

Roles for Statistical Models

- Data Reduction and factorization of the likelihood function.
 - *Sufficient* Statistics
 - *Ancillary* Statistics
- Symmetry and Independence assumptions
 - deFinetti's theorem on *exchangeable sequences*
- Properties of *Maximum Likelihood*

Data Reduction Concepts for Statistical Models

Defn: The (dimensional) random variable $Y = \mathbf{g}(X)$ is *sufficient* for the parameter θ (with respect to X) *iff*

$$\mathbf{P}(X | Y, \theta) = \mathbf{P}(X | Y), \text{ independent of } \theta.$$

Theorem: The likelihood for θ given a sufficient (set of) statistic(s) Y is the same as the likelihood for θ given the (dimensional) variable X for which Y is sufficient.

Proof: $\mathbf{P}(x | \theta) = \mathbf{P}(x, y | \theta)$ as $Y = \mathbf{g}(X)$
 $= \mathbf{P}(x | y, \theta) \mathbf{P}(y | \theta)$ multiplication axiom
 $= \mathbf{P}(x | y) \mathbf{P}(y | \theta)$ by sufficiency of Y
 $\propto \mathbf{P}(y | \theta)$

Corollary (Factorization of the likelihood function):

$Y = \mathbf{g}(X)$ is *sufficient* for the parameter θ (with respect to X) *iff*

The likelihood (probability or density) function can be written as the product of two functions of this form:

$$\mathbf{P}(X | \theta) = \mathbf{h}(X) \mathbf{j}(Y, \theta).$$

Recall: $Y = \mathbf{g}(X)$

Example 1 (coin-tossing, again):

$\mathbf{X} = \langle X_1, \dots, X_n \rangle$ are *iid* Bernoulli trials given θ , with $\mathbf{P}(X_1=1|\theta) = \theta$, $0 < \theta < 1$.

Claim: $\mathbf{g}(\mathbf{X}) = \mathbf{Y} = \langle \sum_i X_i, n - \sum_i X_i \rangle$ is a sufficient reduction to the two statistics, #1's = $\sum_i X_i = k$ and #0's = $n - k$ in the sequence \mathbf{X} .

Proof: $\mathbf{P}(\mathbf{x}, \mathbf{y} | \theta) = \mathbf{P}(\mathbf{x} | \theta) = \theta^k (1-\theta)^{n-k}$

$$\mathbf{P}(\mathbf{y} | \theta) = \binom{n}{k} \theta^k (1-\theta)^{n-k}$$

Thus $\mathbf{P}(\mathbf{x} | \mathbf{y}, \theta) = \mathbf{P}(\mathbf{x}, \mathbf{y} | \theta) / \mathbf{P}(\mathbf{y} | \theta) = \mathbf{P}(\mathbf{x} | \mathbf{y}) = k!(n-k)!/n!$

That is, $\mathbf{P}(\mathbf{X} | \mathbf{y}, \theta)$ is a discrete, uniform distribution over all sequences H_n that begin with k 1's and $(n-k)$ 0's, independent of θ .

Or, use factorization and note that, alternatively $\langle \bar{X}, n \rangle$ are sufficient for θ as

$$\mathbf{P}(\mathbf{X} | \theta) = \theta^{n\bar{X}} (1-\theta)^{n(1-\bar{X})} = \mathbf{h}(\mathbf{X}) \mathbf{j}(Y, \theta)$$

where $\mathbf{h}(\mathbf{X}) = 1$ and $Y = \bar{X} = \sum_i X_i / n$.

Example 2 (Normal distribution, known variance):

$\mathbf{X} = \langle X_1, \dots, X_n \rangle$ are *iid* normal $N(\mu, 1)$ trials.

Claim: The pair $\langle \bar{X}, n \rangle$ is sufficient for μ .

Proof: Write $\mathbf{p}(X | \mu) =$

$$(2\pi)^{-n/2} \exp\left(-\sum_i (X_i - \bar{X})^2 / 2\right) \exp(-n(\mu - \bar{X})^2 / 2),$$

where

$$\underbrace{(2\pi)^{-n/2} \exp\left(-\sum_i (X_i - \bar{X})^2 / 2\right)}_{h(X)} \underbrace{\exp(-n(\mu - \bar{X})^2 / 2)}_{j(Y, \theta)}$$

$$\begin{array}{c} \uparrow \\ h(X) \end{array}$$

$$\begin{array}{c} \uparrow \\ j(Y, \theta) \end{array}$$

Defn: The (dimensional) random variable $Y = \mathbf{g}(X)$ is *ancillary* for the parameter θ (with respect to X) *iff*

$$\mathbf{P}(Y | \theta) = \mathbf{P}(Y), \text{ independent of } \theta.$$

Theorem: The likelihood for θ based on an ancillary (set of) statistic(s) Y is *constant*.

Corollary: The likelihood for θ based on X equals the conditional likelihood for θ based on X , *given* Y .

$$\mathbf{P}(x | \theta) = \mathbf{P}(x | y, \theta)$$

Proof: $\mathbf{P}(x | \theta) = \mathbf{P}(x, y | \theta) = \mathbf{P}(x | y, \theta) \mathbf{P}(y | \theta)$
 $\propto \mathbf{P}(x | y, \theta).$

Example 3 (coin-tossing, again):

$\mathbf{X} = \langle X_1, \dots, X_i, \dots \rangle$ are *iid* Bernoulli trials given θ , with $\mathbf{P}(X_1=1|\theta) = \theta$, $0 < \theta < 1$.

$\mathbf{g}(\mathbf{X}) = \mathbf{Y} = \langle \sum_i X_i, n - \sum_i X_i \rangle$ is a sufficient reduction for inference about θ .

Version 3a: The *stopping rule* is sample to a fixed sample size n . Then N (sample size) is ancillary ($\mathbf{P}(N=n) = 1$) and, **given** $N = n$, $\sum_i X_i$ is sufficient!

Moreover, $\mathbf{P}(\sum_i X_i | n, \theta)$ is given by the *Binomial*(n, θ) distribution.

Version 3b: The *stopping rule* is sample to a fixed number of “heads,” say $\sum_i X_i = k$

Then $\sum_i X_i$ (number of heads) is ancillary ($\mathbf{P}(\sum_i X_i = k) = 1$) and, **given** $\sum_i X_i = k$, the number of flips N is sufficient!

Moreover, $\mathbf{P}(N | k, \theta)$ is given by the *Neg-Binomial*(k, θ) distribution.

However, regardless of the stopping rule, in either version, the *pair* $\langle \sum_i X_i, N \rangle$ is sufficient!

Recapitulation of data-reduction principles for statistical models

Sufficiency principle: A sufficient statistic preserves all the relevant information about the parameter that is in the full data set

Ancillarity principle: All the relevant information in the data set about the parameter is contained in the conditional model, given the ancillary statistic.

Likelihood principle: All the relevant information in the data set about the parameter is contained in the likelihood function given the data.

Birnbaum's Theorem: The *Likelihood* principle is equivalent to the conjunction of the *Sufficiency* and *Ancillarity* principles.

Identifying statistical models by symmetry & independence involving observables
(*deFinetti's Theorem*)

Heuristic Example (coin-tossing yet again!): Let $\mathbf{X} = \langle X_1, \dots, X_i, \dots \rangle$ be an infinite sequence of binary trials, with the σ -algebra (\mathfrak{A}) of events generated by the observable “historical” events $H_n: \langle x_1, \dots, x_n, \{0,1\}, \{0,1\}, \dots \rangle$.

Defn: Say that a probability \mathbf{P} over \mathfrak{A} is:

- *1-exchangeable* if for $\forall(i,j) \mathbf{P}(X_i = 1) = \mathbf{P}(X_j = 1)$

- *2-exchangeable* if $\forall(i_1, i_2, \text{distinct and } j_1, j_2 \text{ distinct})$

$$\mathbf{P}(X_{i_1} = x_1, X_{i_2} = x_2) = \mathbf{P}(X_{j_1} = x_1, X_{j_2} = x_2)$$

- *n-exchangeable* if $\forall(i_1, i_2, \dots, i_n \text{ distinct})$

$\mathbf{P}(X_{i_1} = x_1, X_{i_2} = x_2, \dots, X_{i_n} = x_n)$ does not depend on the n distinct $\langle i_1, i_2, \dots, i_n \rangle$

- *exchangeable* if \mathbf{P} is n -exchangeable for each n ($n = 1, 2, \dots$).

Theorem (deFinetti): \mathbf{P} is exchangeable if and only if \mathbf{P} can be written as

$$\mathbf{P}(E) = \int_{\Theta} \mathbf{P}(E \mid \theta) d\mathbf{Q}(\theta)$$

where

- $\mathbf{P}(\mathfrak{A} \mid \theta)$ is given by *iid* Bernoulli(θ) trials
- $\mathbf{Q}(\theta)$ is a prior probability distribution over Θ determined uniquely by \mathbf{P} over \mathfrak{A} .

Thus, one can use the computational benefits of sampling from an *iid* statistical model, “as if” it were true, given suitable exchangeability (symmetry) assumptions involving only the algebra of the observable random variables.

Remarks:

- This important theorem generalizes to cover both discrete and continuous random variables.
- Also, there is version dealing with finite sequences (N-exchangeability).
- For a thorough discussion of all this, see chapter 1 of Mark Schervish’s book, *Theory of Statistics*, 1995. Springer-Verlag.

Data reduction, Fisher-Information, and Maximum Likelihood

Defn.: Score function: $\mathbf{S}_X(\theta) = \frac{\partial (\ln \mathbf{p}(X | \theta))}{\partial \theta}$

Fisher Information (under general conditions)

$$I_X(\theta) = \text{Var}(\mathbf{S}_X(\theta)) = \mathbb{E}\left[\frac{\partial^2 (\ln \mathbf{p}(X | \theta))}{\partial \theta^2}\right].$$

- Fisher Information is additive for independent data.
- $I_X(\theta) = I_Y(\theta)$ whenever Y is sufficient for θ (with respect to X).
- Fisher Information is a differential form of Kullback-Leiber information.

Defn.: Let θ^* denote the argmax of the likelihood function $\mathbf{p}(X | \theta)$,
the *maximum likelihood estimate (MLE)* of the parameter.

Main Theorem (under general regularity conditions on the statistical model):

$$\mathbf{P}(\theta^* | \theta_0) \approx N(\theta_0, [I_X(\theta^*)]^{-1}) = N(\theta_0, [nI_{X_i}(\theta^*)]^{-1})$$

So (under “regularity” conditions) the *MLE*:

- Has an asymptotic Normal distribution.
- Is asymptotically consistent (converges to θ_0).
- Is asymptotically sufficient.