

# Statistics for Physicists

J. Hardin

IceCube Summer School 6/7/2023

# Intro

- Hi, I'm John Hardin - I'm a postdoc in Janet Conrad's group, and I lead the MEOWS working group
- I've been asked to talk about statistics - I've tried to avoid overlap with yesterday, but please treat any overlap as a refresher.
- I'm going to start pretty basic, I'm sorry to those of you that have seen some or all of this. I hope it's still useful.
- Please feel free to interrupt and ask questions - it works better if you do.

# Outline

- 1) Probability (What questions are we asking?)
- 2) The Loglikelihood and  $\chi^2$  (How do we answer them?)
- 3) Test statistics and confidence regions (What does that answer mean?)
- 4) Real Experimental Considerations (Systematics, etc)
- 5) Physics statistics jargon (“Brazil Plots”, “Feldman Cousins”, “Look Elsewhere”, “Sigma”)

Probability

Or

“What Questions are We Asking?”

# Probability

- Probability of A
- Probability of A given B
  
- Bayes Theorem

$$P(A)$$

$$P(A|B)$$

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

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

# Notation

- Model (parameter set)
- Data
- Probability of Data  
given Model

$$\vec{\theta}$$

$$x_i$$

$$P(\vec{x}|\vec{\theta})$$

## Ultimately, what questions are we asking

- Strict Frequentist: What is the probability of the data given our model (“rule out”)

$$P(\vec{x}|\vec{\theta})$$

- Bayesian: What is the probability of our model (Given our data and previous belief)

$$P(\vec{\theta}|\vec{x})$$

$$P(\vec{\theta}|\vec{x}) = \frac{P(\vec{x}|\vec{\theta})P(\vec{\theta})}{P(\vec{x})}$$

Ultimately, what questions are we asking

Easy

$$P(\vec{x}|\vec{\theta})$$

Hard/Debatable

Not as problematic  
as it seems

$$P(\vec{\theta}|\vec{x}) = \frac{P(\vec{x}|\vec{\theta})P(\vec{\theta})}{P(\vec{x})}$$



# As Physicists, we are lazy

- We will focus on the easy problem, and talk about the “harder” bits later
- It’s relevant to both kinds of interpretation, so you can ignore the statistical philosophical wars for now
- We will be focusing on 2 things:
  - Exclusion
  - Estimation

$$P(\vec{x}|\vec{\theta})$$

The Loglikelihood and  $\chi^2$

Or

“How Do We Answer Them?”

## Start with the easy thing

- The probability of independent statistical events is just the product of their probabilities
- A model, by definition, provides a probability that a given point is observed
- So, we just multiply the probabilities (And sweep the infinitesimals under the rug)

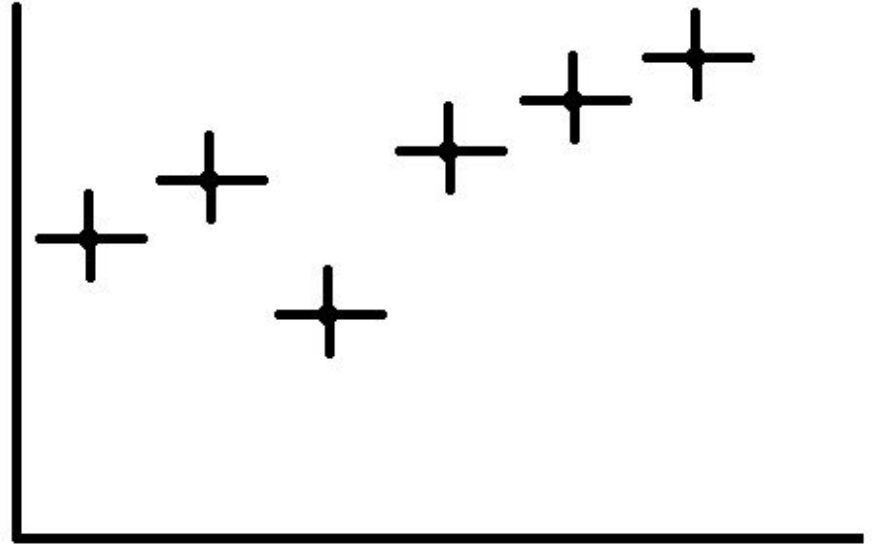
$$L(\vec{\theta}) = \prod_i pdf(x_i | \vec{\theta})$$

Computers, therefore:

$$LL(\vec{\theta}) = \ln\left(\prod_i pdf(x_i | \vec{\theta})\right) = \sum_i \ln(pdf(x_i | \vec{\theta}))$$

# What about Bins?

- “Unbinned” has strictly more information, but we may wish to bin things
- Doing complicated things to various bins for simplicity of models (Ratios, primarily, but be careful)
- Faster



## But How?

- Bin at location  $b_i$ , with count  $n_i$  - basic binned likelihood
- Assume each bin has gaussian error  $\sigma_i$ , expectation  $\mu_i$ , and observed value  $y_i$
- This is the limit of the poisson case, which is formally correct
- We like the  $\chi^2$  (for reasons to come)

$$LL(\vec{\theta}) = \sum_i n_i * \ln(\text{pdf}(b_i | \vec{\theta}))$$

$$\begin{aligned} LL(\vec{\theta}) &= \sum_i \ln\left(e^{-\frac{(y_i - \mu_i(\vec{\theta}))^2}{2\sigma_i^2}}\right) \\ &= -\frac{(y_i - \mu_i(\vec{\theta}))^2}{2\sigma_i^2} \\ &= -\frac{\chi^2(\vec{\theta})}{2} \end{aligned}$$

# What do we do with this

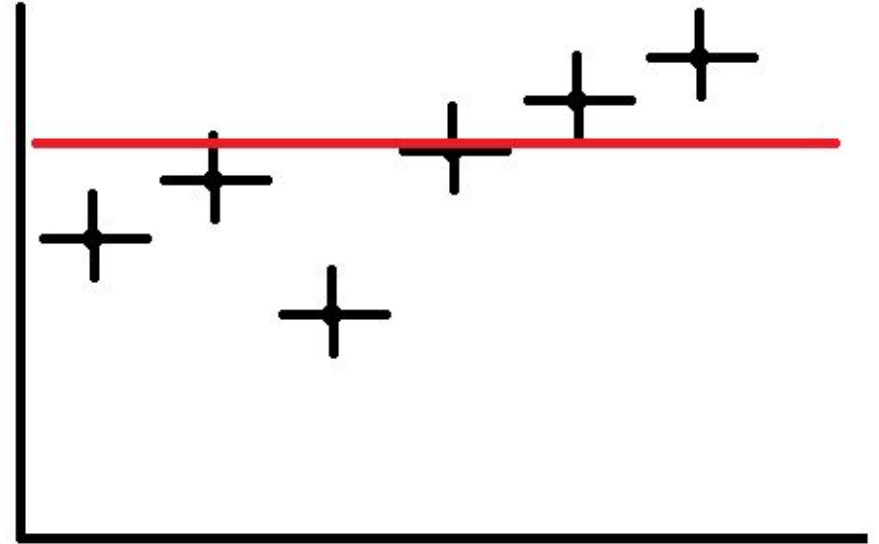
- Optimize!
- We find the model that is most likely to produce our data
- Produce deltas
- The LL is the “best” information you can get
  - $\chi^2$  is nice too I guess

$$\Delta LL = \text{Max}_{\vec{\theta}}(LL(\vec{\theta})) - LL(\vec{\theta}_0)$$

$$\Delta\chi^2 = \chi^2(\vec{\theta}_0) - \text{Min}_{\vec{\theta}}(\chi^2(\vec{\theta}))$$

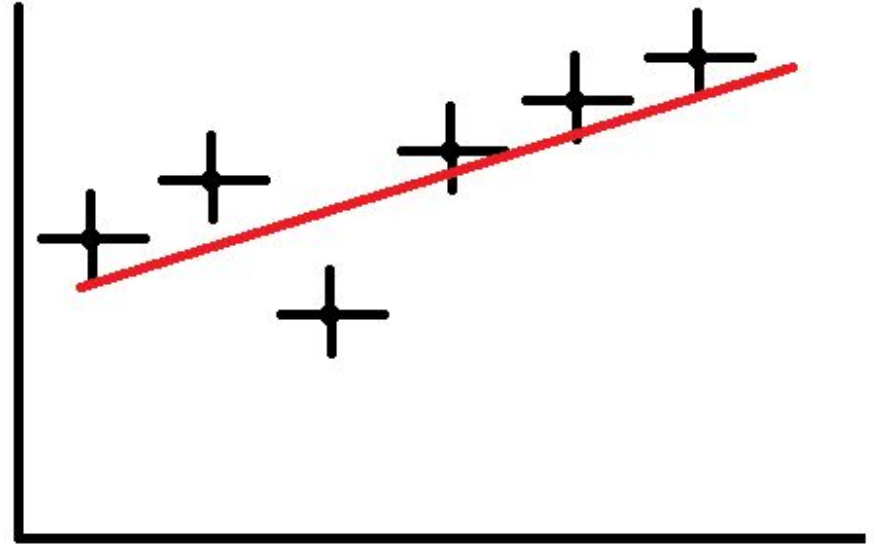
# The null model

- Just a boring model with 0 or no parameters (No slope, for instance)



# The Actual Model

- This has some parameters and dynamism
- You can see the  $\chi^2$  going down as it fits better
- So now what





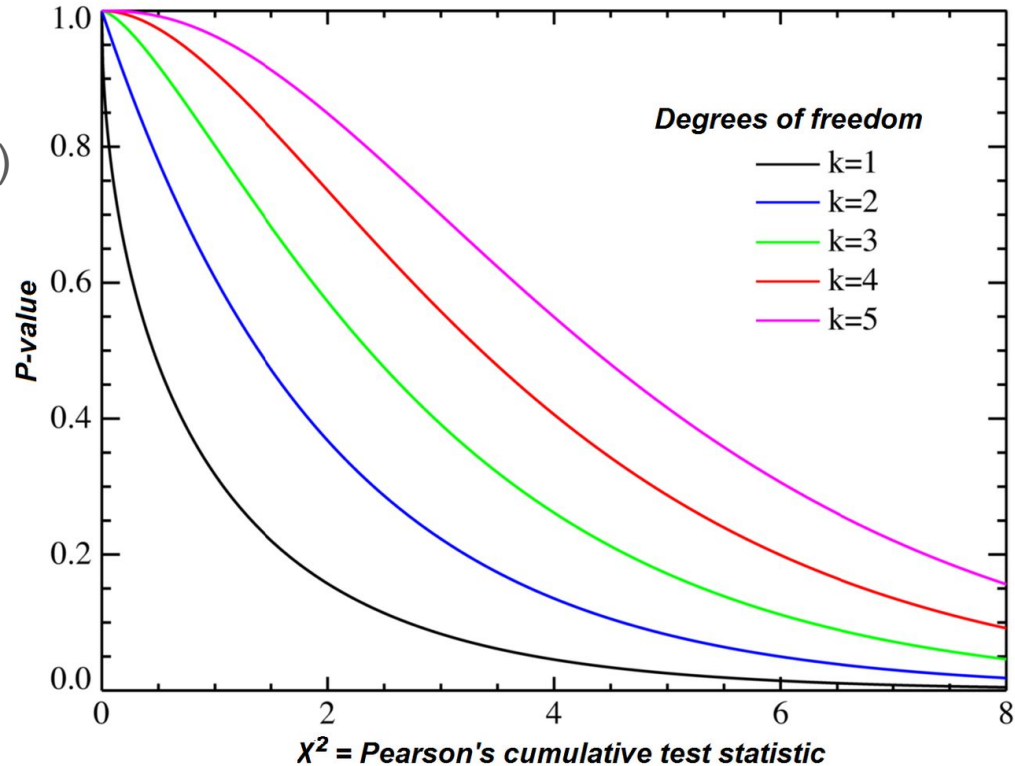
Test Statistics and Confidence Regions

Or

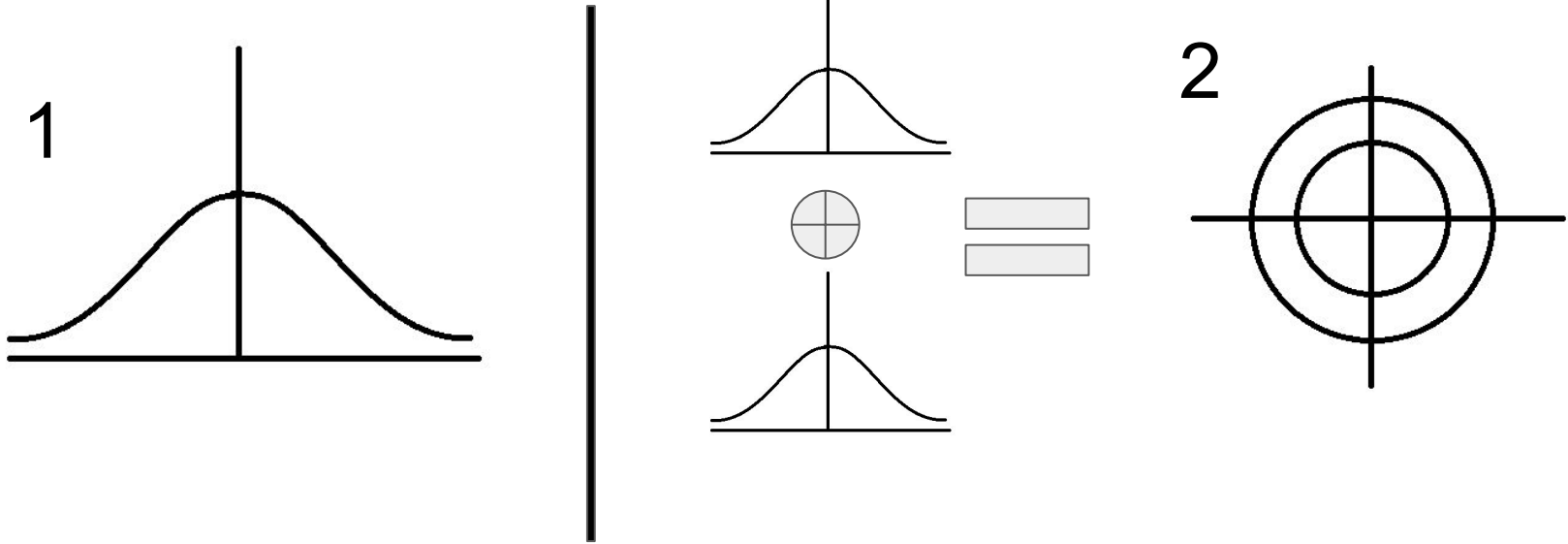
“What Does the Answer Mean?”

# These are what we call “Test Statistics”

- If you take a stats class, you will see others (t-tests, k-tests, Rs, etc)
- We don't care about those - you can usually derive them from the LL (but you shouldn't - just look them up when needed)
- But how do we interpret them?
- Let's start with the  $\chi^2$  - its the easiest because it's a bunch of gaussians



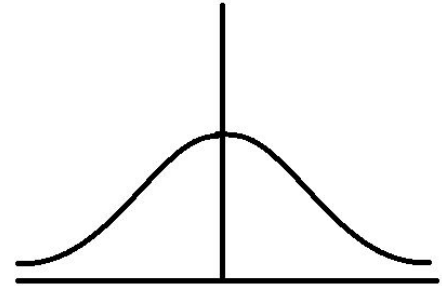
# “Degrees of Freedom”



# “Degrees of Freedom”

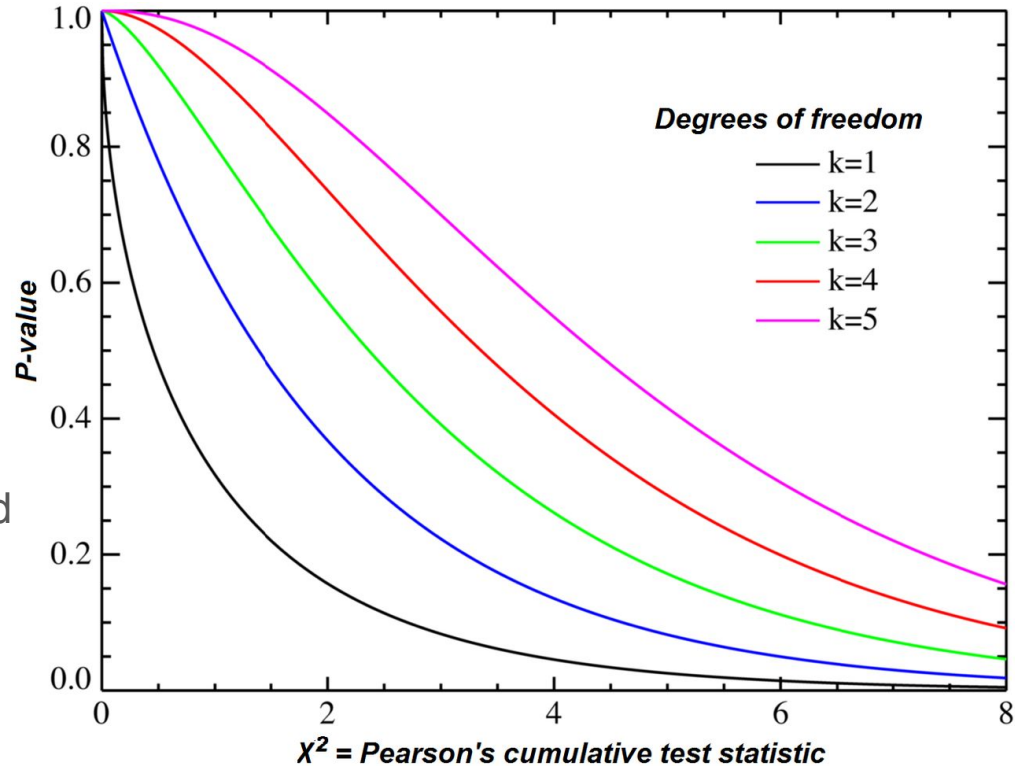
- N degrees of freedom
- Pretty easy to calculate analytically
- Rests on the assumption that everything underlying is gaussian
- Protip: Mean of  $\chi^2/N$  approaches 1 for large N

N



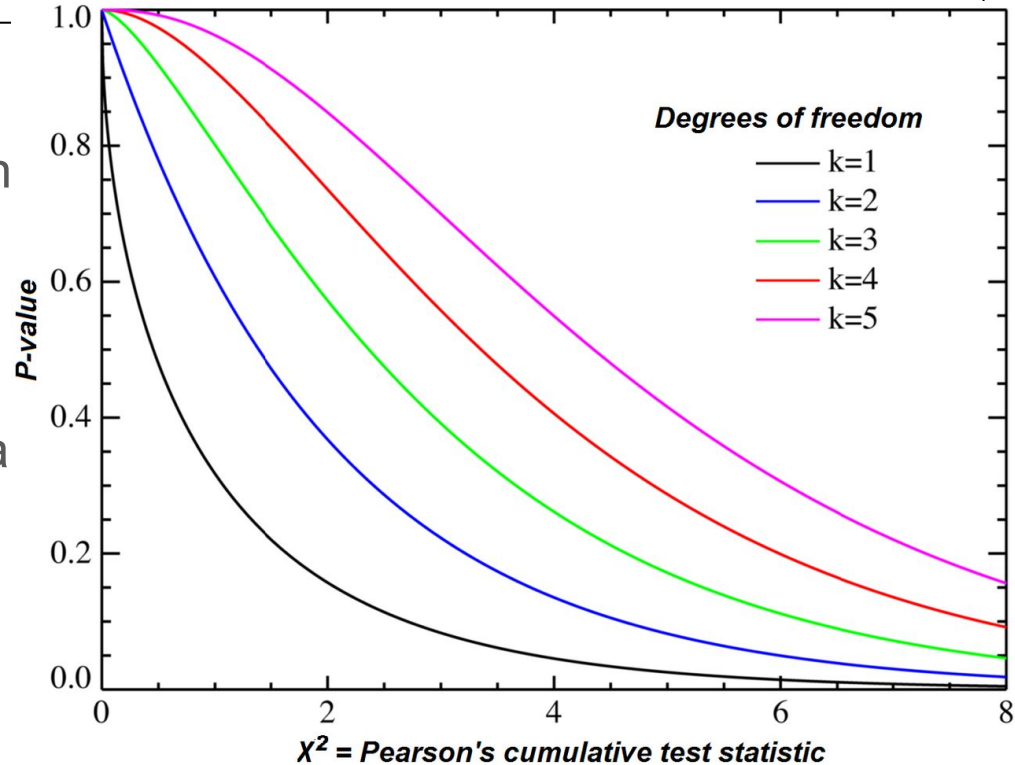
# The “P-value” - what the test statistic is for

- We get a “p-value”
- A p-value is the probability that our test statistic will be this weird or weirder
- In physics, we often convert this p-value to an equivalent  $\sigma$  - how far one would have you be from a standard gaussian to be just as weird
- We mark  $3\sigma$  (1/300) as “evidence” and  $5\sigma$  ( $3e-7$ ) as “discovery”
  - I have many opinions on many things, but especially on these



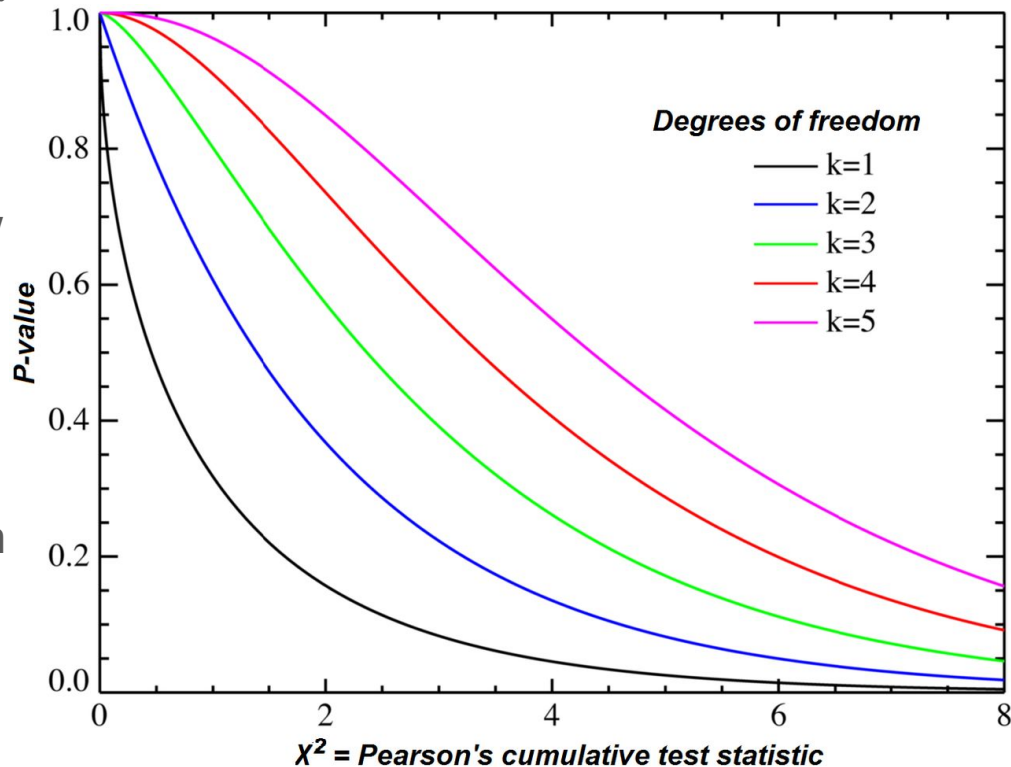
## Why are we talking about $\chi^2$

- I've pulled a fast one on you
- We talk about the  $\chi^2$  distribution differently from the  $\Delta\chi^2$  - The former is a theoretical distribution from a combination of gaussians, and the latter is a common test statistic
- It is important to keep those straight



# Why are we talking about $\chi^2$

- Under certain assumptions, according to a theorem known as “Wilks Theorem”, the  $\Delta\chi^2$  will follow a  $\chi^2$  distribution with  $n$  degrees of freedom where  $n$  is the number of parameters in your model
- Not only that, but a  $2*\Delta LL$  will ALSO follow the  $\chi^2$  distribution with  $n$  degrees of freedom



# Interlude: An Example and Trials

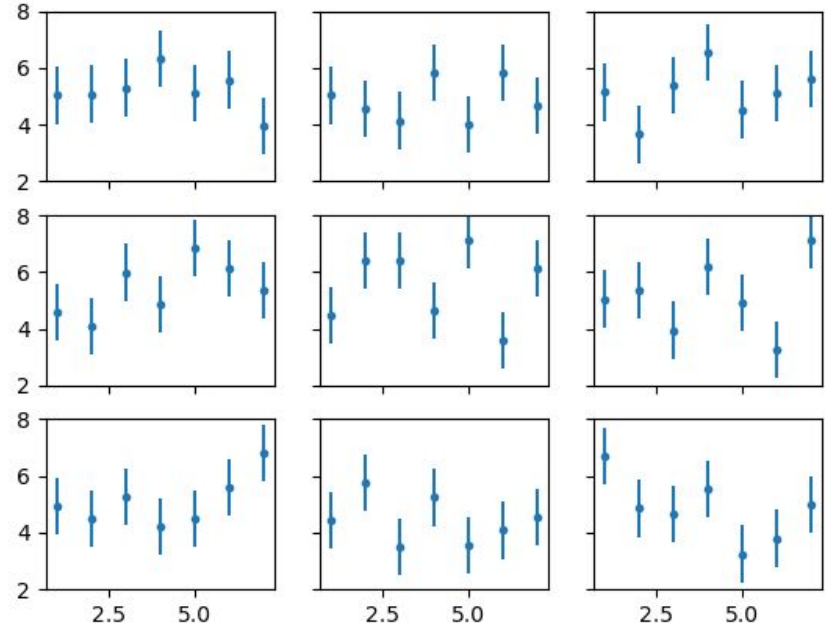


# What is a Trial

- We are using something easy to get something hard
- We know (easy):

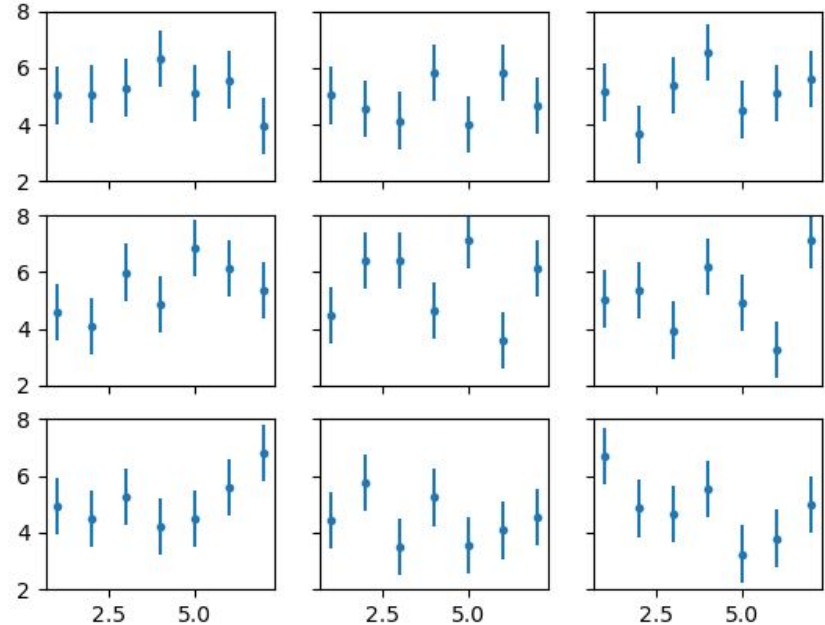
$$pdf(x_i | \vec{\theta})$$

- So we can “throw” “realizations”
- We create fake data that would be the result of our ideal model and examine how our statistical procedures act on them (hard)



# A Toy Model

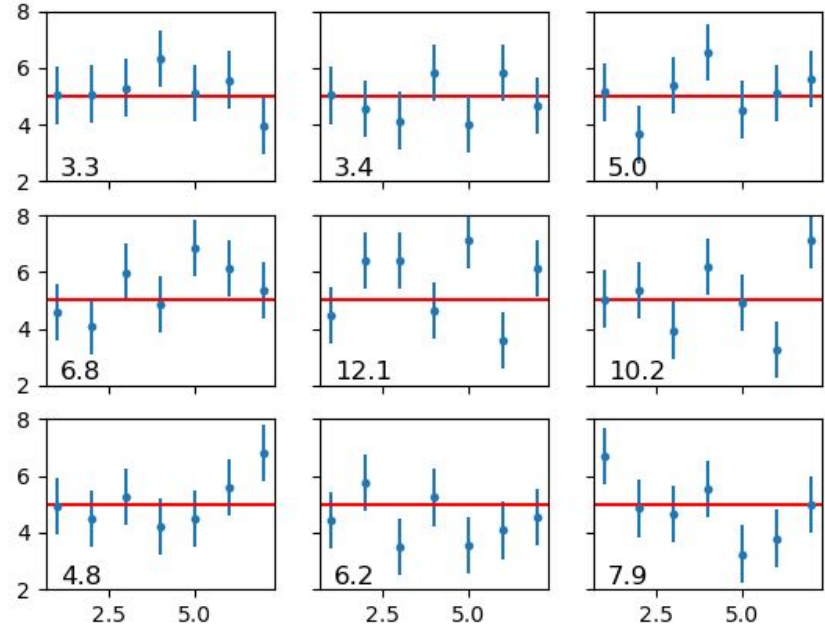
- The x's are the integers 1 through 7
- The y's are thrown with a mean of 5 and standard deviation of 1
  - You can think of this as “binned” or as a model of 2d data - these are equivalent formulations
- I'm going to show 9 “realizations” but report the summary statistics for 1600
- We'll start with the simplest model
  - A flat line at 5, standard deviation 1
- Then we'll add parameters



# Simplest: The Expectation

- The simplest model is just to have a flat line at the model mean
- The  $\chi^2$  is on each plot
- You can see it is a simple, flat model
- Average  $\chi^2$  for the large sample: 7.1

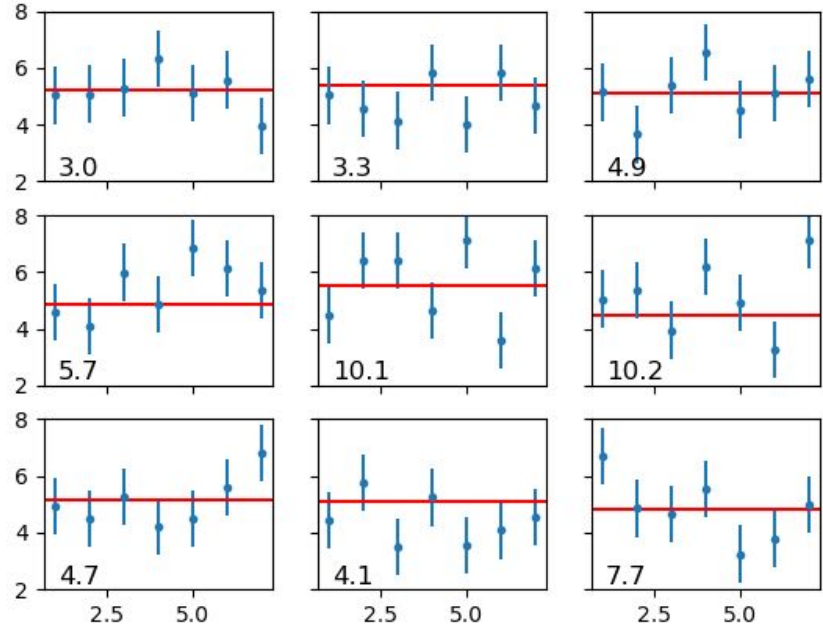
From 5  $\chi^2$



## Next: The Mean

- The next simplest model is just to have a flat line in 1 direction
- But this line is fit per sample
- The  $\chi^2$  is on each plot
- You can see it is a simple, flat model
- Average  $\chi^2$  for the large sample: 6.0

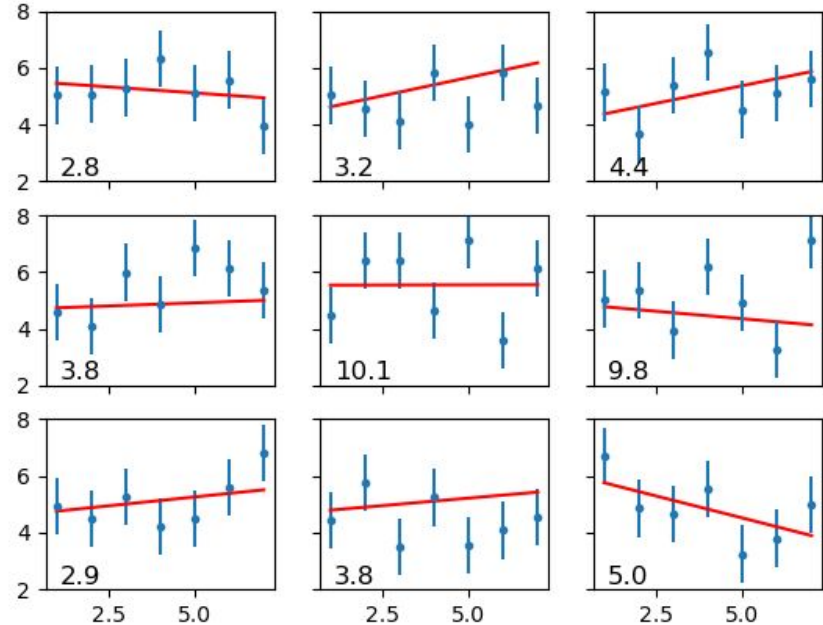
## No-Slope $\chi^2$



# Finally: Add a slope

- Add 1 more parameter to add a slope
- Again, we are fitting
- The  $\chi^2$  is on each plot
- The  $\chi^2$  is smaller - more complex models necessarily fit better
- Average  $\chi^2$  for the large sample: 5.0

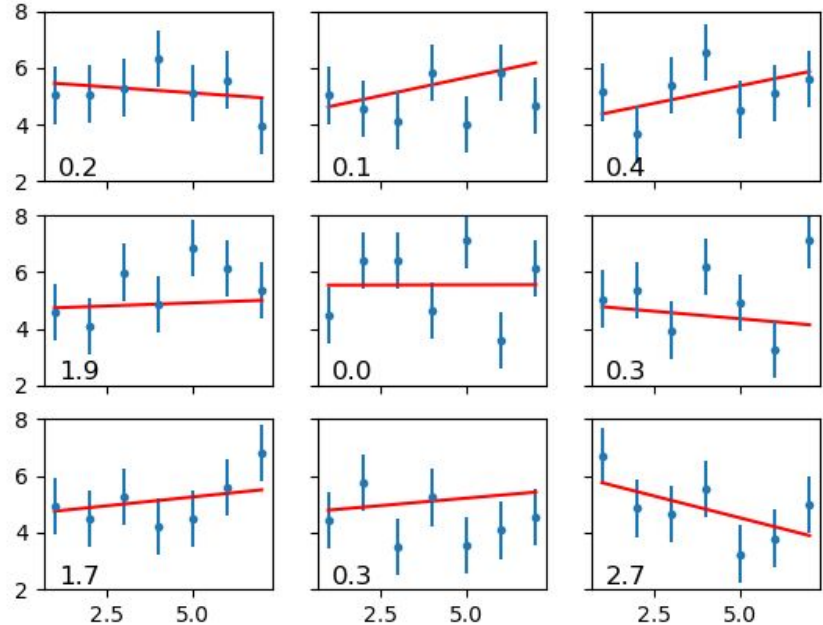
## Slope $\chi^2$



# Finally: How much Better

- Subtract the  $\chi^2$  of the more complex model from the simpler model and get the  $\Delta\chi^2$
- This is the “slope” from the “no slope”
- Average  $\Delta\chi^2$  for the large sample: 1.0

## Slope $\Delta\chi^2$

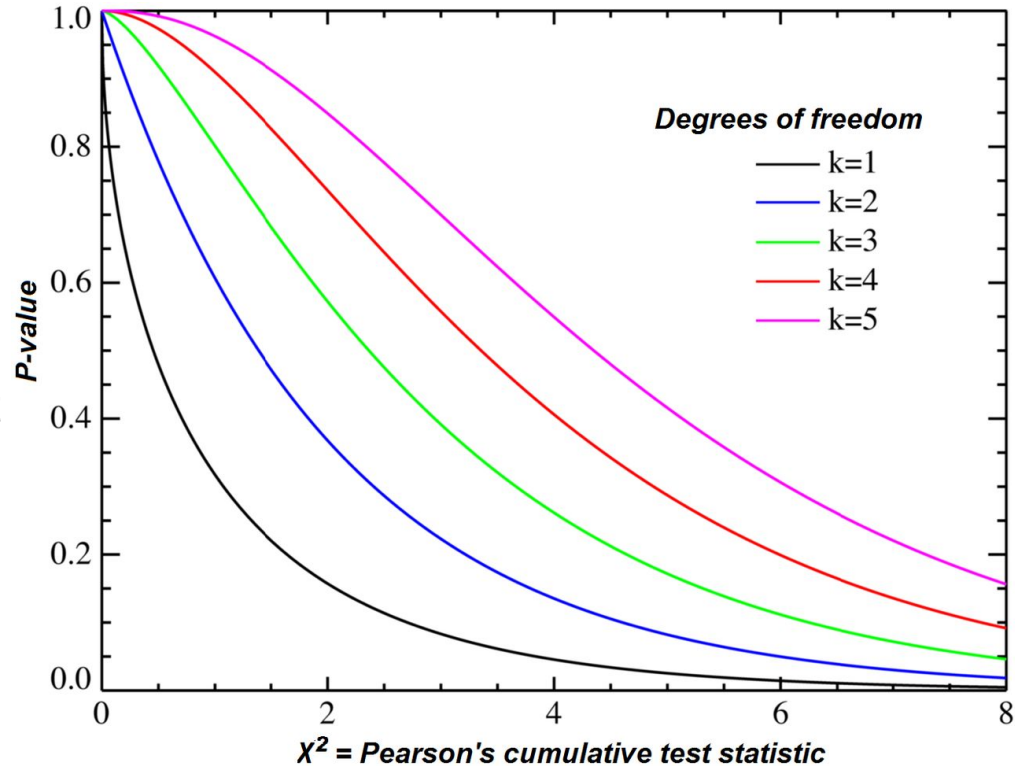


# This was a convenient model

- It is fundamentally built out of gaussians
- This means that the  $\chi^2$  played nicely
- It also fit the assumptions of Wilks' theorem, so the  $\Delta\chi^2$  behaved predictably
- We could have picked any test statistic (LL, or something like a KL-divergence or a KS test) and used the large realization sample to understand their response
- It is important to do this because:

# We often violate the assumptions

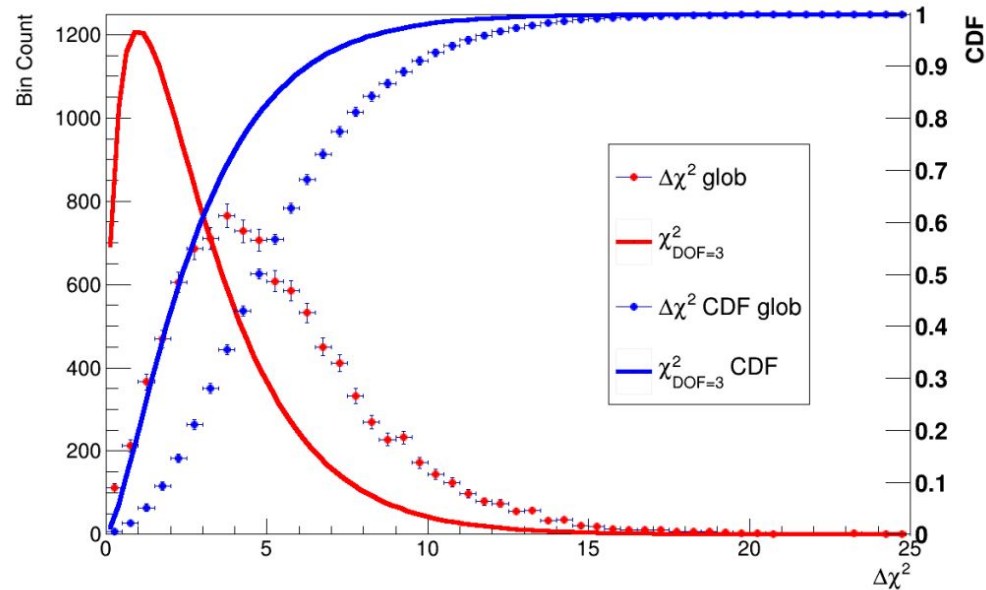
- We OFTEN violate those conditions in physics making us more and less “conservative” - in those cases, the  $\Delta\chi^2$  will NOT follow a  $\chi^2$  distribution
- When this is the case, it is important to build up your own test statistic distribution by throwing from your null mode and applying your fitting procedure. You will have to throw  $O(1/(\text{desired pval}))$  to do it properly





# One such deviation from assumptions

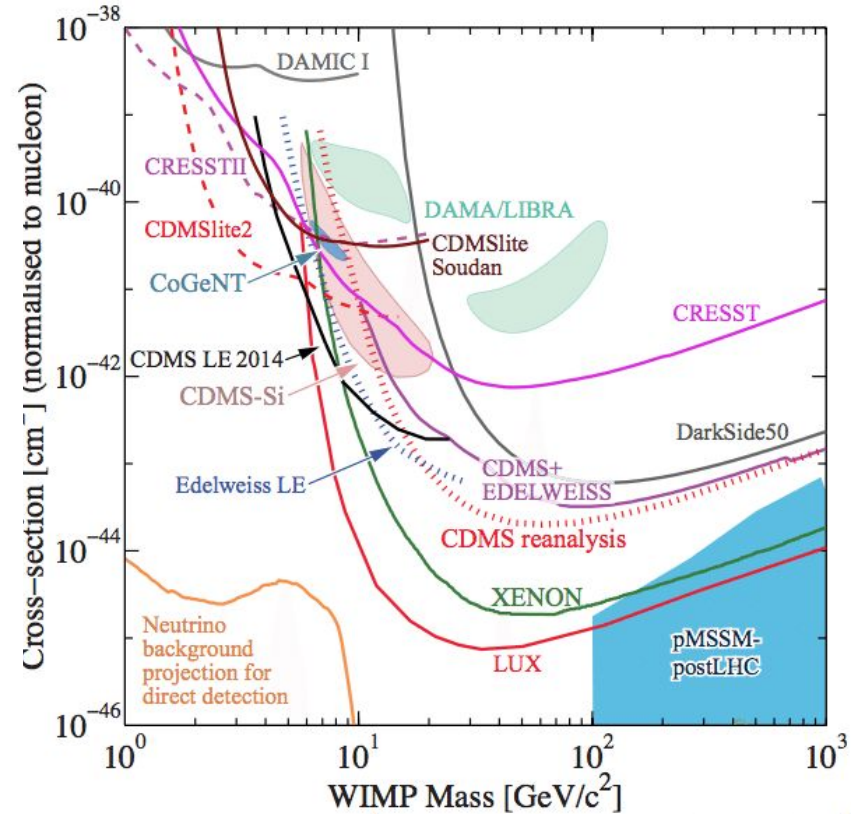
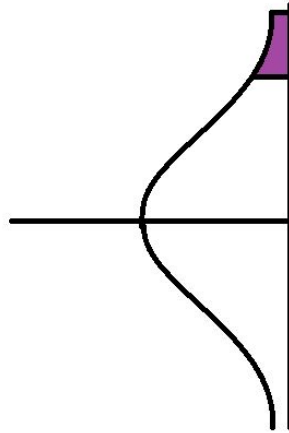
- I wrote it up
- It is here:  
<https://arxiv.org/abs/2211.06347>
- Neutrino oscillation models violate Wilks' theorem in non-conservative ways
- If you would like to tell me anything confusing about it, I would be grateful for any feedback here:  
<https://forms.gle/TcbdsdopMv9APti48>



End Interlude

# Using this for Exclusion

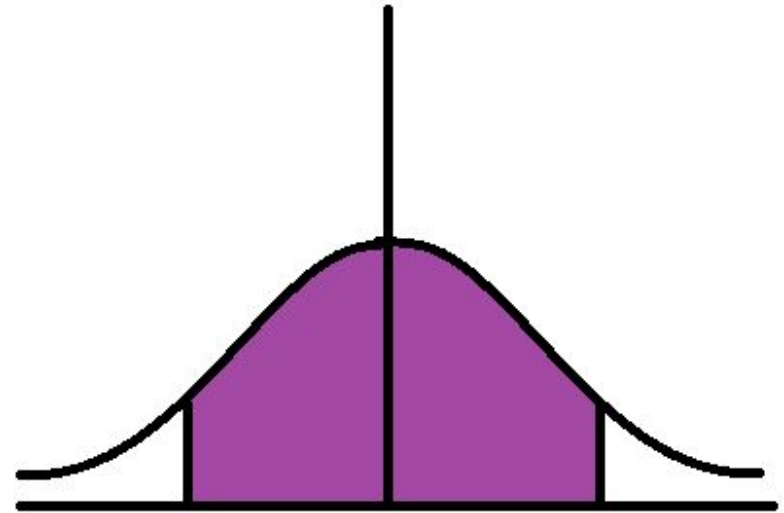
- This is very traditional science
- We reject the null at XX%
- Use the distribution of your test from Wilks' or from trials
- This is the mode we use for discovery
- I like to think about it as a 1-D gaussian, but this is often done in many parameters



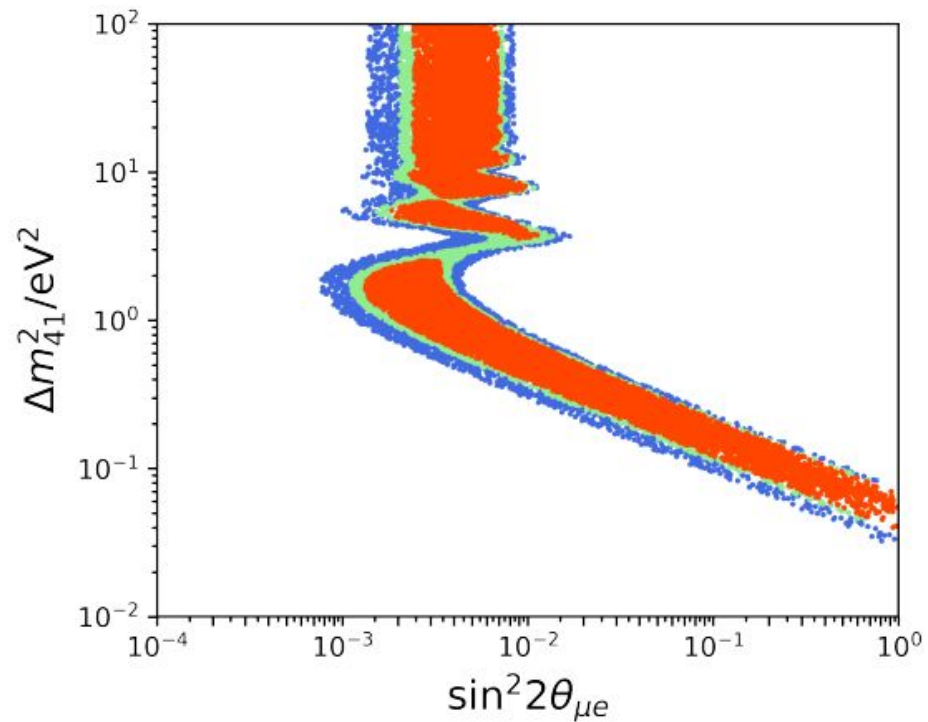
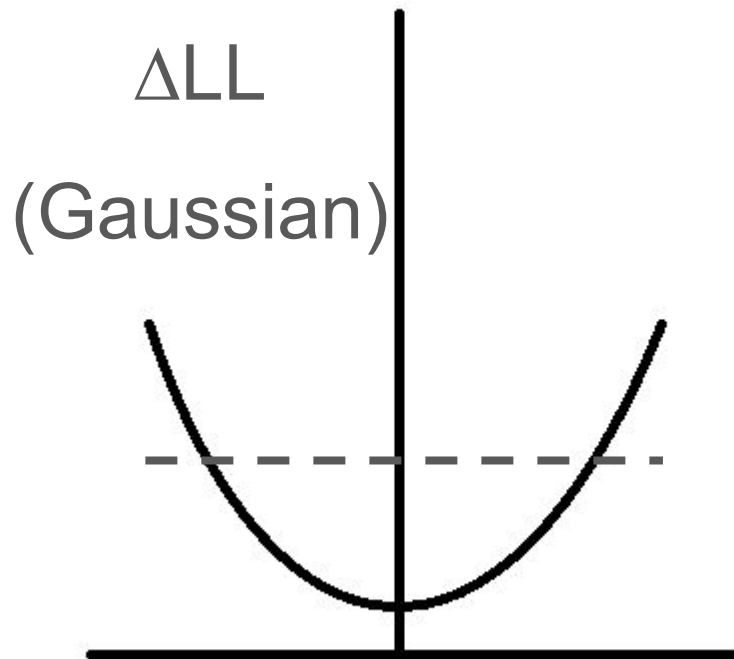
[https://www.researchgate.net/figure/Current-exclusion-limits-and-regions-of-interest-of-dark-matter-searches-for-fig1\\_320890690](https://www.researchgate.net/figure/Current-exclusion-limits-and-regions-of-interest-of-dark-matter-searches-for-fig1_320890690) (2017)

# If we find something

- Parameter estimation!
- As good frequentists, we construct confidence intervals
- Formally, a 90% confidence interval is constructed to contain the true parameter 90% of the time
- Many degenerate ways of doing this, but we normally pick density
- We can again use Wilks' theorem under certain assumptions OR we should check with realizations
- Often we have to only check some of the model space

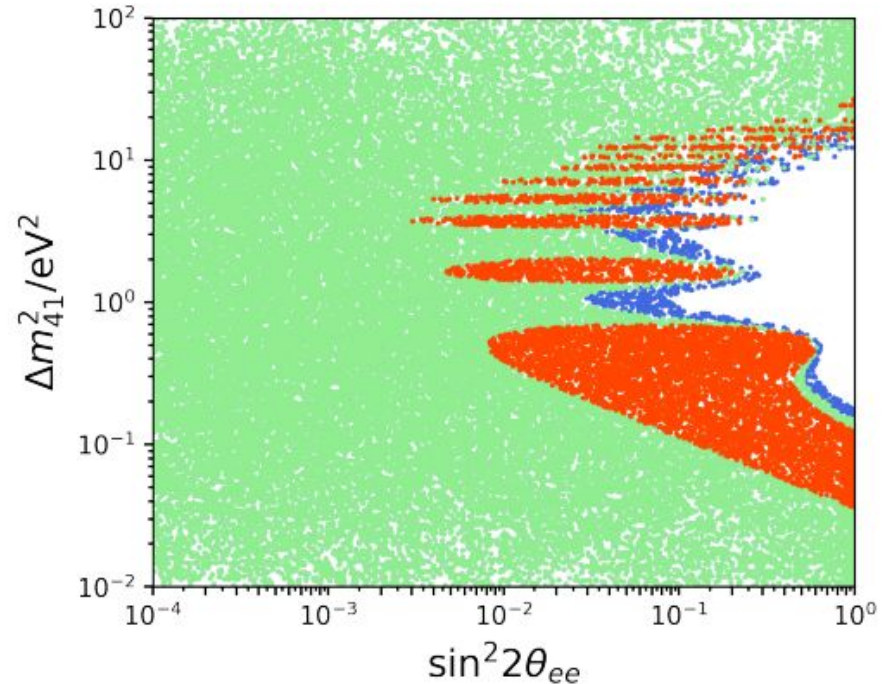


## Other views



# Summary for Exclusions and Estimation

- The  $\Delta LL$  is the gold standard for information about your space
- Wilks' theorem is often useful as a shortcut
- You should check it anyway



## Aside: Bonus Tip

- You may have noticed something squirrely about the normalization in the LL
- If you have the overall normalization as a parameter (as is often convenient), it will want to go to infinity
- You can fix this with a normalized LL
- The math comes from a poisson distribution, but I won't bore you with it

$$LL(\vec{\theta}) \rightarrow LL(\vec{\theta}) - \int_x F(\vec{\theta})$$

Real Experimental Considerations

Or

“What is Everyone Actually Arguing About at Working Group Meetings?”



# What assumptions have I been hiding?

- First, that the data are independent - they might not be for a wide variety of reasons
  - Detector correlations, environmental factors, etc
  - “Known unknowns”
- Second, that we completely understand  $P(x)$ 
  - There may be things about our detector we don’t understand
  - “Unknown unknowns”
- These are controlled for and talked about as “systematics”

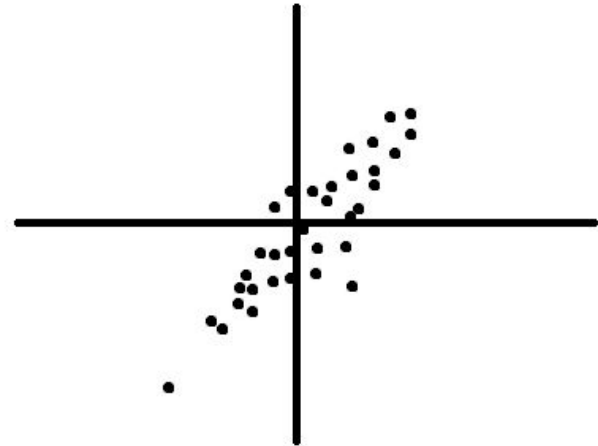
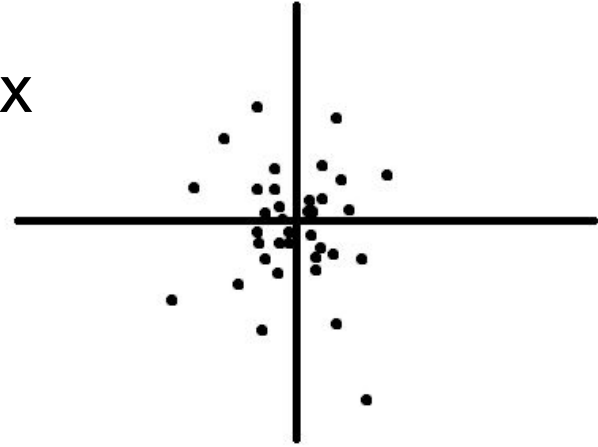
# Systematics V1: Pull terms

- The Loglikelihood version
- “Pull terms”
- We assume we know something, but not everything, about where the non physically relevant parameters ought to be
- So we penalize the likelihood for moving away
- Assume we have two physics parameters, and  $\theta_2$ , is, say, DOM efficiency
- $\theta_2$  controls relationship among data - overall more or less likely

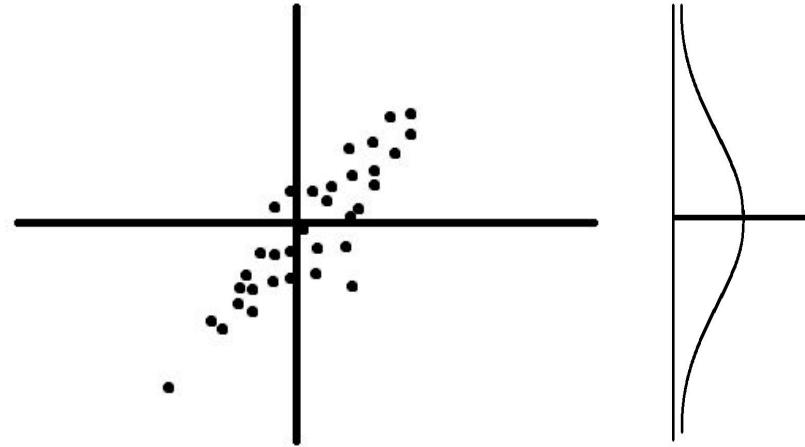
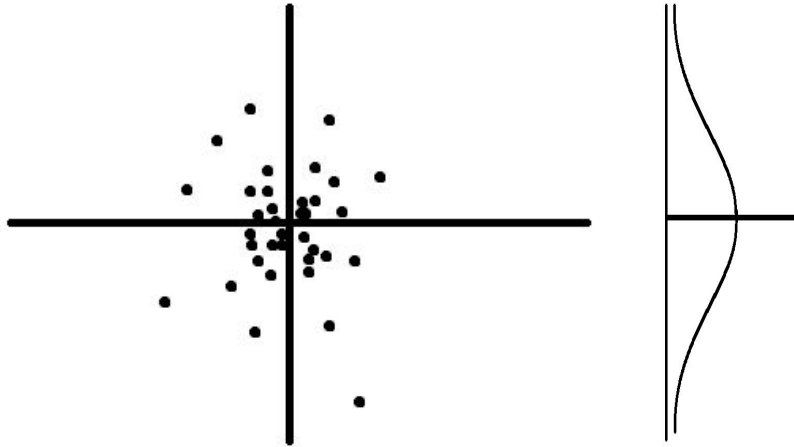
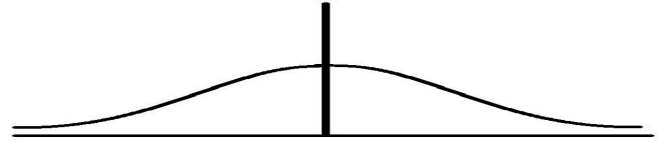
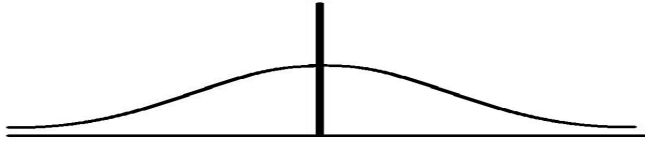
$$LL(\vec{\theta}) \rightarrow LL(\vec{\theta}) - \frac{(\theta_2 - \mu_{\theta_2})^2}{2\sigma_{\theta_2}^2}$$

# Systematics V2: Covariance Matrix

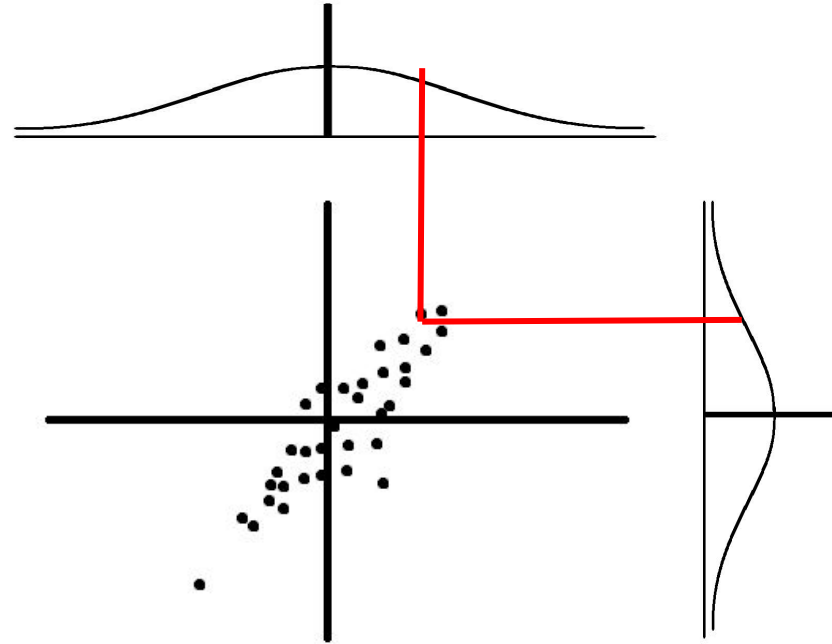
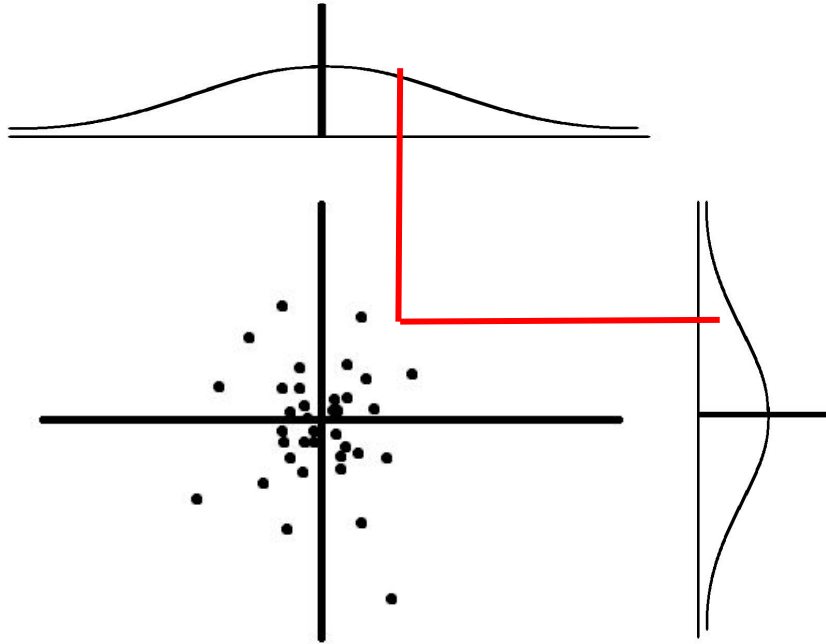
- Plot your multidimensional data against itself
- Check to see if it is truly “uncorrelated”
- If it isn't, you find the correlations, and you've better understood your space



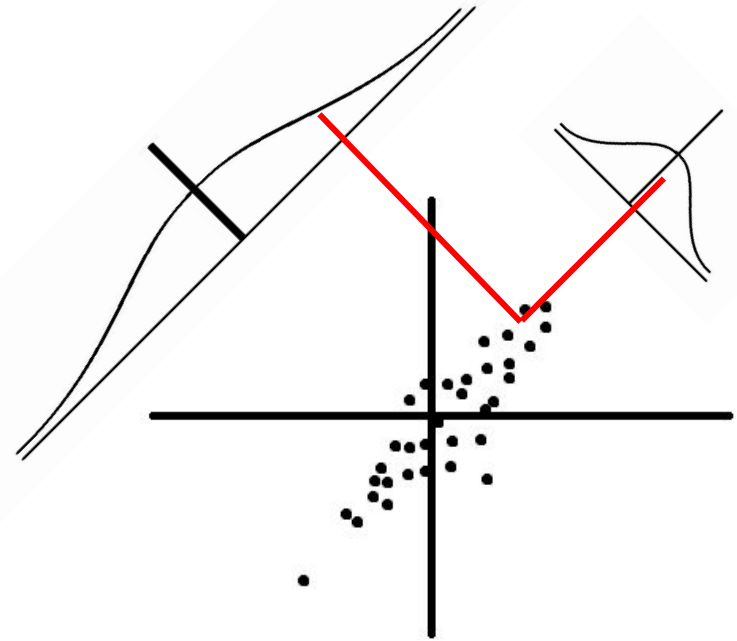
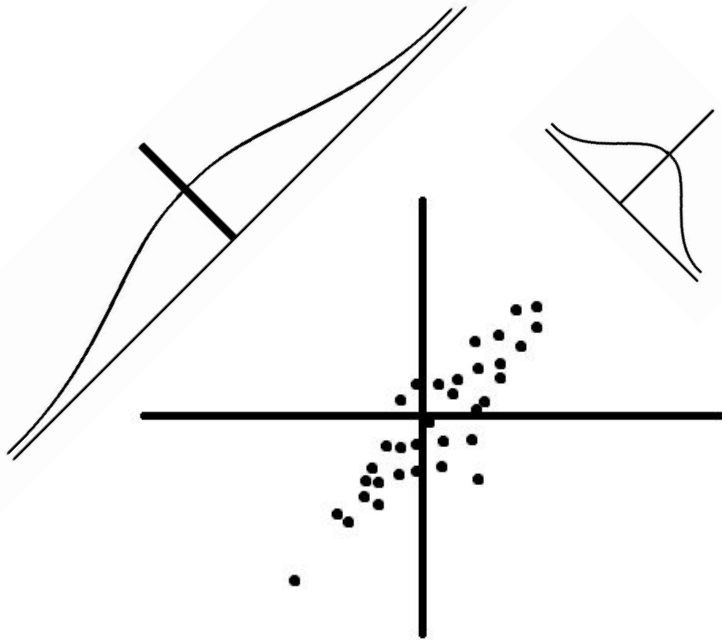
# Systematics Part 2: Covariance



The Question: How weird is the same point



# What do we do



# The Covariance Matrix

- When there is no correlation, it is diagonal
- Can be “block diagonal”
- Makes it easy to convert a set of correlated residuals ( $r$ ) to a  $\chi^2$
- This is useful, but relies on a gaussian assumption
- In principle, the correlation matrix is a function of model parameters - be careful

$$r_i(\vec{\theta}) = n_i - \mu_i(\vec{\theta})$$

$$\chi^2(\vec{\theta}) = \vec{r}(\vec{\theta})^T C(\vec{\theta}) \vec{r}(\vec{\theta})$$

# General Systematics Warnings

- Systematics can be as much art as science
- You will mostly see them as pull terms or correlation matrices
- In principle, you are worried about things that look like your signal, even a little bit
- Check carefully - trials are your best friend





## Aside: penalties and regularization

- Occasionally, you will have a very general model
- You might want to avoid overfitting by penalizing amplitudes
- These things are sort of Bayesian (the math tends to be the same)
- This looks like systematics, but this is usually not what is meant

$$F(\vec{A}) = \sum_j^n A_j x^j$$

$$LL(\vec{A}) \rightarrow LL(\vec{A}) - \sum_j^n |A_j|^m$$

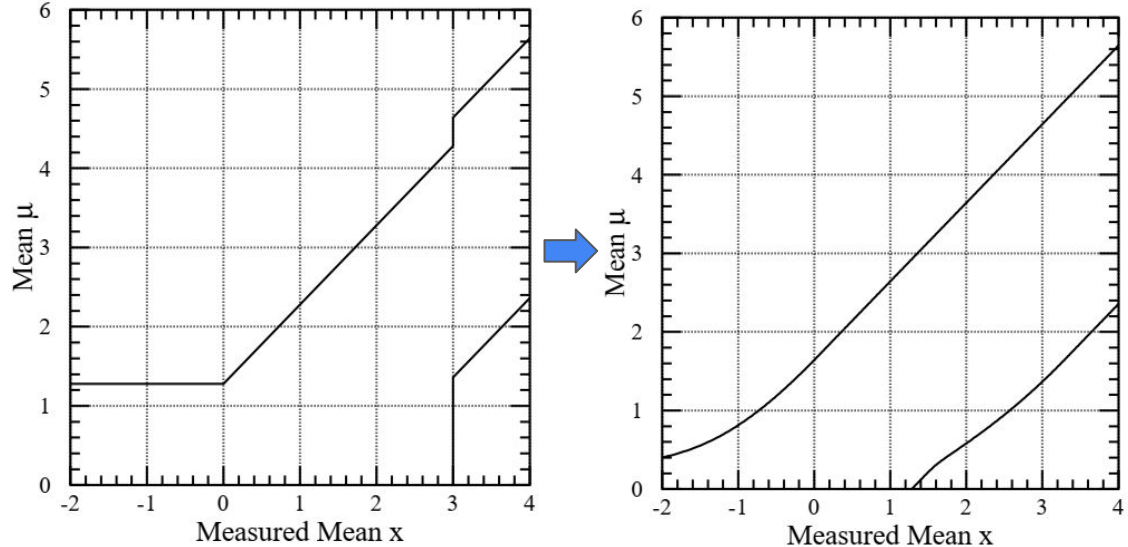
Physics Jargon  
Or  
“What Did They Mean by That Barely  
Googlable Term?”

# Sigma

- We talked about this already
- My pet peeve: “ $\sigma$ ”  $> \sim 8$  - that is not a real p-value that you can check
- Look at [arxiv.org/ftp/arxiv/papers/1103/1103.5672](https://arxiv.org/ftp/arxiv/papers/1103/1103.5672)
- $7 \sigma$  is  $1/7.8e11$
- $10 \sigma$  is  $1/1.3e21$
- $25 \sigma$  is  $1/3.3e135$
- Tails don't stay that gaussian that long
- “We were seeing things that were 25-standard deviation moves, several days in a row,”
  - David Viniar, just before 2008
- “This fits to 11 Sigma”
  - A talk on the (real, but not that real) pentaquark

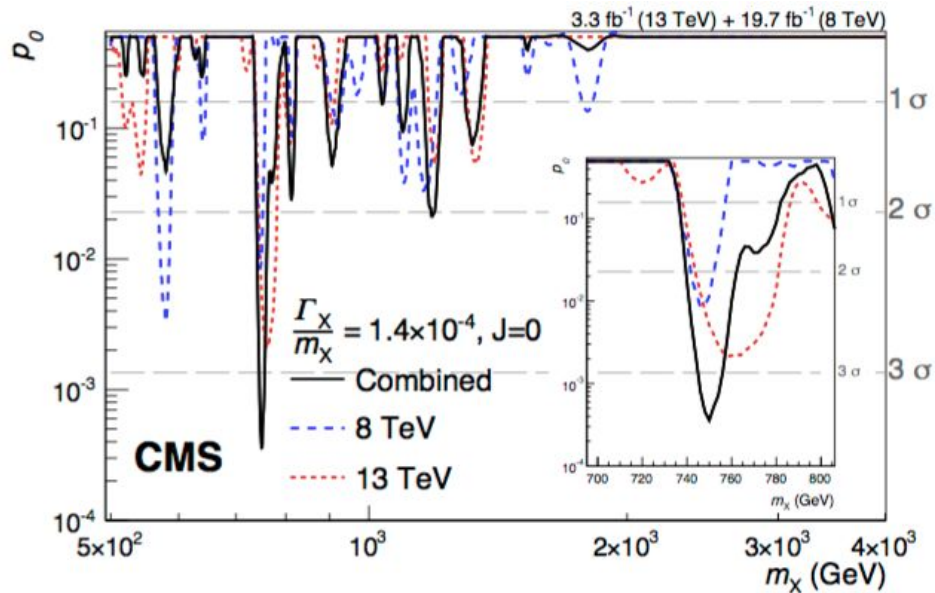
# Feldman-Cousins

- You should go read this paper
- Essentially, it aims to unify exclusion and estimation
- Other meanings: Trials



# Look-Elsewhere

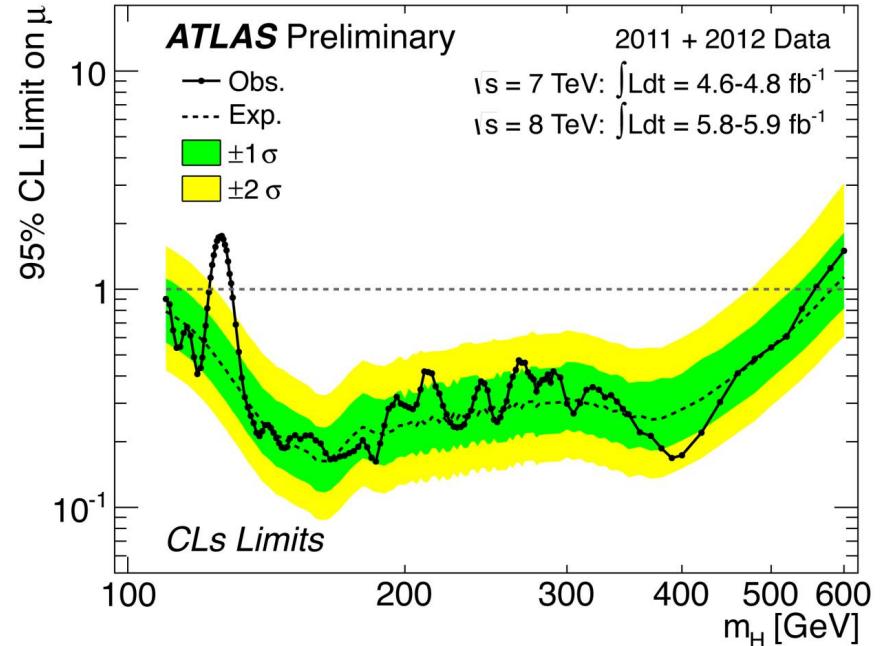
- One of the ways that Wilks fails
- If you are looking for a bump on noise, you can focus on each slice of the plot individually
- This gives you more bites at the apple than Wilks assumes
- ALWAYS run trials



<https://arxiv.org/pdf/1606.04093.pdf>

# Brazil Plots

- Popularized by the Higgs Boson
- Run null trials, then general exclusion limits
- Put the middle  $1\sigma$  of bands in green and the middle  $2\sigma$  in yellow
- Then overlay the exclusion from actual data
- Where the data fails to exclude as strong as you expect is where you expect your new physics



<https://home.cern/resources/image/physics/infographics-gallery>

# Final Summary

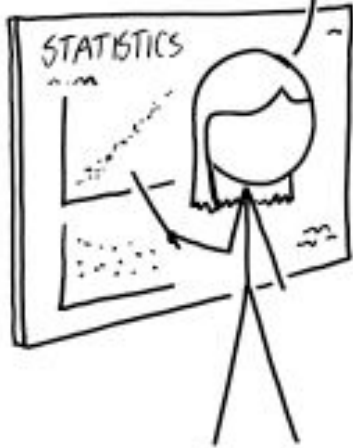
- 1) Probability (Bayesian? Frequentist? Strange notation?)
- 2) The Loglikelihood and  $\chi^2$  (How I learned to stop worrying and love LL)
- 3) Test statistics and confidence regions (How I learned to stop worrying and love running null realizations)
- 4) Real Experimental Considerations (What is really going on)
- 5) Physics statistics jargon (How we talk about things)

Thank you  
Or  
“Any Questions?”



# Rejected General Systematics Warnings

IF YOU DON'T CONTROL FOR  
CONFOUNDING VARIABLES,  
THEY'LL MASK THE REAL  
EFFECT AND MISLEAD YOU.



BUT IF YOU CONTROL FOR  
TOO *MANY* VARIABLES,  
YOUR CHOICES WILL SHAPE  
THE DATA, AND YOU'LL  
MISLEAD YOURSELF.



SOMEWHERE IN THE MIDDLE IS  
THE SWEET SPOT WHERE YOU DO  
BOTH, MAKING YOU DOUBLY WRONG.  
STATS ARE A FARCE AND TRUTH IS  
UNKNOWNABLE. SEE YOU NEXT WEEK!

