_{K,S}(k, t|\boldsymbol{\Theta}^*)}{f_{C,S|W,A}(\mathcal{C}, t|w, \alpha, \boldsymbol{\Theta}^*)}. \quad (\text{c})$$

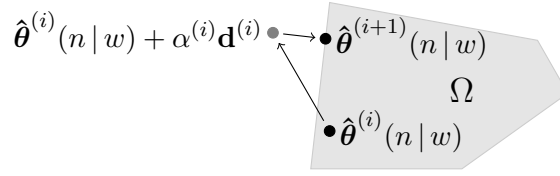
By eqs. (3.8) and (3.9), the above equation may be rewritten as

$$p_{K|C,S,W,A}(k|\mathcal{C}, t, w, \alpha, \boldsymbol{\Theta}^*) = \frac{\frac{R_S(t|\boldsymbol{\Theta}^*)}{\binom{m-1}{w}} \left((1-\alpha)h_S(t|\boldsymbol{\Theta}^*) - (1-\alpha\frac{m}{w}) \sum_{j \in \mathcal{C}} h_j(t|\boldsymbol{\Theta}_j^*) \right) h_k(t|\boldsymbol{\Theta}_k^*)}{\frac{R_S(t|\boldsymbol{\Theta}^*)}{\binom{m-1}{w}} \left((1-\alpha)h_S(t|\boldsymbol{\Theta}^*) - (1-\alpha\frac{m}{w}) \sum_{k \in \mathcal{C}} h_k(t|\boldsymbol{\Theta}_k^*) \right) h_S(t|\boldsymbol{\Theta}^*)}. \quad (\text{d})$$

The above equation may be simplified to

$$p_{K|C,S,W,A}(k|\mathcal{C}, t, w, \alpha, \boldsymbol{\Theta}^*) = \frac{h_k(t|\boldsymbol{\Theta}_k^*)}{h_S(t|\boldsymbol{\Theta}^*)} \quad (\text{e})$$

which is the same as $p_{K|S}(k|t, \boldsymbol{\Theta}^*)$ and thus K given S is conditionally independent of C , W , and A . □

Figure 3: Projection onto convex set Ω .

D Numerical solutions to the MLE

The function $\ell(\theta | w, \mathbf{m}_n)$ is the log-likelihood with respect to $\theta \in \Omega$ where $\Omega \subset \mathbb{R}^{m \cdot q}$. This function has a surface in an $(m \cdot q + 1)$ dimensional space, where a particular point on this surface represents the log-likelihood of observing \mathbf{m}_n with respect to θ . By definition 5.3, $\hat{\theta}(n | w)$ is the point on this surface that is at a maximum,

$$\hat{\theta}(n | w) = \arg \max_{\theta \in \Omega} \ell(\theta | w, \mathbf{m}_n). \quad (5.9 \text{ revisited})$$

Generally there is no closed-form solution that solves $\hat{\theta}(n | w)$, in which case iterative search methods may be used to numerically approximate a solution.

The general version of iterative search that numerically approximates a solution to eq. (5.10), subject to the constraint $\theta^* \in \Omega$, is shown in algorithm 3. Since iterative search is a local search method, it may fail to converge to a global maximum.

Algorithm 3: Iterative maximum likelihood search

Result: an approximate solution to the stationary points of the maximum likelihood equation

Input:

Θ^* , the true parameter index.

Output:

$\hat{\theta}(n | w)$, an approximation of the maximum likelihood estimate.

```

1 Model find_mle( $\mathbf{m}_n$ )
2    $\hat{\theta}^{(0)}(n | w) \leftarrow$  an initial starting point in  $\Omega$ 
3    $i \leftarrow 0$ 
4   while stopping criteria not satisfied do
5      $\hat{\theta}^{(i+1)}(n | w) \leftarrow \text{project} \left( \hat{\theta}^{(i)}(n | w) + \alpha^{(i)} \mathbf{d}^{(i)}, \Omega \right)$ 
6      $i \leftarrow i + 1$ 
7   end
8   return  $\hat{\theta}^{(i)}(n | w)$ 

```

Assume parameter space Ω is convex. Then, the function $\text{project}(\theta, \Omega)$ in algorithm 3 projects any point θ to the nearest point in the parameter space Ω as depicted by fig. 3. Thus, the search method is restricted to searching over the feasible parameter space.

In algorithm 3, $\mathbf{d}^{(i)}$ is a unit vector denoting a *promising* direction in which to search for a better solution than $\hat{\theta}^{(i)}(n | w)$. The Fisher scoring algorithm is a slight variation of Newton-Raphson³ in

³If the *observed* information matrix is used instead of the *expected* information matrix, the Fisher scoring algorithm is equivalent to Newton-Raphson.

which

$$\mathbf{d}^{(i)} = \mathcal{I}^{-1} \left(\hat{\boldsymbol{\theta}}^{(i)}(n|w) \right) \mathbf{s} \left(\hat{\boldsymbol{\theta}}^{(i)}(n|w) \right). \quad (\text{D.1})$$

The scalar $\alpha^{(i)} > 0$ is the distance to move from point $\hat{\boldsymbol{\theta}}^{(i)}(n|w)$ to generate the next point, $\hat{\boldsymbol{\theta}}^{(i+1)}(n|w)$. Most naturally, the solution to

$$\alpha^{(i)} = \arg \max_{\alpha > 0} \ell(\hat{\boldsymbol{\theta}}^{(i)}(n|w) + \alpha \cdot \mathbf{d}^{(i)}) \quad (\text{D.2})$$

is desired, e.g., using *Golden Section* line search to find an approximate solution. Finally, the iterations are repeated until some *stopping criteria* is satisfied. Most naturally, this is given by

$$\| \mathbf{s}(\hat{\boldsymbol{\theta}}^{(i)}(n|w)) \| < \epsilon, \quad (\text{D.3})$$

signifying a stationary point has been reached.

E General joint distribution

By the axioms of probability, the joint distribution of C, S, W, and A may be decomposed into the chain

$$f_{C,S,W,A}(\mathcal{C}, t, w, \alpha | \boldsymbol{\theta}^*) = p_{C|S,W,A}(\mathcal{C}|t, w, \alpha, \boldsymbol{\theta}^*) p_{W,A|S}(w, \alpha | t, \boldsymbol{\theta}^*) f_S(t | \boldsymbol{\theta}^*). \quad (\text{E.1})$$

If we remove ??, the distribution of W and A may be dependent on S but we still they do not carry information about the true parameter index $\boldsymbol{\theta}^*$, thus we may rewrite the above equation as

$$f_{C,S,W,A}(\mathcal{C}, t, w, \alpha | \boldsymbol{\theta}^*) = p_{C|S,W,A}(\mathcal{C}|t, w, \alpha, \boldsymbol{\theta}^*) p_{W,A|S}(w, \alpha | t) f_S(t | \boldsymbol{\theta}^*). \quad (\text{E.2})$$

Since W and A carry no information about the true parameter index $\boldsymbol{\theta}^*$, the maximum likelihood estimate $\hat{\boldsymbol{\theta}}_{\mathbf{n}}$ is independent of the distribution of W and A.

Proof.

$$\begin{aligned} \hat{\boldsymbol{\theta}}_{\mathbf{n}} &= \arg \max_{\boldsymbol{\theta}} \ell_n(\boldsymbol{\theta}) \\ &= \arg \max_{\boldsymbol{\theta}} \left[\sum_{i=1}^n \ln p_{C|S,W,A}(\mathcal{C}_i | t_i, w_i, \alpha_i, \boldsymbol{\theta}) \right. \\ &\quad \left. + \sum_{i=1}^n \ln f_S(t_i | \boldsymbol{\theta}) + \underbrace{\sum_{i=1}^n \ln p_{W,A|S}(w_i, \alpha_i | t_i)}_{\text{constant}} \right], \end{aligned} \quad (\text{a})$$

where $w_i = |\mathcal{C}_i|$ and the last term is constant with respect to $\boldsymbol{\theta}$ and thus may be dropped from the maximum likelihood equation without changing the point that maximizes the likelihood. \square

The sampling distribution of $\hat{\boldsymbol{\theta}}_{\mathbf{n}}$ is a function of $\mathbf{M}_{\mathbf{n}}$ (as opposed to a particular realization), and since $\mathbf{M}_{\mathbf{n}}$ is a function of W and A, the sampling distribution of $\hat{\boldsymbol{\theta}}_{\mathbf{n}}$ is also a function of W and A.

Given a statistical model, the asymptotic sampling distribution of the maximum likelihood estimator is nearly automatic. The most taxing (and uncertain) part is the *model selection*. Since we may not have much confidence in any particular model of the joint distribution of W and S , the *expected* information matrix is sketchy. In general, the *observed* information matrix is preferred. The sampling distribution of $\hat{\boldsymbol{\theta}}_n$ converges in distribution to a multivariate normal with a mean given by $\boldsymbol{\theta}^*$ and a variance-covariance given by the inverse of the *observed* information matrix, written

$$\mathbf{Y}_n \xrightarrow{d} \text{MVN}(\boldsymbol{\theta}^*, \mathbf{J}_n^{-1}(\boldsymbol{\theta}^*)) . \quad (\text{E.3})$$

If an accurate distribution of W and A given S could be constructed such that it is dependent on $\boldsymbol{\theta}^*$, then the distribution of W and A in a sample carries extra information about $\boldsymbol{\theta}^*$. A plausible model may depend upon $H(K|S = t)$.

F Sampling distribution of functions of parameters

Suppose we have a characteristic of interest $\mathbf{g}: \mathbb{R}^{m \cdot q} \mapsto \mathbb{R}^p$ that is a function of the parametric model. The *true* value of the characteristic is given by $\mathbf{g}(\boldsymbol{\theta}^*)$. By the invariance property of maximum likelihood estimators, if $\hat{\boldsymbol{\theta}}_n$ is the maximum likelihood estimator of $\boldsymbol{\theta}^*$ then $\mathbf{g}(\hat{\boldsymbol{\theta}}_n)$ is the maximum likelihood estimator of $\mathbf{g}(\boldsymbol{\theta}^*)$.

By postulate 5.4, \mathbf{Y}_n is a random vector drawn from a multivariate normal distribution,

$$\mathbf{Y}_n \sim \text{MVN}(\boldsymbol{\theta}^*, \text{Var}[\mathbf{Y}_n]) , \quad (5.24 \text{ revisited})$$

therefore $\mathbf{g}(\mathbf{Y}_n)$ is a random vector. Under the regularity conditions (see assumption 10), $\mathbf{g}(\mathbf{Y}_n)$ is asymptotically normally distributed with a mean given by $\mathbf{g}(\boldsymbol{\theta}^*)$ and a variance-covariance given by $\text{Var}[\mathbf{g}(\mathbf{Y}_n)]$, written

$$\mathbf{g}(\mathbf{Y}_n) \xrightarrow{d} \text{MVN}(\mathbf{g}(\boldsymbol{\theta}^*), \text{Var}[\mathbf{g}(\mathbf{Y}_n)]) . \quad (\text{F.1})$$

The generative model may be used to generate samples of statistics that are functions of the true parameter index $\boldsymbol{\theta}^*$. A sample drawn from $\mathbf{g}(\mathbf{Y}_n)$ may be generated by drawing r maximum likelihood estimates from the sampling distribution

$$\langle \hat{\boldsymbol{\theta}}_n^{(1)}, \dots, \hat{\boldsymbol{\theta}}_n^{(r)} \rangle , \quad (\text{F.2})$$

and applying \mathbf{g} to each, resulting in the sample

$$\langle \mathbf{g}(\hat{\boldsymbol{\theta}}_n^{(1)}), \dots, \mathbf{g}(\hat{\boldsymbol{\theta}}_n^{(r)}) \rangle . \quad (\text{F.3})$$

An estimate of the variance-covariance $\text{Var}[\mathbf{g}(\mathbf{Y}_n)]$ is given by the sample covariance

$$\begin{aligned} \hat{\text{Var}}[\mathbf{g}(\mathbf{Y}_n)] = \\ \frac{1}{r} \sum_{i=1}^r \left(\mathbf{g}(\hat{\boldsymbol{\theta}}_n^{(i)}) - \mathbf{g}(\hat{\boldsymbol{\theta}}_n) \right) \left(\mathbf{g}(\hat{\boldsymbol{\theta}}_n^{(i)}) - \mathbf{g}(\hat{\boldsymbol{\theta}}_n) \right)^T . \end{aligned} \quad (\text{F.4})$$

Thus,

$$\mathbf{g}(\mathbf{Y}_n) \sim \text{MVN}(\mathbf{g}(\hat{\boldsymbol{\theta}}_n), \hat{\text{Var}}[\mathbf{Y}_n]) . \quad (\text{F.5})$$

G Bootstrap of sample covariance

Another estimator of the variance-covariance of \mathbf{Y}_n is given by the sample covariance.

Definition G.1. Given a sample of r maximum likelihood estimates, $\hat{\theta}_n^{(1)}, \dots, \hat{\theta}_n^{(r)}$, the maximum likelihood estimator of the variance-covariance of the sampling distribution of $\hat{\theta}_n$ is given by the sample covariance,

$$\hat{\text{Var}}[\mathbf{Y}_n] = \frac{1}{r} \sum_{i=1}^r \left(\hat{\theta}_n^{(i)} - \theta^* \right) \cdot \left(\hat{\theta}_n^{(i)} - \theta^* \right)^\top. \quad (\text{G.1})$$

The sample covariance depends upon a random sample of maximum likelihood estimates, and thus the sample covariance is a random matrix. However, by the property that maximum likelihood estimators converge in probability to the true value,

$$\hat{\text{Var}}[\mathbf{Y}_n] \xrightarrow{P} \text{Var}[\mathbf{Y}_n], \quad (\text{G.2})$$

it follows that

$$\mathbf{Y}_n \xrightarrow{d} \text{MVN}(\theta^*, \text{Var}[\mathbf{Y}_n]). \quad (\text{G.3})$$

If only one maximum likelihood estimate is realized, $\mathbf{Y}_n = \hat{\theta}_n$, then the sample covariance given by ?? cannot be computed. In the *parametric Bootstrap*, the sample covariance is approximated by using $\hat{\theta}_n$ as an estimate of θ^* in the generative model described by ??, i.e.,

$$\hat{\theta}_n^{(1)}, \dots, \hat{\theta}_n^{(r)} \leftarrow \text{generate_mle}(\hat{\theta}_n), \quad (\text{G.4})$$

and the sample covariance of these approximate maximum likelihood estimates is an asymptotically unbiased estimator of $\text{Var}[\mathbf{Y}_n]$.

Algorithm 4: Sample covariance of a bootstrapped sample of maximum likelihood estimates.

Input:

$\hat{\theta}_n$, an estimate of θ^* .

Output:

sample covariance of a bootstrapped sample of maximum likelihood estimates.

1 **Model** *BootstrapCovariance*

2 | $(\hat{\theta}_n)$

3 **for** $i = 1$ to r **do**

4 | $\hat{\theta}_n^{(i)} \leftarrow \text{generate_mle}(\hat{\theta}_n)$

5 **end**

6 **return** sample covariance of $\hat{\theta}_n^{(1)}, \dots, \hat{\theta}_n^{(r)}$
