

Searches and Likelihood Ratio Tests

8 Oct 2021

Reed Essick

Learning Objectives

Basics of Decision Theory (how to build optimal searches)

Neyman-Pearson Lemma

Sufficient Statistics

Marginalization vs. maximization (max-likelihood estimators)

Bayes Factors vs. Odds Ratios (Odds Factors)

Computational Methods

Savage-Dickey Density Ratio

Thermodynamic Integration (power sampling)

Nested Sampling

Decision Theory as classification, hypothesis testing

Consider 2 classes of events, each of which will produce a set of observables x with some probability

$p(x|A)$: prob. of x given event is of type A

$p(x|B)$: prob. of x given event is of type B

We wish to determine which type a particular event is based on the observables for that event. What is the "most powerful test" to distinguish between the classes?

The Neyman-Pearson Lemma demonstrates an optimal statistic is the Likelihood Ratio Test.

$$\Lambda_B^A(x) \equiv \frac{p(x|A)}{p(x|B)}$$

This is derived by optimizing the true alarm prob. at a fixed false alarm prob., both defined as integrals over regions of observable space (R)

$$P_{\text{true alarm}}(R) = \int_R dx p(x|A)$$

$$P_{\text{false alarm}}(R) = \int_R dx p(x|B)$$

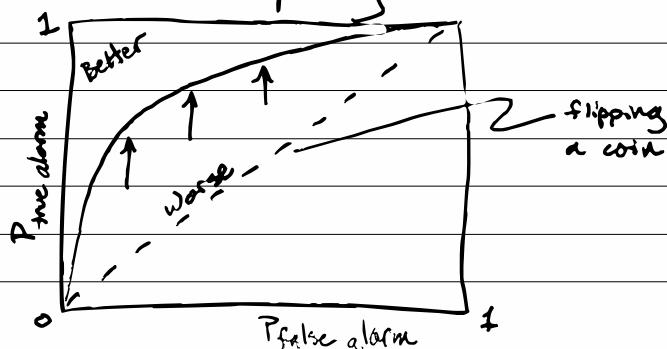
Defining regions ($R(\lambda_0)$: $\Lambda_B^A(x) \geq \lambda_0$) maximizes

$$\int_R dx p(x|A) - \lambda \left(P_0 - \int_R dx p(x|B) \right)$$

That is, this maximizes the prob. of correctly identifying an event of class A at a fixed prob. of misidentifying an event of class B.

This is often plotted as a Receiver Operating Characteristic (ROC) curve:

Λ_B^A maximizes the weight of this curve at every point in its domain.



The likelihood ratio test is of fundamental importance because it tells you how to order multi-dimensional observations. It is also a sufficient statistic.

Consider a set of observables: x . We can ask whether adding additional observables y will carry more information. That is, are we better able to distinguish events when we consider $y \oplus x$ than when we only consider x ? First, we need a concept of information.

KL-Divergence between
your distribution & product of
marginal distributions

The mutual information is defined as:

$$I(a, b) = \int d\alpha d\beta p(a, b) \log \left[\frac{p(a, b)}{p(a)p(b)} \right]$$

$$\text{w/ } p(a) = \int d\beta p(a, b) ; \quad p(b) = \int d\alpha p(a, b)$$

larger values of I imply more "correlation" between a & b .

Let us consider some function $c(b)$. The data processing inequality states

$$I(a, b) \geq I(a, c)$$

with inequality iff c is a sufficient statistic for b . That is, sufficient statistics contain all the information available.

The likelihood ratio test maps an N -dim space onto the real line. It retains all available information from the original space!

Bayes Factors (ratios of marginal likelihoods) and Odds Factors (ratios of posterior probabilities) are examples of likelihood ratio tests.

$$B_B^A = \frac{p(h \mid H_A)}{p(h \mid H_B)} = \frac{\int ds p(h \mid s) p(s \mid H_A)}{\int ds p(h \mid s) p(s \mid H_B)} \quad \begin{matrix} \text{evidence} \\ \text{ratio} \end{matrix}$$

$$\Omega_B^A = \frac{p(H_A \mid h)}{p(H_B \mid h)} = B_B^A \frac{p(H_A)}{p(H_B)}$$

↑ prior odds

∴ under the assumption that all modeling is correct (noise distn, signal model, etc.) then ordering events by the Bayes Factor (or Odds ratio) is an optimal test in that there are sufficient statistics to distinguish b/w hypotheses A & B.

Note, however, that the calculation of B often involves a high-dimensional integral over the signal space. This can be difficult & expensive in practice.

Searches, thus, often define maximum likelihood estimates as their search statistics.

Note! Maximization over Θ is not equivalent to minimization of L_{ML} over Θ

Max-likelihood estimators are often computationally efficient ways to obtain good detection statistics or estimates of signal parameters. In general, these are defined as

$$\hat{\Theta}_{MLE}(h) = \underset{\Theta}{\operatorname{argmax}} \{ p(h|\Theta) \} = \underset{\Theta}{\operatorname{argmax}} \{ \log p(h|\Theta) \}$$

$$L_{MLE} = p(h|\hat{\Theta}_{MLE}(h))$$

One can then use L_{MLE} as a detection statistic ; $\hat{\Theta}_{MLE}$ as an estimate of the parameters. MLE estimators are usually well behaved (unbiased, Fisher efficient).

In general, to build an Optimal Search you should

- ① write down the likelihood
- ② write down your signal model & prior
if you don't have a prior, skip to ③ b)
- ③ see if you can efficiently compute Bayes Factors between signal and noise hypotheses
 - a) if yes, use the Bayes Factor as detection stat.
 - b) if not, define max. likelihood estimate and use that as detection statistic
- ④ measure distribution of detection statistic under null hypothesis (noise-only model) to determine map from detection-statistic \rightarrow false alarm prob.
(or more generally the full ROC curve)

Worked Example: Matched Filtering in Stationary Gaussian Noise

Let's assume we know the signal's shape (h) perfectly but do not know its amplitude ($A \in \mathbb{R}$). This means our likelihood is

$$\log p(h|A) = -\frac{1}{2} \left(4 \int_0^{\infty} df \frac{|h - A|_f|^2}{S_n} \right)$$

where S_n is the One-Sided Power Spectral Density defined as

$$\langle n(t) n(t+\tau) \rangle = \frac{1}{2} \int df c^{2\pi i f \tau} S_n(f)$$

Note: if we define a prior for A under the signal model then we may be able to define a Bayes Factor for signal ($A > 0$) vs. noise ($A = 0$).

Let's instead define a MLE for A

$$\frac{\partial \log P}{\partial A} = -2 \int_0^{\infty} df \frac{2A|h|^2 - 2\operatorname{Re}\{h^* s\}}{S_n} = 0$$

$$\Rightarrow \hat{A} = \frac{4 \int_0^{\infty} df \operatorname{Re}\{h^* s\} / S_n}{4 \int_0^{\infty} df |h|^2 / S_n}$$

(the factors of 4 are conventional)

Now, we are free to choose the scale of the template \mathbf{h} and it is convenient to normalize it so that

$$4 \int_0^{\infty} df \frac{|\mathbf{h}|^2}{S_n} = 1$$

so that $\hat{A} = 4 \int_0^{\infty} df \frac{\operatorname{Re} \{ \mathbf{h}^* \mathbf{s} \}}{S_n} = \underline{\text{Signal-to-Noise Ratio}} \\ (\text{SNR})$

In matched filtering, the SNR is just the max. likelihood estimate for the normalized signal amplitude.

Let's plug this back into the likelihood

$$\begin{aligned} \log p(\mathbf{h} | \mathbf{x}, \omega) &= -2 \int_0^{\infty} df \frac{|\mathbf{h}|^2 - 2\operatorname{Re} \{ \mathbf{h}^* \mathbf{s} \} \hat{A} + \hat{A}^2 |\mathbf{s}|^2}{S_n} \\ &= -2 \int_0^{\infty} df \frac{|\mathbf{h}|^2}{S_n} - \frac{1}{2} \hat{A}^2 \left(4 \int_0^{\infty} df \frac{|\mathbf{s}|^2}{S_n} \right) + \hat{A} \left(4 \int_0^{\infty} df \frac{\operatorname{Re} \{ \mathbf{h}^* \mathbf{s} \}}{S_n} \right) \end{aligned}$$

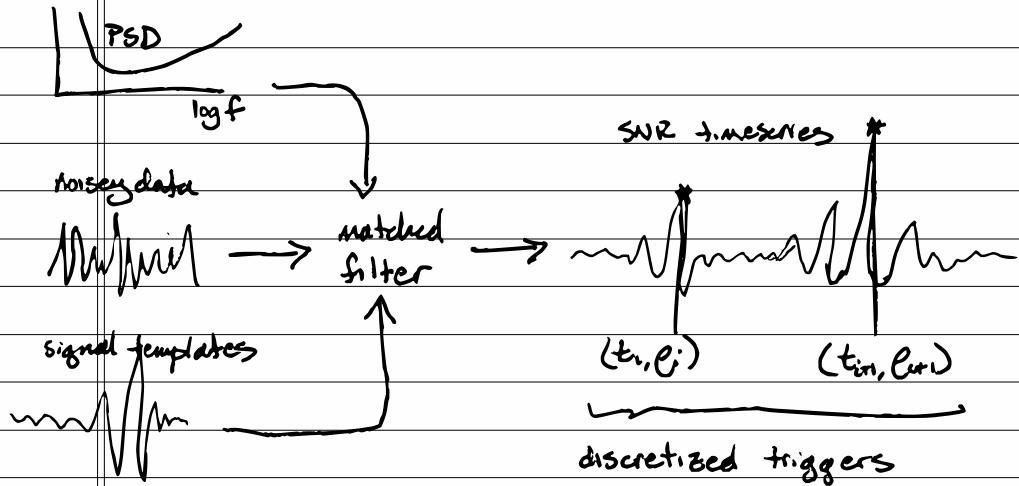
$$\log p(\mathbf{h} | \hat{A}, \omega) = -\frac{1}{2} \left(4 \int_0^{\infty} df \frac{|\mathbf{h}|^2}{S_n} \right) + \frac{1}{2} \hat{A}^2$$

$\underbrace{\quad}_{\text{noise-only likelihood}}$ \uparrow $\boxed{\log L_{\text{noise}} \sim \frac{1}{2} (\text{SNR})^2}$

\hookrightarrow does not depend on signal

EXERCISE: derive the max. likelihood estimator for an unmodulated signal (\mathbf{s}) and the corresponding L_{MLE}

Matched filter searchers work by defining a template bank of known signals \mathcal{A} and computing the SNR for each template as a function of time. Peaks in the SNR timeseries correspond to higher likelihoods that a signal is present. Selecting local maxima in the SNR timeseries produces a sequence of discretized "triggers," which are the candidate events.



Most (all?) matched filter searchers follow this basic approach. Many of the differences come in after triggers have been produced (ie, how one estimates the false alarm prob. or rate from the set of triggers).

Computational Techniques for Bayes Factors

① analytic integration (often not possible)

if you can, compute the marginal likelihood (evidence) analytically. This can often provide insight & be computationally very efficient. Sometimes scaling of B with SNR can be surprising.

EXERCISE: consider a univariate Gaussian variable

$$x \sim N(\mu, 1) = \frac{1}{\sqrt{2\pi}} \exp(-\frac{1}{2}(x-\mu)^2)$$

w) noise model: $\mu=0$

signal model: $\mu \neq 0$

Compute the Bayes factor under signal model where you know the sign & magnitude of μ & when you only know the magnitude of μ (marginalize over the sign).

② Savage-Dickey Density Ratio (SDDR)

for nested models (one hypothesis corresponds to a subset of the alternative hypotheses) so that one model is recovered by a specific set of parameters in the other, we can compute the Bayes factor from the posterior & the prior.

$$p(\theta | h, H_A) = \frac{p(h|\theta) p(\theta | H_A)}{\int_{\Theta} p(h|\theta) p(\theta | H_A)}$$

now, if H_B corresponds to H_A w/ $\theta = \theta_B$, then

$$p(\theta_B | h, H_A) = \frac{p(h|\theta_B) p(\theta_B | H_A)}{\int_{\Theta} p(h|\theta) p(\theta | H_A)}$$

$$\Rightarrow \frac{p(h|\theta_B)}{\int_{\Theta} p(h|\theta) p(\theta | H_A)} = \frac{p(h|H_B)}{p(h|H_A)} = B_A^B = \frac{p(\theta_B | h, H_A)}{p(\theta_B | H_A)}$$

Savage-Dickey Density Ratio

Because it is often easier to draw samples from $p(\theta | h, H)$ than to estimate $p(h | H)$, this can be very efficient.

(3) Thermodynamic Integration (Power Sampling) and Nested Sampling

More general techniques exist to compute B when we cannot marginalize analytically or we do not have nested models. These tend to be (much) more expensive & opaque. Open source libraries exist in several languages, but **BE CAREFUL** and make sure you know what you're computing.

Example: Power Sampling

Consider the integral defined for $\beta \in [0, 1]$

$$Z_\beta = \int d\theta g(\theta) \left(\frac{p(\theta)}{g(\theta)} \right)^\beta$$

$$\text{now, } \frac{\partial \ln Z_\beta}{\partial \beta} = \frac{1}{Z_\beta} \frac{\partial}{\partial \beta} \int d\theta g(\theta)^\beta$$

$$= \int d\theta \left[\frac{g(P/q)^\beta}{Z_\beta} \right] \left(\frac{q}{P} \right)^\beta \frac{\partial}{\partial \beta} \left(\frac{q}{P} \right)^\beta$$

$$= \int d\theta \frac{g(P/q)^\beta}{Z_\beta} \frac{\partial}{\partial \beta} \beta \log(P/q)$$

$$= \int d\theta r(\theta | \beta) \log(p(\theta)/g(\theta))$$

$$\text{where } r(\theta | \beta) = \frac{g(\theta) (p(\theta)/g(\theta))^\beta}{Z_\beta}$$

$$\text{we note that } \int_0^1 d\beta \frac{\partial}{\partial \beta} \ln \frac{Z_\beta}{Z_0} = \ln \int d\theta p(\theta) - \ln \int d\theta g(\theta)$$

and if $p = p(h|\theta) p(s|h)$; $g = p(h|\theta) p(s|B)$ this becomes

$$\int_0^1 d\beta \frac{\partial}{\partial \beta} \ln \frac{Z_\beta}{Z_0} = \ln B^{\frac{1}{\beta}}$$

we can estimate $\frac{\partial}{\partial \beta} \ln \frac{Z_\beta}{Z_0}$ via a Monte Carlo sum

$$\frac{\partial}{\partial \beta} \ln \frac{Z_\beta}{Z_0} \approx \frac{1}{N} \sum_i^N \log \frac{p(\theta_i)}{g(\theta_i)} \quad | \theta_i \sim r(\theta | \beta)$$

we proceed by defining a "temperature ladder"

$$\beta = 1/T \in [0, 1] \Rightarrow T \in [1, \infty)$$

and sampling from $g_\beta(\theta)$ separately for each rung on the ladder. If we consider N_T temperatures, we then obtain N_T sets of samples. A can use each to estimate $\frac{\partial^2}{\partial \beta^2} \ln Z_\beta$ at that temperature.

With the set of Mark Carlo estimates for $\frac{\partial^2}{\partial \beta^2} \ln Z_\beta$, we can then numerically integrate the curve to obtain $\ln B_B^A$.

Note, Thermodynamic Integration is an example of power sampling where

$$p(\theta) = p(h|\theta) p(\theta|H) \quad (\text{signal model})$$

$$q(\theta) = p(\theta|H^\perp) \quad (\text{noise model})$$

and is often used to estimate B_N^S : Bayes factor in favor of the signal model compared to just noise. B_N^S is a powerful detection statistic.