

# What have we measured?

Daniel Greenwald  
Technische Universität München

December 1, 2022

Let's look at an analysis of

$$\tau^- \rightarrow XY$$

in which we

- identify *all* sources of events that mimic the signal (backgrounds) and
- determine criteria that
  - select for signal events and
  - reject background ones.

Explicitly, for each source  $i$ , we **observe** that

- for  $M_i$  simulated events,
- $m_i$  events pass the selection criteria.

We then look at the real data and **observe** that

- for data of a known size,
- $N$  events pass the same selection criteria.

What's the goal?

To measure the branching fraction

$$\mathcal{B}_s \equiv \frac{\Gamma(\tau^- \rightarrow XY)}{\Gamma(\tau^- \rightarrow \text{anything})}$$

What have we observed?

- $M_i \rightarrow m_i$  in simulation
- $N$  in data (of known size)

How do we connect the two?

A model ... a **physics** model.

Let's think like physicists first and then like statisticians.

## From collision to observation:

For signal,

$$e^+e^- \rightarrow \tau^+\tau^-$$

$\hookrightarrow \tau^\pm \rightarrow 1 \text{ prong}$  occurs  $\mathcal{B}_t$  fraction of the time (tag)

$\hookrightarrow \tau^\mp \rightarrow XY$  occurs  $\mathcal{B}_s$  fraction of the time (signal)

- $2\mathcal{B}_s\mathcal{B}_t$  of  $\tau\tau$  events are the type of event we reconstruct.
- And we have an efficiency  $\epsilon_s$  to detect them.

So we expect to observe  $2\epsilon_s\mathcal{B}_s\mathcal{B}_t$  of  $\tau\tau$  events.

How many  $\tau\tau$  events are there?

$$\sigma_{\tau\tau} L$$

- $\sigma_{\tau\tau} \equiv$  cross section for  $e^+e^- \rightarrow \tau^+\tau^-$
- $L \equiv$  integrated luminosity of the data set

So we expect to measure  $2\epsilon_s\mathcal{B}_s\mathcal{B}_t\sigma_{\tau\tau}L$  signal events

## From collision to observation:

For each background,

$e^+e^- \rightarrow (\text{background source})_i$

- with cross section  $\sigma_i$
- and efficiency  $\epsilon_i$

So we expect to measure  $\sum_i \epsilon_i \sigma_i L$  background events

In total, we expect

$$\left[ 2\epsilon_s \mathcal{B}_s \mathcal{B}_t \sigma_{\tau\tau} + \sum_i \epsilon_i \sigma_i \right] L \text{ events}$$

## From simulation to observation

We simulated  $M_i$  events, and observe  $m_i$  events after selection criteria.

This informs us about the **efficiency**.

Naively:

$$\epsilon_i = \frac{m_i}{M_i} \pm \frac{m_i}{M_i} \sqrt{\frac{1}{m_i} + \frac{1}{M_i}}$$

Naive ... and wrong ...  $X \pm Y$  implies a normal distribution,

$$P(\epsilon_i | M_i, m_i) = \mathcal{N}\left(\frac{m_i}{M_i}, \frac{m_i}{M_i} \sqrt{\frac{1}{m_i} + \frac{1}{M_i}}\right)$$

which is not the case.

And what happens when  $m_i = 0$ :

$$\epsilon_i(m_i = 0) = 0 \text{ exactly.}$$

This states we're 100% confident the **efficiency is 0** ... regardless of  $M_i$ 's value.

This is also not a definition of an **efficiency**.  
It's just a ratio of two numbers.

We need a statistical model.

Let's try again using Bayesian statistics.

An ultra quick refresher:

$$P(A \cap B) = P(B \cap A)$$

$$P(A|B) P(B) = P(B|A) P(A)$$

$$P(A|B) = \frac{P(B|A) P(A)}{P(B)}$$

We relate the conditional probability of A given B (the left side)  
to the conditional probability of B given A (the right side)

$$P(A|B) = \frac{P(B|A) P(A)}{P(B)}$$

Let's look at this more concretely in terms of data and parameters:

We want to know the probability of parameters,  $\vec{\lambda}$ , given measured data

$$P(\vec{\lambda}|\text{data}) = \frac{P(\text{data}|\vec{\lambda}) P(\vec{\lambda})}{P(\text{data})}$$

So we use Bayes' theorem.

These four parts have common names (and symbols):

$P(\vec{\lambda} \text{data})$	$\equiv$	<i>posterior</i> prob. distribution
$P(\text{data} \vec{\lambda}) \equiv \mathcal{L}(\text{data} \vec{\lambda})$	$\equiv$	likelihood of the data
$P(\vec{\lambda}) \equiv P_0(\vec{\lambda})$	$\equiv$	<i>a priori</i> prob. distribution
$P(\text{data}) \equiv Z(\text{data}) = \int \mathcal{L}(\text{data} \vec{\lambda}) P_0(\vec{\lambda}) d\vec{\lambda}$	$\equiv$	evidence



$$P(\vec{\lambda}|\text{data}) \propto \mathcal{L}(\text{data}|\vec{\lambda}) P_0(\vec{\lambda})$$

Since we can normalize the posterior from the above proportionality, we don't need to

- concern ourselves with the evidence,
- normalize the likelihood, or
- normalize the prior

$$P(\vec{\lambda}|\text{data}) \propto \mathcal{L}(\text{data}|\vec{\lambda}) P_0(\vec{\lambda})$$

Let's look now at our efficiency,  $\epsilon_i$ :

$$P(\epsilon_i|M_i, m_i) \propto \mathcal{L}(m_i|M_i, \epsilon_i) P_0(\epsilon_i)$$

What is the likelihood for selecting  $m_i$  events from a total of  $M_i$  when the efficiency for selecting each event is  $\epsilon_i$ ?

$$\mathcal{L}(m_i|M_i, \epsilon_i) \propto \epsilon_i^{m_i} (1 - \epsilon_i)^{M_i - m_i}$$

⇒ the Binomial distribution:

What about  $P_0(\epsilon_i)$ ? A general approach that parameterizes two common choices:

$$P_0(\epsilon_i|\delta_i) = \begin{cases} (\epsilon_i)^{-\delta_i} (1 - \epsilon_i)^{-\delta_i}, & \text{if } \epsilon_i \in [0, 1], \\ 0, & \text{else} \end{cases}$$

$\delta_i = 0 \Rightarrow P_0(\epsilon_i) = \text{constant prior}$

$\delta_i = 1 \Rightarrow P_0(\epsilon_i) = \text{double-inverse prior}$

Putting it all together

$$P(\epsilon_i | M_i, m_i) \propto \epsilon_i^{m_i - \delta_i} (1 - \epsilon_i)^{M_i - m_i - \delta_i}$$

if  $\epsilon_i \in [0, 1]$ , else it's zero.

In terms of  $\epsilon_i$ , this is the Beta distribution ... with properties:

mean  $\bar{\epsilon}_i = \frac{m_i + 1 - \delta_i}{M_i + 2 - 2\delta_i}$

mode  $\epsilon_i^* = \frac{m_i - \delta_i}{M_i - 2\delta_i}$

variance  $\text{Var}[\epsilon_i] = \frac{2}{\bar{\epsilon}_i} \frac{M_i - m_i + 1 - \delta_i}{(m_i + 1 - \delta_i)(M_i + 3 - 2\delta_i)}$   
 $\frac{\text{Var}[\epsilon_i]}{\bar{\epsilon}_i^2} = \frac{1}{m_i + 1 - \delta_i} - \frac{1}{M_i + 3 - 2\delta_i} - \frac{1}{(m_i + 1 - \delta_i)(M_i + 3 - 2\delta_i)}$

skew = ...

kurtosis = ...

Let's make the full model now:

$$P(\mathcal{B}_s, \mathcal{B}_t, \vec{\epsilon}, \vec{\sigma}, L | N, \vec{M}, \vec{m}) \propto \mathcal{L}(N, \vec{M}, \vec{m} | \mathcal{B}_s, \mathcal{B}_t, \vec{\epsilon}, \vec{\sigma}, L) P_0(\mathcal{B}_s, \mathcal{B}_t, \vec{\epsilon}, \vec{\sigma}, L)$$

Let's focus on the likelihood first. It factorizes:

$$\mathcal{L}(N, \vec{M}, \vec{m} | \mathcal{B}_s, \mathcal{B}_t, \vec{\epsilon}, \vec{\sigma}, L) = \mathcal{L}(N | \mathcal{B}_s, \mathcal{B}_t, \vec{\epsilon}, \vec{\sigma}, L) \prod_i \mathcal{L}(m_i | M_i, \epsilon_i)$$

And we already know  $\mathcal{L}(m_i | M_i, \epsilon_i)$ .

The likelihood to observe  $N$  events, given an expectation is

$$\mathcal{L}(N | \mathcal{B}_s, \mathcal{B}_t, \vec{\epsilon}, \vec{\sigma}, L) \propto \nu^N e^{-\nu}$$

$\Rightarrow$  the Poisson distribution ... with

$$\nu \equiv \left[ 2\epsilon_s \mathcal{B}_s \mathcal{B}_t \sigma_{\tau\tau} + \sum_i \epsilon_b \sigma_b \right] L$$

## What about the prior?

It factorizes:

$$P_0(\mathcal{B}_s, \mathcal{B}_t, \vec{\epsilon}, \vec{\sigma}, L) = P_0(\mathcal{B}_s) P_0(\mathcal{B}_t) P_0(L) \prod_i P_0(\epsilon_i) P_0(\sigma_i)$$

We've already discussed  $P_0(\epsilon_i)$ .

What about the others?

- $P_0(\mathcal{B}_t) = \mathcal{N}(\mathcal{B}_t | \text{mean, variance})$  [PDG. . . ]
- $P_0(L) = \mathcal{N}(L | \text{mean, variance})$  [Belle II's measured value]
- $P_0(\sigma_i) = \mathcal{N}(\sigma_i | \text{mean, variance})$ 
  - measured value @  $\sqrt{s} = 10.58$  GeV
  - calculated value @  $\sqrt{s} = 10.58$  GeV

If necessary,  $\sqrt{s}$  could itself be a parameter.

And what about  $P_0(\mathcal{B}_s)$ ?

This is our parameter of interest . . . we have several choices.

## $P_0(\mathcal{B}_s)$ Choices

- $P_0(\mathcal{B}_s) = \text{flat in } [0,1] \rightarrow \text{flat in the value of } \mathcal{B}_s$
- $P_0(\mathcal{B}_s) \propto 1/\mathcal{B}_s \rightarrow \text{flat in the scale of } \mathcal{B}_s$
- $P_0(\mathcal{B}_s) = \text{previous measurement result} \rightarrow \text{most informative}$

This may lead to the lowest upper bound,  
but requires work to figure out the proper likelihood.

Ordinarily we don't do this, so that our measurement can be used with others—cf. the PDG—but since experiments generally don't report *useful* results on upper limits, the PDG doesn't average them anyway!

The results should not vary greatly (at least between the first two choices).

Such a check is easy—and should be done.

## Putting this aaaaall together ...

$$P(\mathcal{B}_s, \mathcal{B}_t, \vec{\epsilon}, \vec{\sigma}, L | N, \vec{M}, \vec{m}) \propto \mathcal{L}(N, \vec{M}, \vec{m} | \mathcal{B}_s, \mathcal{B}_t, \vec{\epsilon}, \vec{\sigma}, L) P_0(\mathcal{B}_s, \mathcal{B}_t, \vec{\epsilon}, \vec{\sigma}, L)$$

$$\mathcal{L}(N, \vec{M}, \vec{m} | \mathcal{B}_s, \mathcal{B}_t, \vec{\epsilon}, \vec{\sigma}, L) = \mathcal{L}(N | \mathcal{B}_s, \mathcal{B}_t, \vec{\epsilon}, \vec{\sigma}, L) \prod_i \mathcal{L}(m_i | M_i, \epsilon_i)$$

$$\mathcal{L}(N | \mathcal{B}_s, \mathcal{B}_t, \epsilon_s, \sigma_{\tau\tau}, L) = \text{Poisson}(N | 2\epsilon_s \mathcal{B}_s \mathcal{B}_t L + \sum_b \epsilon_b \sigma_b L)$$

$$\mathcal{L}(m_i | M_i, \epsilon_i) = \text{Binom}(m_i | M_i, \epsilon_i)$$

$$P_0(\mathcal{B}_s, \mathcal{B}_t, \vec{\epsilon}, \vec{\sigma}, L) = P_0(\mathcal{B}_s) P_0(\mathcal{B}_t) P_0(L) \prod_i P_0(\epsilon_i) P_0(\sigma_i)$$

$$P_0(\mathcal{B}_t) = \text{world average}$$

$$P_0(L) = \text{Belle II measurement from other channels}$$

$$P_0(\epsilon_i) = \epsilon_i^{-\delta_i} (1 - \epsilon_i)^{-\delta_i}$$

$$P_0(\sigma_i) = \text{measurement or calculation}$$

$$P_0(\mathcal{B}_s) = \mathcal{B}_s^{-\delta_{\mathcal{B}}}$$

a semantic matter:

$$\mathcal{L}(m_i | M_i, \epsilon_i) \times P_0(\epsilon_i) = \text{Binom}(m_i | M_i, \epsilon_i) \times P_0(\epsilon_i)$$

can instead be taken as a prior on  $\epsilon_i$ ,  
with no likelihood component for it then:

$$P_0(\epsilon_i) = \text{Beta}(\epsilon_i | m_i + 1 - \delta_i, M_i - m_i + 1 - \delta_i)$$



So we have a model:

$$P(\mathcal{B}_s, \mathcal{B}_t, \vec{\epsilon}, \vec{\sigma}, L | N, \vec{M}, \vec{m}) = \dots$$

But we don't care about  $\mathcal{B}_t, \vec{\epsilon}, \vec{\sigma}, L \dots$  they're *nuisance* parameters.

So we integrate them out:

$$P(\mathcal{B}_s | N, \vec{M}, \vec{m}) = \int P(\mathcal{B}_s, \mathcal{B}_t, \vec{\epsilon}, \vec{\sigma}, L | N, \vec{M}, \vec{m}) d\mathcal{B}_t d\vec{\epsilon} d\vec{\sigma} dL$$

This is called the “marginalization” of  $\mathcal{B}_s$ .

The integral is difficult (perhaps impossible?) to analytically calculate.

But it's very easy to calculate on the computer via sampling:

- Sample from the full posterior.
- Histogram of  $\mathcal{B}_s$  = marginalization of  $\mathcal{B}_s$

Then we have the probability distribution for  $\mathcal{B}_s$  and upper limits are easy:

$$\int_0^{\mathcal{B}_s^{\alpha UL}} P(\mathcal{B}_s | N, \vec{M}, \vec{m}) = \alpha$$

so let's do it ...

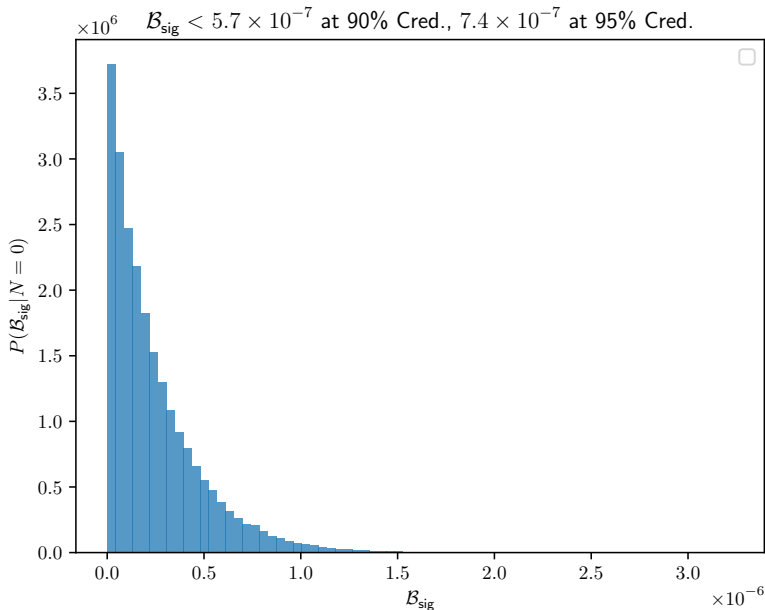
Let's say we have the following information

BG Source	$M [10^6]$	$m$	$\sigma$ [nb]
$q\bar{q}$	372	0	$3.720\,00 \pm 0.003\,72$
$B\bar{B}$	105	0	$1.100 \pm 0.011$
$ee\gamma$	2940	0	$294 \pm 2$
$\mu\mu\gamma$	575	0	$1.148 \pm 0.005$
$ee\mu\mu$	1890	0	$18.97 \pm 0.02$
eeee	397	0	$39.74 \pm 0.02$
$\tau\tau$ (bg)	91.9	0	$0.919 \pm 0.003$

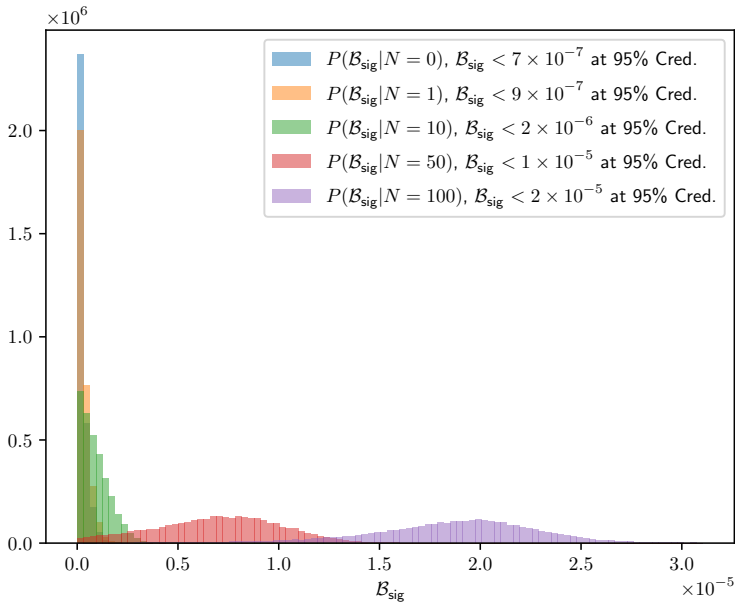
$$\tau^- \rightarrow XY \text{ (sig)} \quad 1 \quad 53\,489$$

$$\mathcal{B}_t = (84.71 \pm 0.06) \times 10^{-2} \quad L = (100.0 \pm 1.4) \text{ fb}^{-1}$$

$N = 0$ , all  $\delta = 0$  (flat priors)



## Let's look at different measurement possibilities:



an interesting aside:

Let's look at a Poisson distribution with two sources:

$$P(\nu, \beta | N) \propto (\nu + \beta)^N e^{-(\nu + \beta)} P_0(\nu) P_0(\beta)$$

The marginalized distribution for  $\nu$  is

$$P(\nu | N) \propto e^{-\nu} P_0(\nu) \int (\nu + \beta)^N e^{-\beta} P_0(\beta) d\beta$$

And when  $N = 0$ :

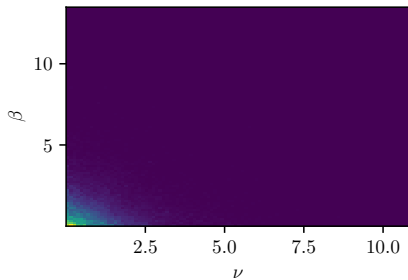
$$P(\nu | N = 0) \propto e^{-\nu} P_0(\nu) \int e^{-\beta} P_0(\beta) d\beta$$

That is,

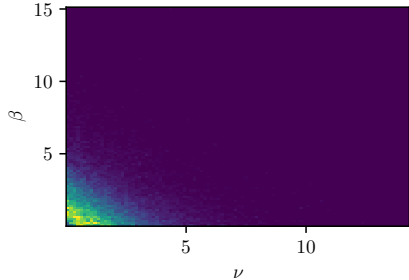
$$P(\nu | N = 0) \propto e^{-\nu} P_0(\nu)$$

$\Rightarrow P(\nu | 0)$  is independent of prior knowledge on  $\beta$ !

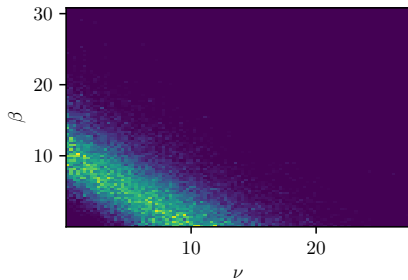
$N = 0$ , correlation = 1%



$N = 1$ , correlation = -16%



$N = 10$ , correlation = -63%



$N = 100$ , correlation = -94%

