

Statistics in Physics Analysis

Lecture 2

Harrison B. Prosper
Florida State University

Workshop on LHC Physics
TIFR, Mumbai
21 – 27 October, 2009

Outline

- Lecture 2 – **Foundations & Applications**
 - The Bayesian Approach
 - Decisions & Loss
 - Hypothesis Tests
 - Summary

The Bayesian Approach



The Bayesian Approach

A Bayesian calculation requires the following ingredients:

$p(\mathbf{D} | \theta, \varphi)$ the **probability model** that represents the mechanism that gave rise to the observed data \mathbf{D} , given some *unknown* values of the parameters θ, φ .

$p(\theta, \varphi)$ the **prior probability density** over the parameter space of the probability model

The Bayesian Approach

Then one calculates the posterior density as follows:

posterior

marginal
likelihood

prior

$$p(\theta | D) = \frac{p(D | \theta) p(\theta)}{\int p(D | \theta) p(\theta) d\theta}$$

where the **marginal (or integrated) likelihood** is given by

$$p(D | \theta) = \int_{\Phi} p(D | \theta, \phi) p(\phi | \theta) d\phi$$

and $p(\theta, \phi) = p(\phi | \theta) p(\theta)$ is the full prior density.

Example – D0 Top Discovery Data

D0 1995 Top Discovery Data

$$n = 17 \text{ events}$$

$$b_0 = 3.8 \pm 0.6 \text{ events}$$

Example – D0 Top Discovery Data

Likelihood Functions

$$\begin{aligned} p(n|b+s, H_1) &= \text{Poisson}(n|b+s) = \exp[-(b+s)] (b+s)^n / n! \\ p(n|b, H_0) &= \text{Poisson}(n|b) = \exp[-b] b^n / n! \end{aligned}$$

Prior Density

$$\begin{aligned} p(b, s) &= p(b|s) p(s) \\ p(b|s) &= \text{Gamma}(kb|B+1) = k \exp(-kb) (kb)^B / \Gamma(B+1) \end{aligned}$$

where the effective scale factor k and count B are

$$\begin{aligned} b_0 &= B / k & B &= (b_0 / \delta b)^2 = (3.8 / 0.6)^2 = 41.11 \\ \delta b &= \sqrt{B} / k & k &= b_0 / \delta b^2 = 3.8 / 0.6^2 = 10.56 \end{aligned}$$

Example – Integrated Likelihoods

The **integrated likelihoods** are

$$\begin{aligned} p(n | s, H_1) &= \int_0^{\infty} \text{Poisson}(n | b + s) \text{Gamma}(kb | B + 1) db \\ &= \left(\frac{k}{1+k} \right)^{B+1} \sum_{r=0}^n \frac{1}{(1+k)^r} \frac{\Gamma(B+1+r)}{\Gamma(B+1)r!} \text{Poisson}(n-r | s) \end{aligned}$$

and

$$p(n | H_0) = p(n | s = 0, H_1) = \left(\frac{k}{1+k} \right)^{B+1} \frac{1}{(1+k)^n} \frac{\Gamma(B+1+n)}{\Gamma(B+1)n!}$$

Exercise 1: Compute these integrated likelihoods

Example – Posterior Density

Given the integrated likelihood

$$p(n | s, H_1) = \left(\frac{k}{1+k} \right)^{B+1} \sum_{r=0}^n c_r(k, B) \text{Poisson}(n-r | s)$$

where

$$c_r(k, B) \equiv \frac{1}{(1+k)^r} \frac{\Gamma(B+1+r)}{\Gamma(B+1)r!}$$

we can compute

$$p(s | n, H_1) = \frac{p(n | s, H_1) p(s | H_1)}{\int_0^{\infty} p(n | s, H_1) p(s | H_1) ds}$$

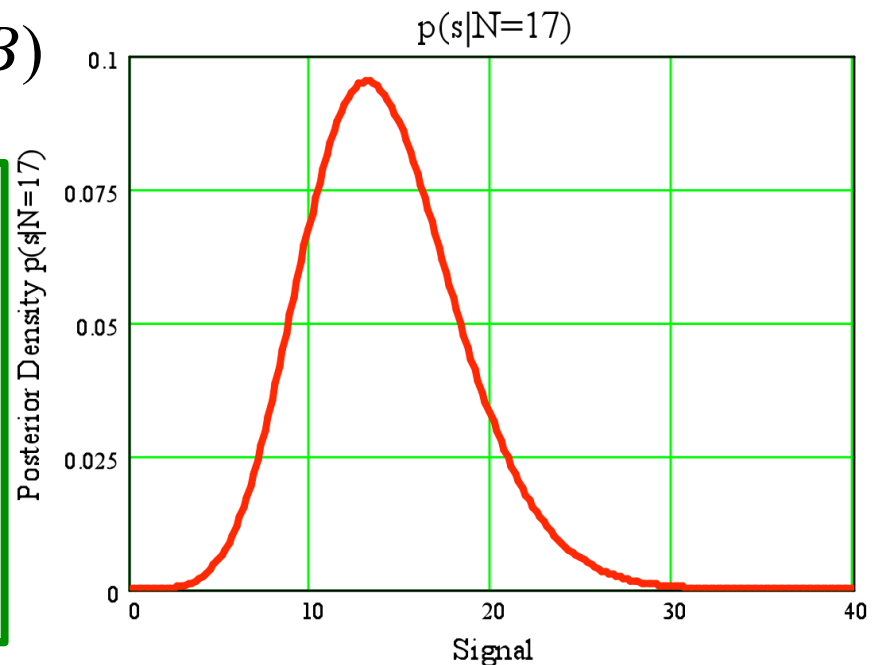
Example – Posterior Density

Assuming a *flat prior* for the signal $p(s|H_1) = \text{constant}$, the posterior density is given by

$$p(s | n, H_1) = \frac{\sum_{r=0}^n c_r(k, B) \text{Poisson}(n - r | s)}{\sum_{r=0}^n c_r(k, B)}$$

Exercise 2: Compute $p(s|n, H_1)$.

Repeat assuming the prior density $p(s|H_1) = \text{Gamma}(qs|S+1)$, where $S = (s_0 / \delta s)^2$ and $q = s_0 / \delta s^2$



Decisions & Loss



Decisions and Loss

The posterior density $p(\theta|D)$ is the *complete* answer to an inference about the parameter θ .

However, it is often of interest to summarize this answer with a **point estimate** θ^* (a measurement) and, or, an **interval estimate** $[\theta_L, \theta_U]$.

Or, we wish to decide which of two or more competing models is preferred by the data.

Decision theory provides a general way to model such problems.

Decisions and Loss

One way to render a decision about the value of θ is to implement the decision as a function d that returns an **estimate** θ^* of θ . A function d that returns estimates is called an **estimator**.

In principle, we also need to specify a **loss function** $L(d, \theta)$ that quantifies what we lose should the estimate turn out to have been a bad one.

Decisions and Loss

In practice, since our knowledge of the parameter θ is encoded in the posterior density $p(\theta | D)$, our decisions will be more *robust* if we average ($E[*]$) the loss $L(d, \theta)$ with respect to $p(\theta | D)$

$$\begin{aligned} R(d) &= E[L(d, \theta)] \\ &= \int L(d, \theta) p(\theta | D) d\theta \end{aligned}$$

The quantity $R(d)$ is called the **risk function**.

By definition, the **optimal estimate** of θ is the one that minimizes the risk

$$\theta^* = \arg \min_d R(d)$$

Comments

In general, different loss functions will yield different estimates.

Therefore, even with *exactly the same data* one should not be surprised to obtain different results.

Reasonable people can disagree about the results simply because they disagree about what properties of the results are thought to be most useful.

For example, many insist that a result should always be *unbiased*, while others do not!

Comments

Consider a loss function $L(d, m)$ to extract a value for the Higgs mass, m , from a posterior density $p(m|D)$.

Suppose $L(d, m)$ is invariant in the following sense: it yields an estimate m^* of m which, when inserted into the prediction $\sigma = g(m)$ for the Higgs cross section, yields an estimate of the cross section $\sigma^* = g(m^*)$ that is *identical* to the one obtained using the loss function $L(d, \sigma)$.

$L(d, \sigma)$ is the loss function $L(d, m)$ with m replaced by σ .

In general, either m^* or σ^* (or both) will be **biased**.

Comments

To see this, expand $\sigma^* = g(m^*)$ about the *true* Higgs mass m

$$\sigma^* \approx g(m) + (m^* - m) g' + \frac{1}{2} (m^* - m)^2 g''$$

and average both sides over an **ensemble** of estimates. This gives

$$E[\sigma^*] \approx \sigma + \mathbf{bias} g' + \frac{1}{2} \mathbf{mse} g'',$$

$$E[\sigma^*] \approx \sigma + \mathbf{bias} g' + \frac{1}{2} [\mathbf{bias}^2 + \mathbf{variance}] g'',$$

where $\mathbf{bias} = E[m^*] - m$ and $\mathbf{variance} = E[m^{*2}] - E[m^*]^2$.

\mathbf{mse} : mean square error (note: $\mathbf{rms} = \sqrt{\mathbf{mse}}$)

Decisions and Loss

Point Estimation

quadratic loss

$$L(d, \theta) = (d - \theta)^2$$

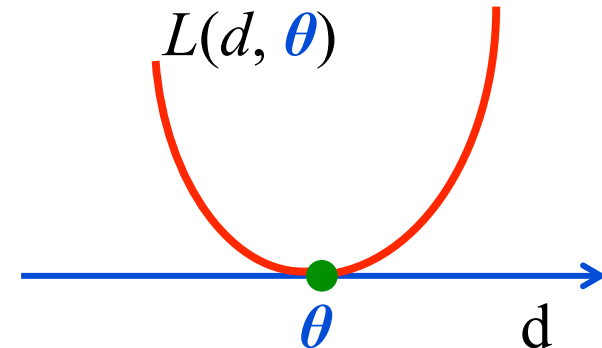
Average with respect to $p(\theta|D)$

$$\begin{aligned} \text{risk } R(d) &= E[(d - \theta)^2] \\ &= E[d^2] - 2E[\theta d] + E[\theta^2] \\ &= d^2 - 2E[\theta]d + E[\theta^2] \end{aligned}$$

minimize with respect to d

$$dR/dd = 2d - 2E[\theta] = 0$$

$$\text{obtaining, } \theta^* = E[\theta]$$



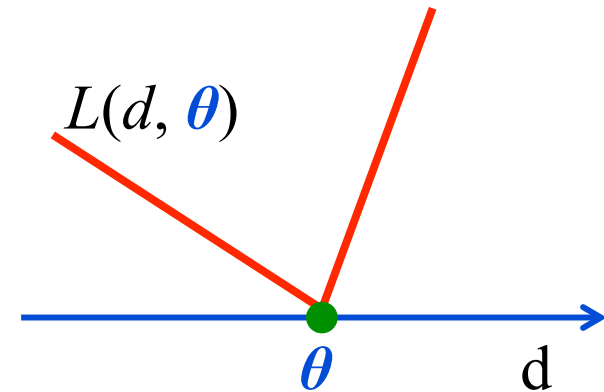
Note: quadratic loss is *not invariant*. If $a = g(\theta)$, then $L(d, a) = (d - a)^2$ gives $a^* = E[a] \neq g(\theta^*)$

Decisions and Loss

Point Estimation

bilinear loss

$$L(d, \theta) = \begin{cases} a(\theta - d), & d < \theta \\ b(d - \theta), & d \geq \theta \end{cases}$$



$$\begin{aligned} \text{risk } R(d) &= a \int H(\theta - d)(\theta - d)p(\theta | D)d\theta \\ &+ b \int H(d - \theta)(d - \theta)p(\theta | D)d\theta \end{aligned}$$

$$H(x) = 1 \text{ if } x > 0 \text{ else } 0$$

Decisions and Loss

Point Estimation

bilinear loss

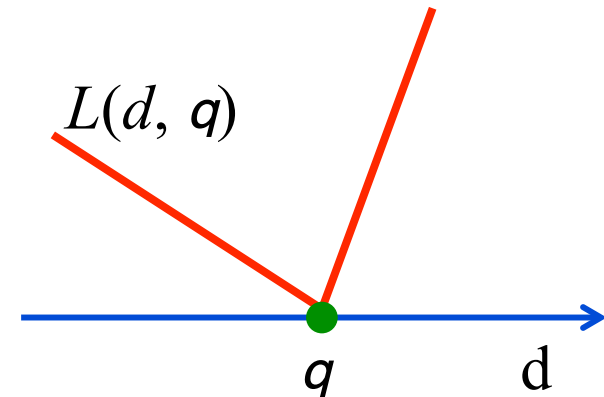
The optimal estimate is

$$\theta^* = \arg \min_d R(d)$$

where θ^* is the $a/(a+b)$ **quantile**

$$\int_{-\infty}^{\theta^*} p(\theta | D) d\theta = a/(a+b)$$

of $p(\theta | D)$. If we set $a = b$, $\theta^* =$ **median** of $p(\theta | D)$



Note: estimates based on quantiles are *invariant*.

Decisions and Loss

Point Estimation

zero-one loss

$$L(d, \theta) = \begin{cases} 0, & |d - \theta| \leq b \\ 1, & |d - \theta| > b \end{cases}$$

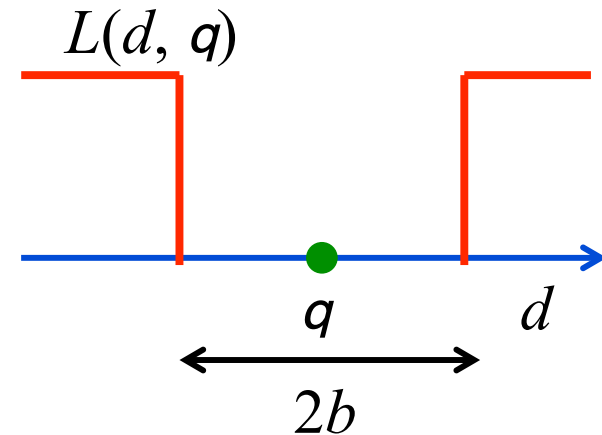
Its risk function is

$$R(d) = \int [H(\theta - b - d) + H(d - \theta - b)] p(\theta | D) d\theta$$

and the optimal estimate $\theta^* = \min_d R(d)$ is the solution of

$$p(\theta^* + b | D) = p(\theta^* - b | D).$$

In the limit $b \rightarrow 0$, one obtains $\theta^* = \mathbf{mode}$ of $p(\theta | D)$. The mode is *not invariant*.



Example – Posterior Mean

Compute the moments of $p(s|n, H_1)$ about zero

$$M_m = \int_0^{\infty} s^m p(s | n, H_1) ds$$
$$= \sum_{r=0}^n c_r(k, B) (n - r + m)! / (n - r)! / \sum_{r=0}^n c_r(k, B)$$

$p(s|N=17)$

For the D0 top quark discovery data we find:

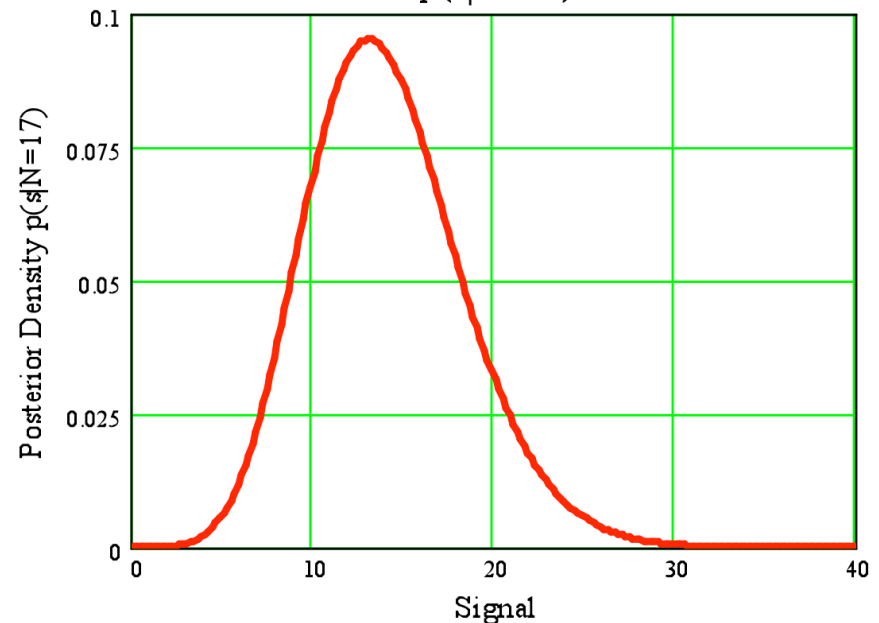
mean

$$M_1 = 14.0 \text{ events}$$

standard deviation

$$\sqrt{(M^2 - M_1^2)} = 4.3 \text{ events}$$

Exercise 3: Compute M_m .

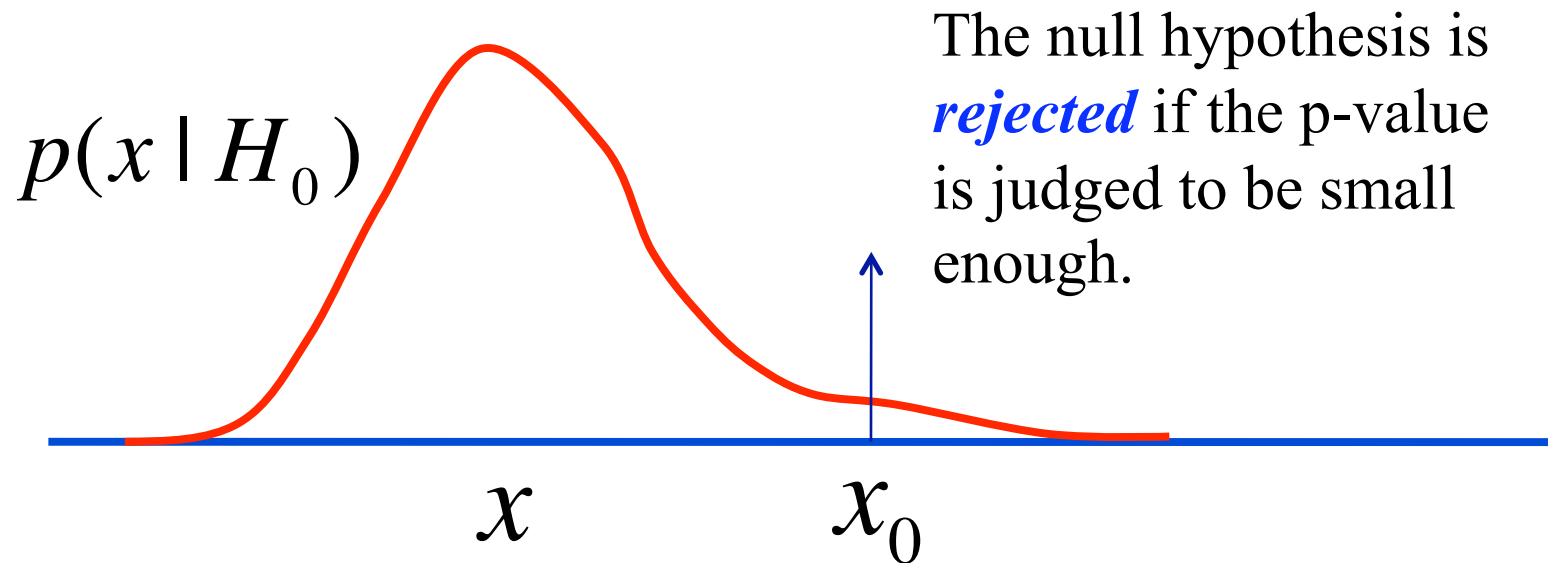


Hypothesis Testing



P-Values

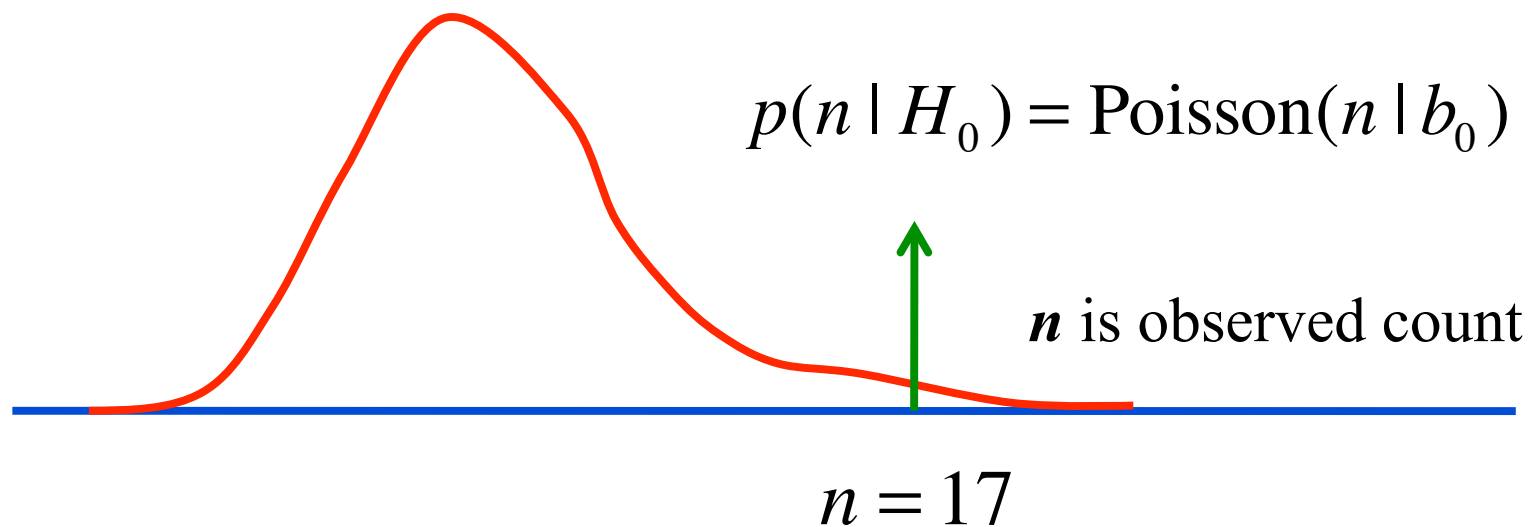
Null hypothesis (H_0): **background-only**



$$\text{p-value} \equiv \int_{x_0}^{\infty} p(x | H_0) dx$$

Example – Top Discovery p-value (a)

Background, $b_0 = 3.8$ events (*ignoring uncertainty*)

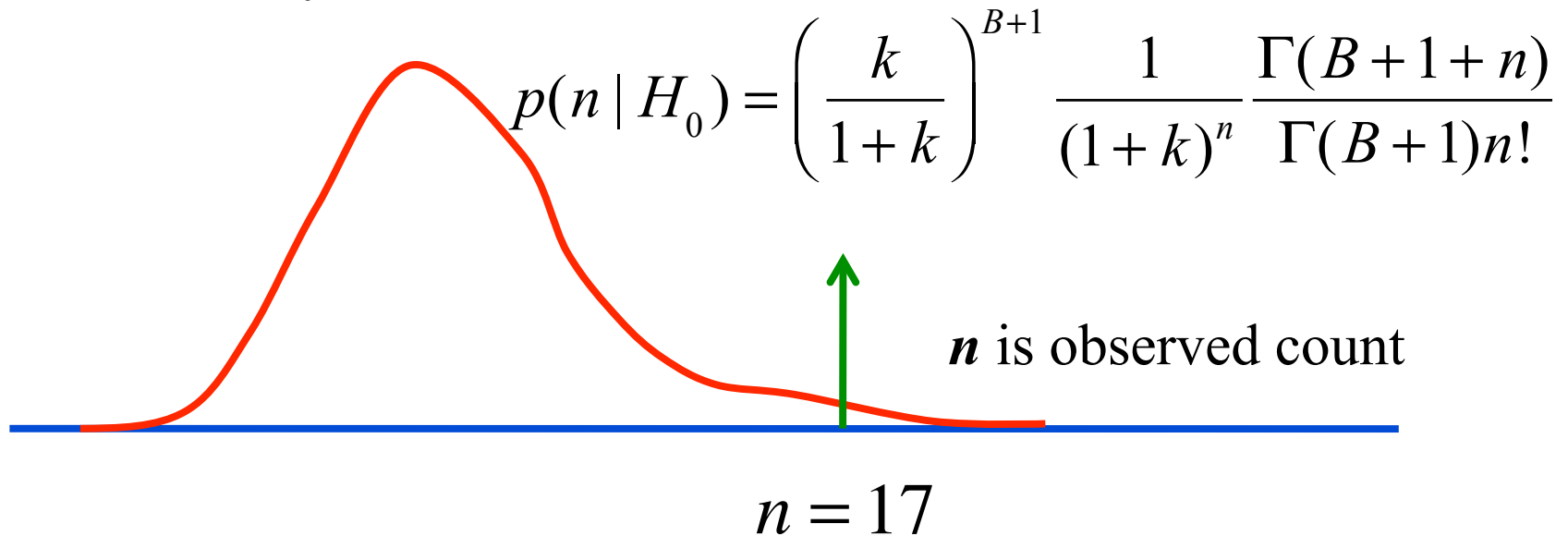


$$\text{p-value} = \sum_{n=17}^{\infty} \text{Poisson}(n | 3.8) = 5.7 \times 10^{-7}$$

This is equivalent to **4.9 σ**

Example – Top Discovery p-value (b)

Background, $b_0 = 3.8 \pm 0.6$ events



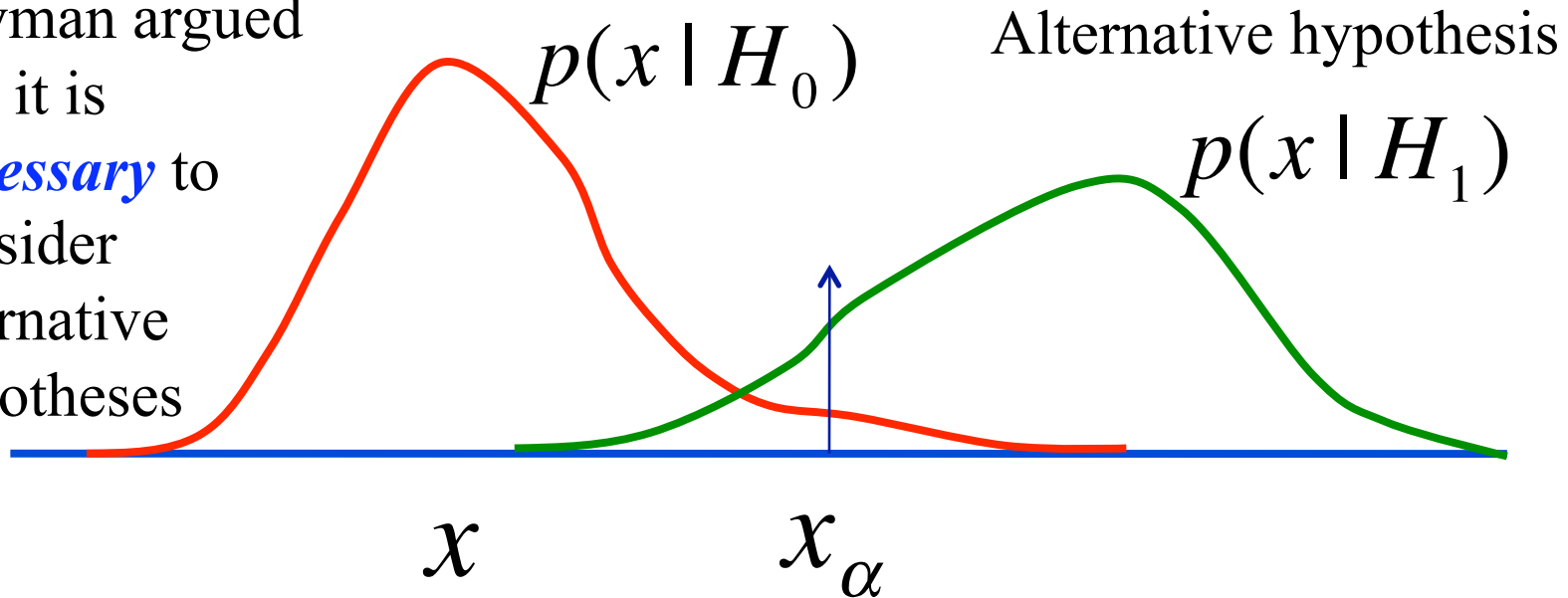
$$\text{p-value} = \sum_{n=17}^{\infty} p(n | H_0) = 5.4 \times 10^{-6}$$

This is equivalent to **4.4 σ**

The Neyman-Pearson Test

Neyman argued that it is *necessary* to consider alternative hypotheses

H_1

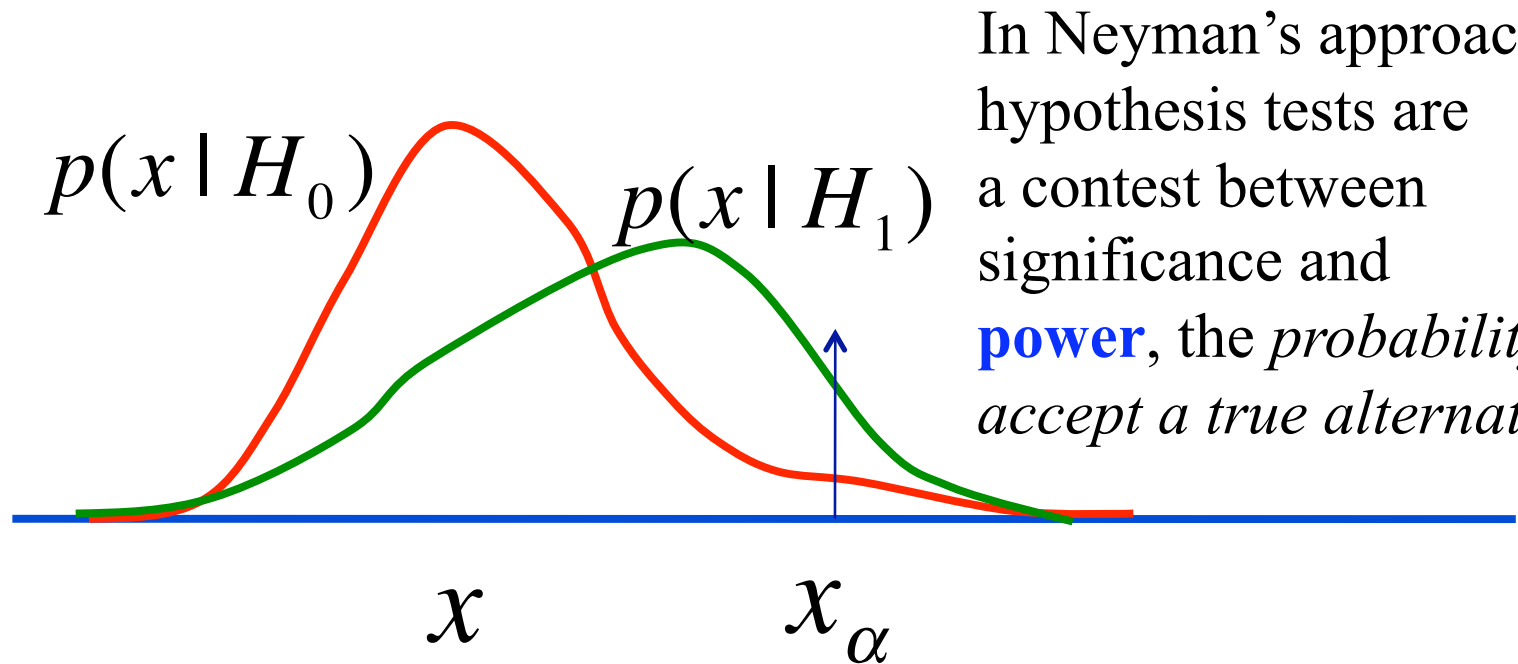


$$\alpha = \int_{x_\alpha}^{\infty} p(x | H_0) dx$$

significance of test

A *fixed* significance, the *probability to reject a true null*, is chosen before data are analyzed.

The Neyman-Pearson Test



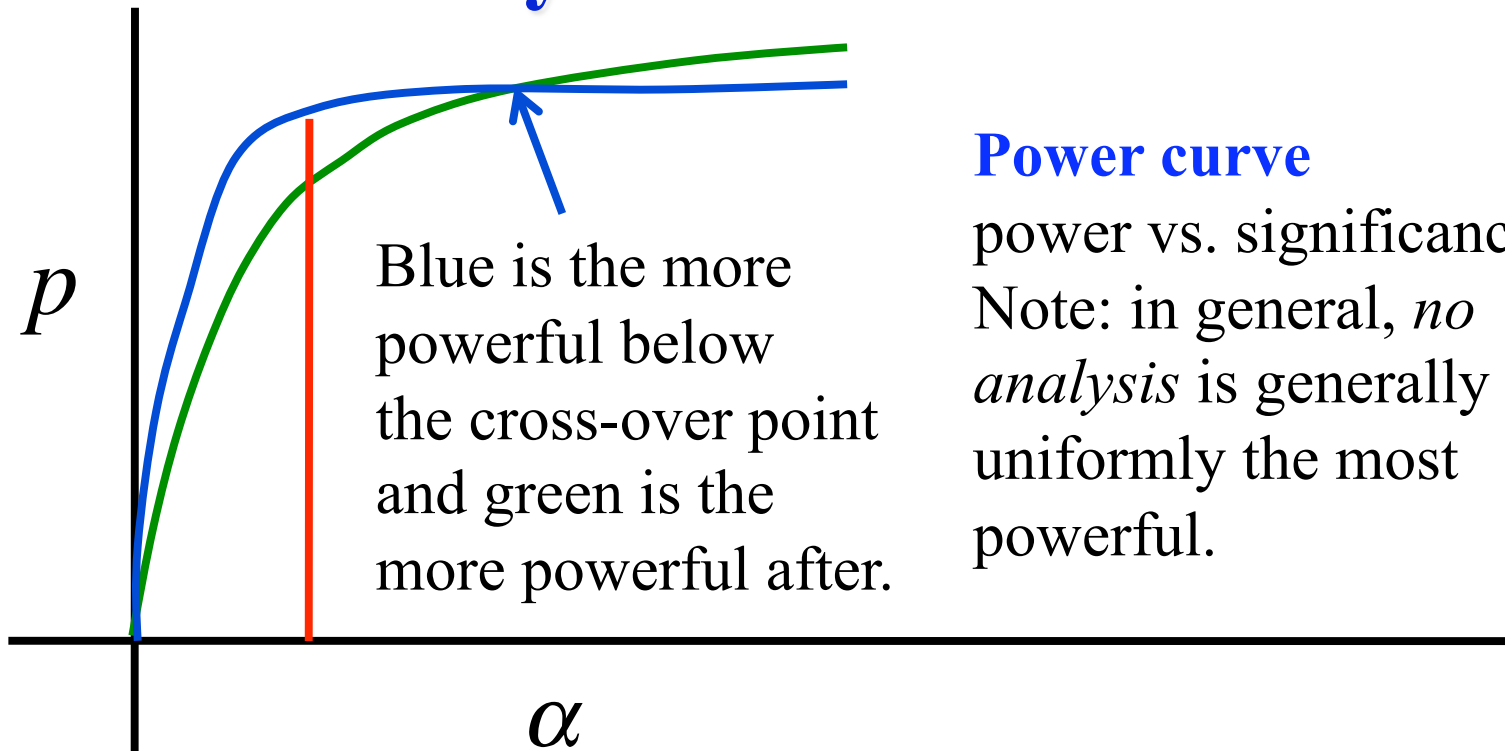
$$\alpha = \int_{x_\alpha}^{\infty} p(x | H_0) dx$$

significance of test

$$p = \int_{x_\alpha}^{\infty} p(x | H_1) dx$$

power

The Neyman-Pearson Test



$$\alpha = \int_{x_\alpha}^{\infty} p(x | H_0) dx$$

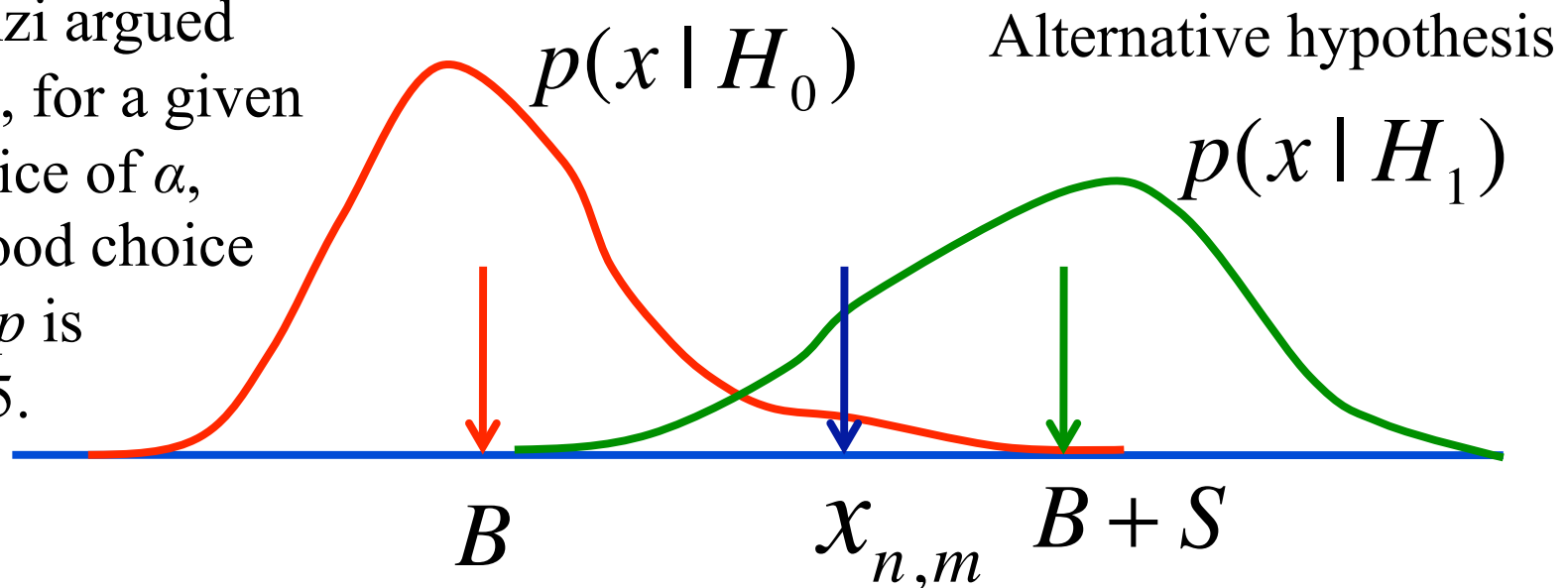
significance of test

$$p = \int_{x_\alpha}^{\infty} p(x | H_1) dx$$

power

Punzi's Test

Punzi argued that, for a given choice of α , a good choice for p is 0.95.



$$S = n\sqrt{B} + m\sqrt{B + S}$$

Exercise 4: Write as $Q = S/\sqrt{(B+a)}$ and find a

Summary

- **Decision Theory**

- The basic insight is that optimal decision making entails combining a *loss function* with a *posterior density*. Since loss functions can differ, it is unsurprising that results can differ even when using the same data.

- **Hypothesis Tests**

- The standard non-Bayesian approach is that of Neyman and Pearson, plus the calculation of p-values. Neyman argued (in agreement with Bayesians) that it is necessary to consider pairs of hypotheses.