

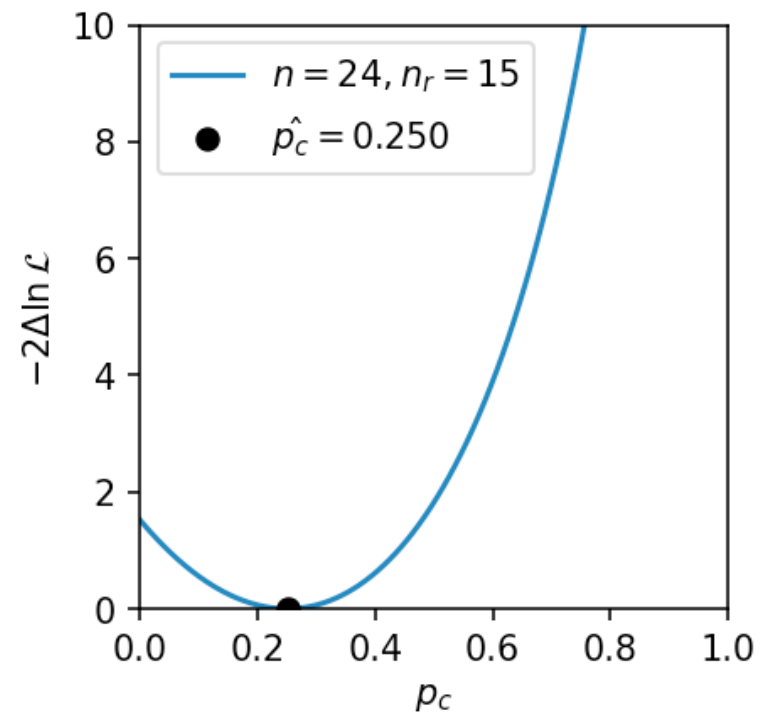
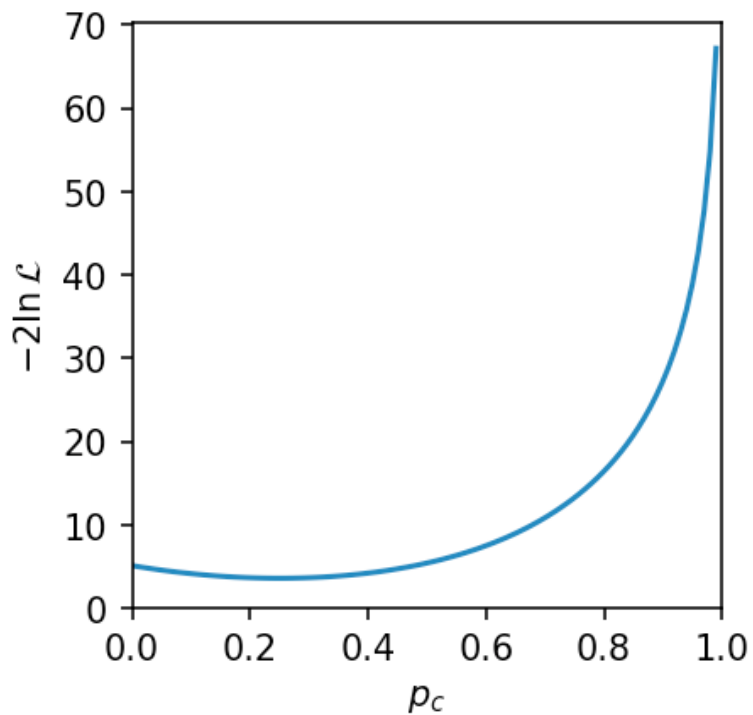
Intervals

Topics: intervals

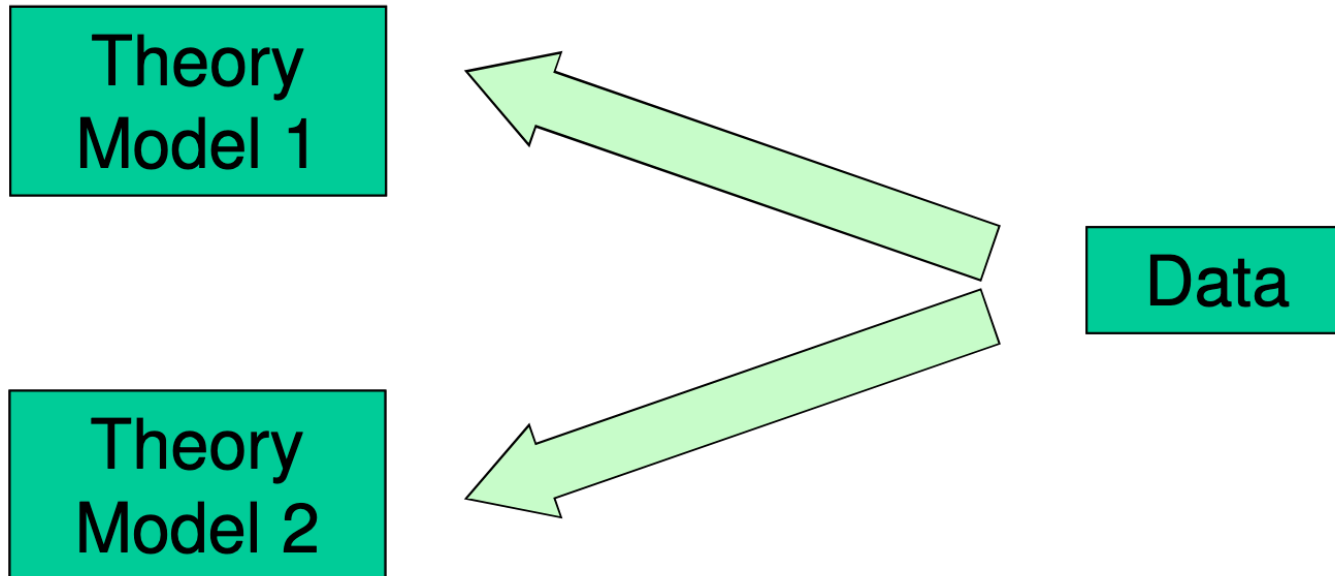
- Hypothesis tests
- Neyman intervals
- Likelihood ratio test statistic
 - Neyman-Pearson lemma
- Under-fluctuation and CLs
- Asymptotic behavior - Wilk's theorem

Inferring P(cheat) as a frequentist

- Scan $-2 \ln \mathcal{L}(p_c) = -2 \ln f(n_r; n, p_c)$
- Find minimum
- ???
- Profit



Hypothesis tests



Which hypothesis is the most consistent with the experimental data?

Hypothesis tests

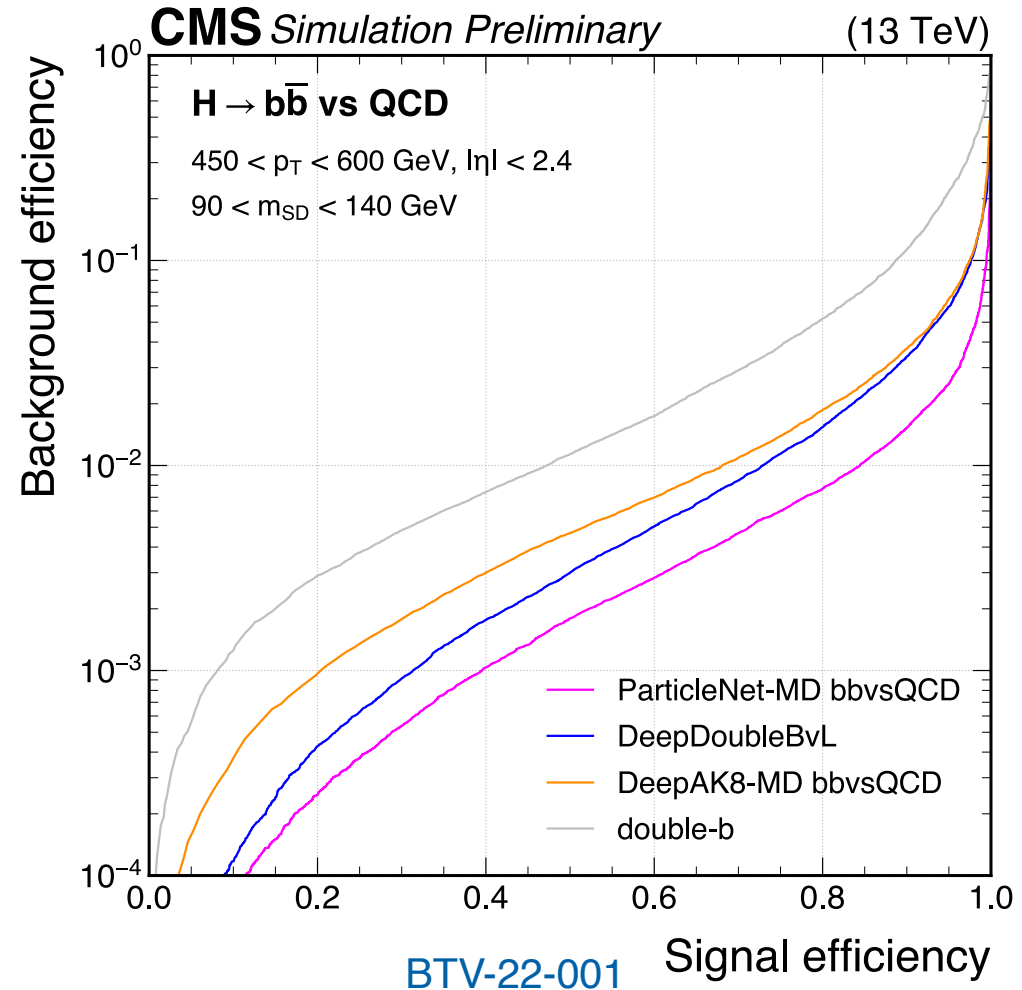
- Simple test parameterized by two probabilities: α , β
 - β with respect to an alternate hypothesis

Table of error types		Null hypothesis (H_0) is	
		TRUE	FALSE
Decision about null hypothesis (H_0)	Don't reject	Correct inference (true negative) (probability = $1-\alpha$)	Type II error (false negative) (probability = β)
	Reject	Type I error (false positive) (probability = α)	Correct inference (true positive) (probability = $1-\beta$)

Hypothesis test example

- ROC plot: let H_0 be QCD jet
 - DNN tagger score is the test statistic
 - Lower values more consistent with H_0
- Background efficiency
 - True negative: accept H_0 when true
 - $1 - \alpha$ (where α is the *size*)
- Signal efficiency
 - True positive: reject H_0 when false
 - *Power* $1 - \beta$
- Prefer tests with higher *power* for a given *size*
 - Here, test size is large
 - Usually test size is *very small*

Comparison of the performance of the $X \rightarrow b\bar{b}$ identification algorithms ... in the 2018 data-taking conditions



Coverage (confidence) sets

- One can ask, “assuming a value of θ , would x_{obs} be a likely outcome?”
 - This is a hypothesis test (sufficiently likely or not)
- We have $P(x;\theta)$, so we can answer this if we:
 - Choose a *size* (significance level) α of the test (e.g. 0.05)
 - Define a test statistic (ordering) of possible outcomes
 - Run pseudo-experiments (toys) for each θ to determine distribution of test statistic
 - Perform experiment, report set of θ where test statistic is below the $1-\alpha$ quantile
- A good ordering will lead to
 - Good *coverage*: in repeated experiments the (unknown) θ_{true} will be in the set with probability at least $1-\alpha$, though it may over-cover
 - High *power* ($1-\beta$): the set does not contain θ_{alt} for some specified alternative hypothesis

Neyman interval

- Neyman construction:

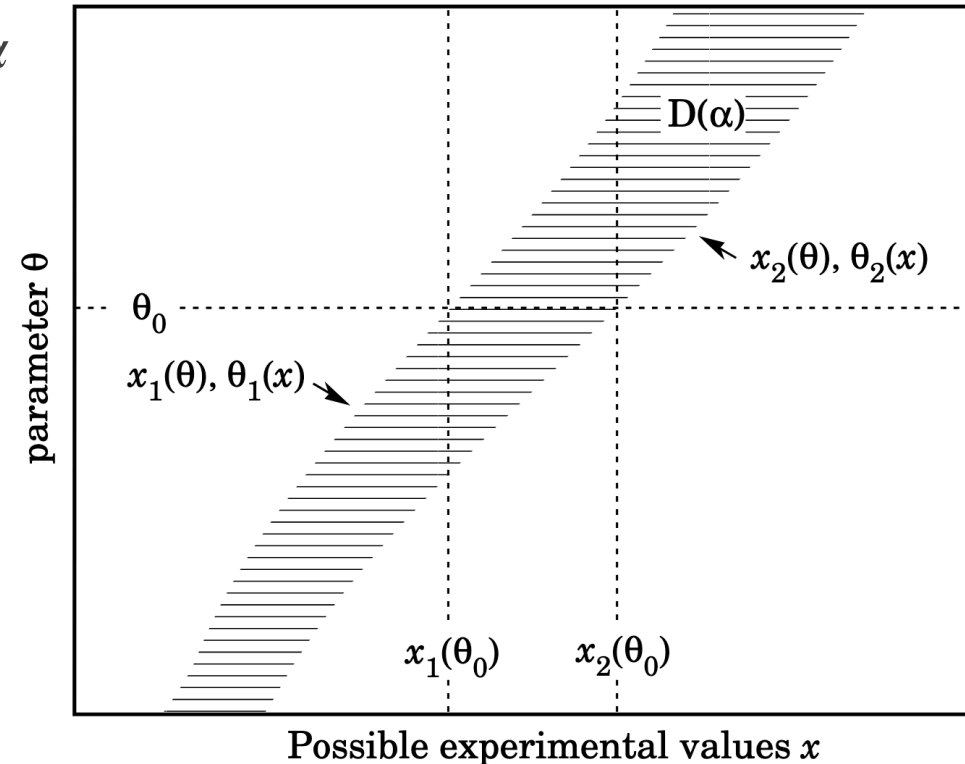
- For each θ , find range $[x_1, x_2]$ s.t.

$$P(x_1 < x < x_2; \theta) = \int_{x_1}^{x_2} P(x; \theta) dx \geq 1 - \alpha$$

- Perform experiment
- Report confidence interval: $[\theta_1, \theta_2]$ where $x_{\text{obs}} \in [x_1(\theta), x_2(\theta)]$ for all $\theta \in [\theta_1, \theta_2]$

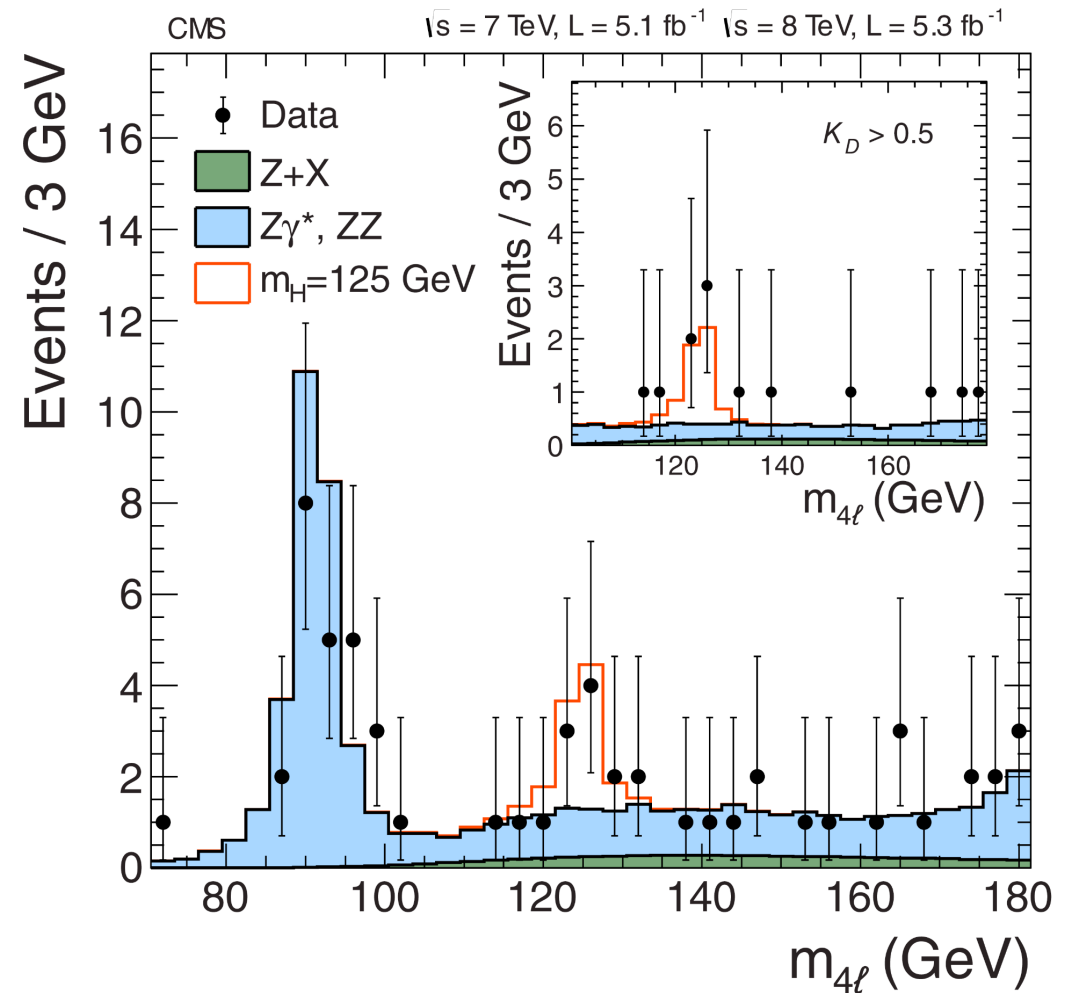
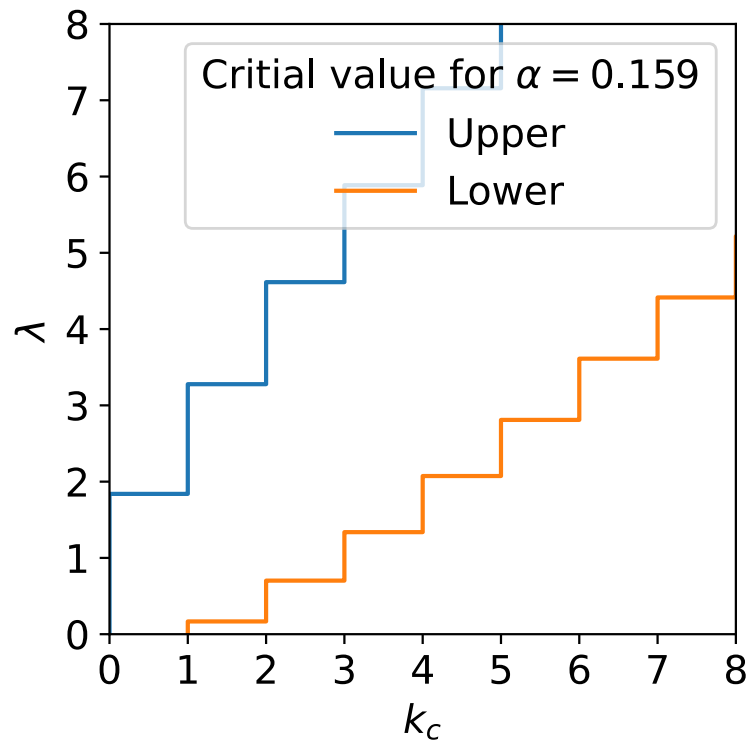
- Interval has coverage $1-\alpha$.

- For an ensemble of experiments, the interval $[\theta_1, \theta_2]$ will contain (unknown) θ_{true} with probability $1-\alpha$. This is a statement about the distribution of θ_1 and θ_2 , NOT θ_{true} .



Neyman interval

- You have all seen these: error bars on data points are the Neyman intervals for a Poisson distribution
 - With $\alpha = 0.159\dots$
 - Also referred to as Garwood intervals



Likelihood ratio test statistic

1. Define $t_\theta(x) = -2 \ln \frac{\mathcal{L}(\theta)}{\mathcal{L}(\hat{\theta})} = -2 \ln \frac{f(x; \theta)}{f(x; \hat{\theta})} \geq 0$

2. Compute associated pdf (change of variables)

$$f(t_\theta; \theta') = \int \delta(t_\theta - t_\theta(x)) dP(x; \theta')$$

3. For each θ , find the critical value $t_{\theta,c}$ that covers

$$P(t_\theta < t_{\theta,c}) = \int_0^{t_{\theta,c}} f(t_\theta; \theta' = \theta) dt_\theta \geq 1 - \alpha$$

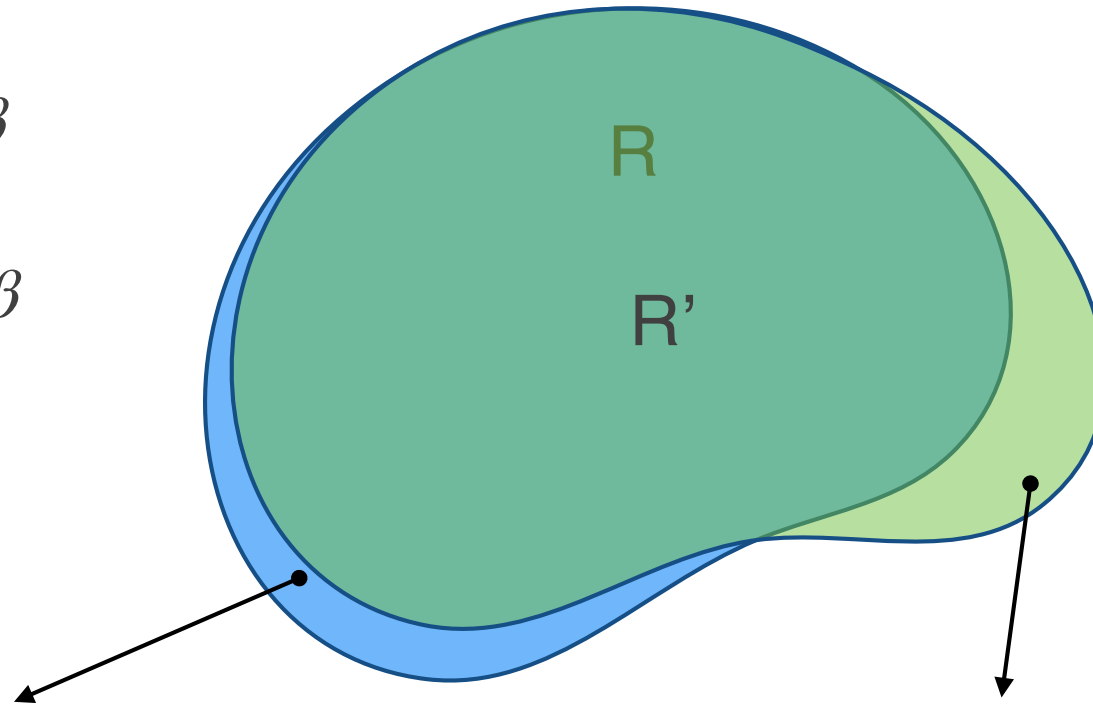
4. Perform experiment, get x_{obs} , report confidence set $\{\theta \mid t_\theta(x_{obs}) < t_{\theta,c}\}$

Again, the set is the random variate, and will contain (unknown) θ_{true} with probability at least $1-\alpha$. Dimension of θ and x are arbitrary.

If θ is 1-d and t_θ is monotone, can make a [Feldman-Cousins](#) interval.

Neyman-Pearson lemma

- $R = \{x \mid t_\theta(x) \geq t_{\theta,c}\}$, or $\{x \mid f(x; \hat{\theta}) \geq e^{t_{\theta,c}/2} f(x; \theta) = cf(x; \theta)\}$
- $P(x \in R; \theta) = \alpha$
- $P(x \in R; \hat{\theta}) = 1 - \beta$
- $P(x \in R'; \theta) = \alpha$
- $P(x \in R'; \hat{\theta}) < 1 - \beta$
 - We reject θ less often



$$f(x; \hat{\theta}) \geq cf(x; \theta)$$

$$f(x; \hat{\theta}) < cf(x; \theta)$$

$$P(x \in R \setminus R'; \hat{\theta}) \geq cP(x \in R \setminus R'; \theta)$$

$$P(x \in R' \setminus R; \hat{\theta}) < cP(x \in R' \setminus R; \theta)$$

$$P(x \in R \setminus R'; \hat{\theta}) > P(x \in R' \setminus R; \hat{\theta})$$

The (log-)likelihood ratio test is the most powerful test for a given size

Frequentists...



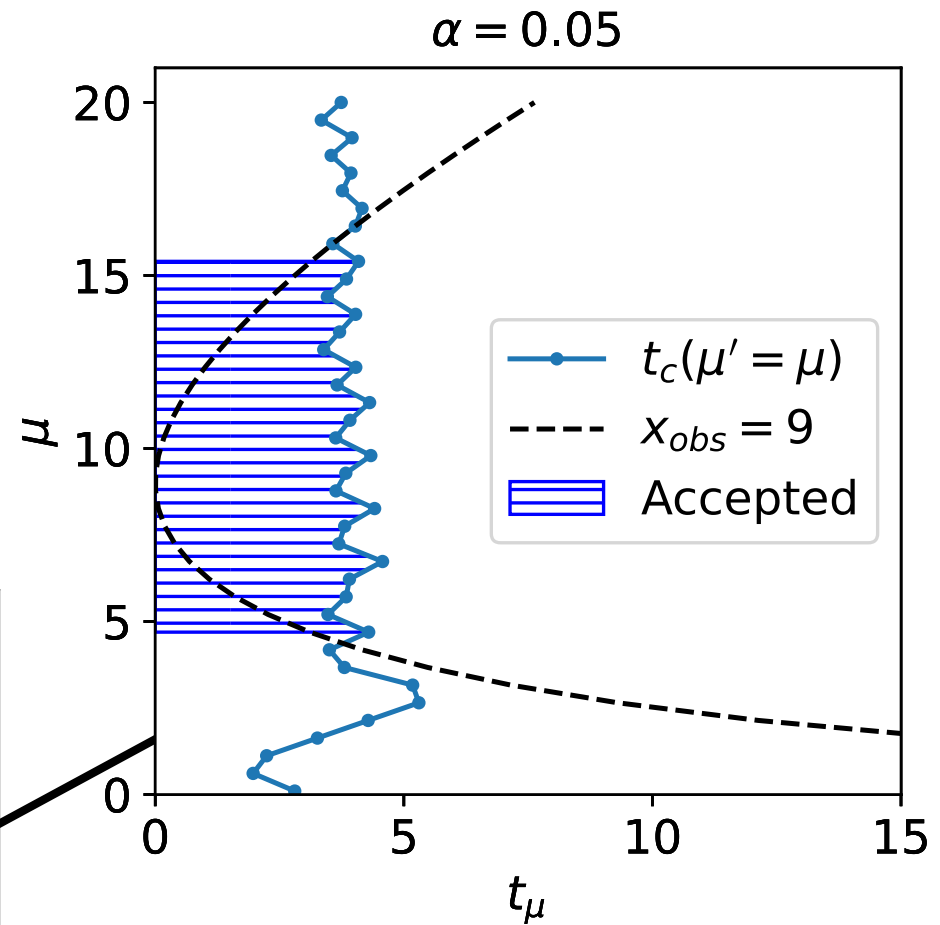
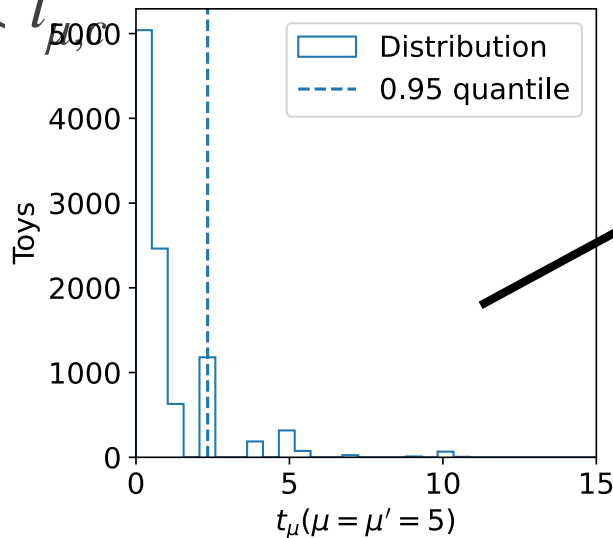
Likelihood ratio examples

- Poisson example

$$P(x; \mu) = \frac{\mu^x e^{-\mu}}{x!}$$

- For each μ :
 - Throw 10k toys
 - Compute $\hat{\mu}, t_\mu$
 - Find 0.95 quantile in distribution
- Draw $t_\mu(x_{obs})$ contour
- Accept $t_\mu(x_{obs}) < t_{\mu}$

Note: jagged behavior is due to discrete nature of t , not limited toy statistics



$$t_\mu = -2 \ln \frac{\mathcal{L}(\mu)}{\mathcal{L}(\hat{\mu})}$$

Likelihood ratio examples

- Poisson with background example

$$P(x; \mu s + b) = \frac{(\mu s + b)^x e^{-(\mu s + b)}}{x!}$$

- s=5, b=10 fixed, x=20

- Plan: set upper limit on μ

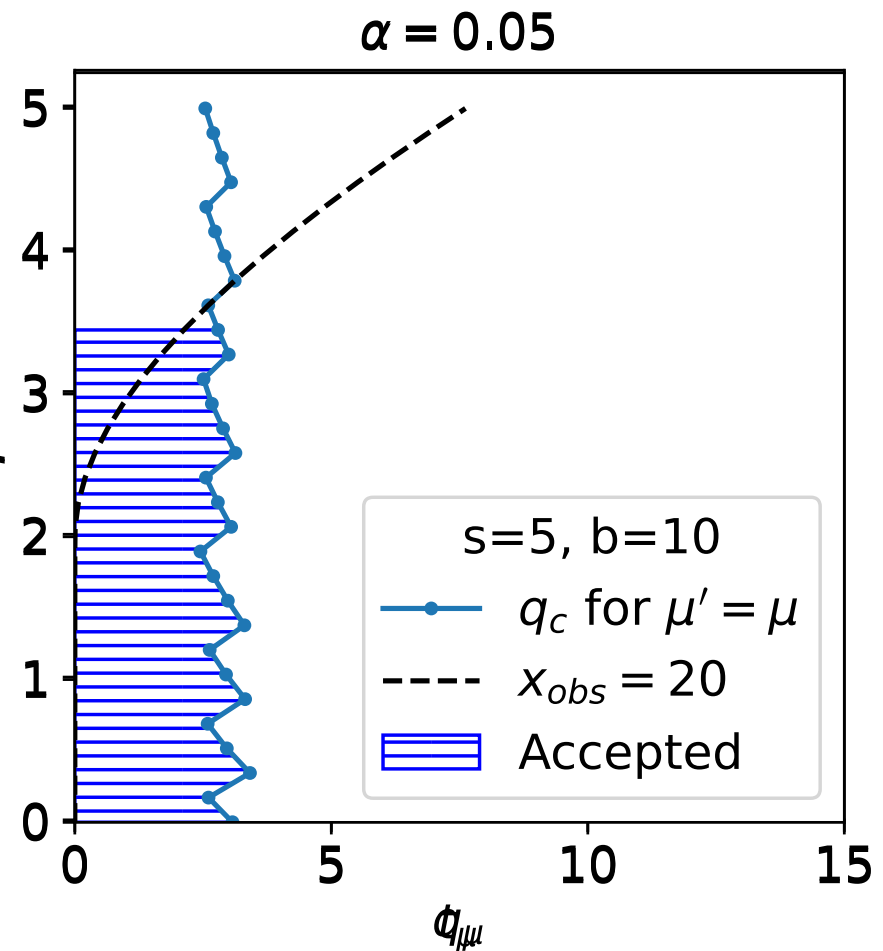
- Problem: two-sided region

- We should not consider $\hat{\mu} > \mu$ to indicate less compatibility with a model that assumes a rate μ .

- Solution: modify test statistic

- Define $q_\mu = -2 \ln \frac{\mathcal{L}(\mu)}{\mathcal{L}(\min(\mu, \hat{\mu}))}$

- i.e. over-fluctuations are “not extreme”



$$t_\mu = -2 \ln \frac{\mathcal{L}(\mu)}{\mathcal{L}(\hat{\mu})}$$

Likelihood ratio examples

- Poisson with background example

$$P(x | \mu s + b) = \frac{(\mu s + b)^x e^{-(\mu s + b)}}{x!}$$

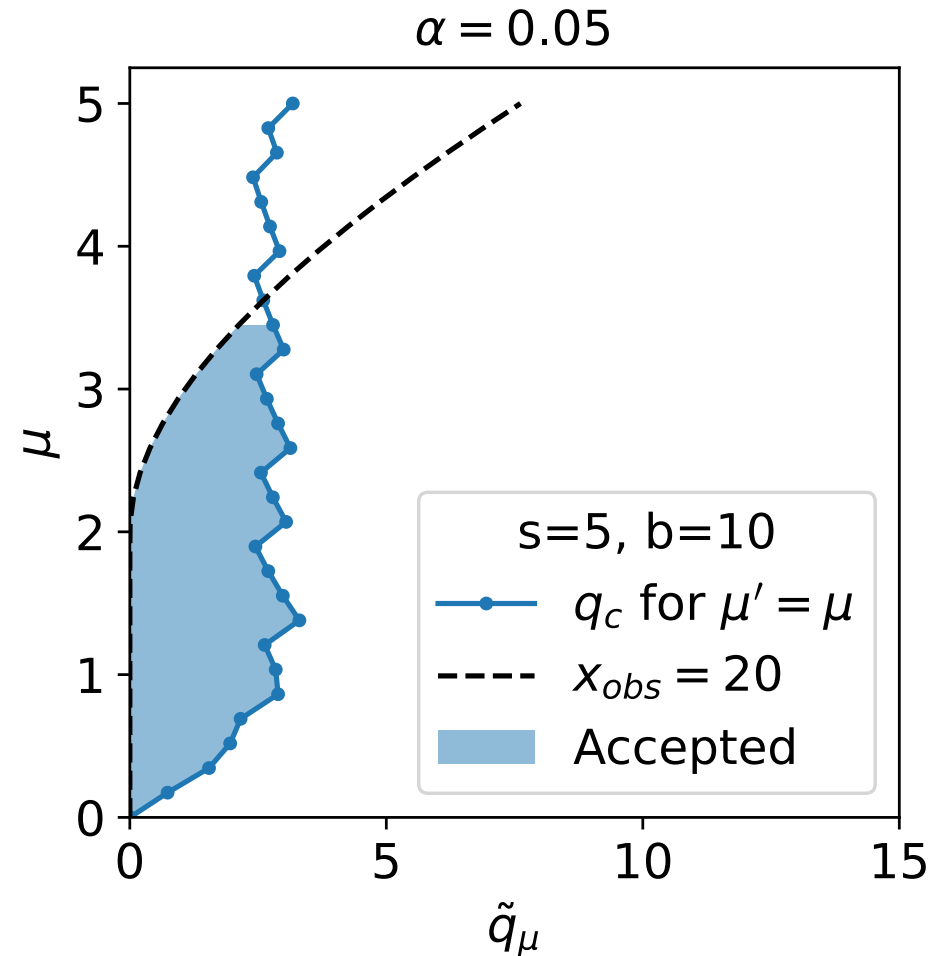
- s=5, b=10 fixed, x=20

- Problem: negative $\hat{\mu}$

- Test stat distribution at 0 should collapse

- Define $\tilde{q}_\mu = -2 \ln \frac{\mathcal{L}(\mu)}{\mathcal{L}(\max(0, \min(\mu, \hat{\mu})))}$

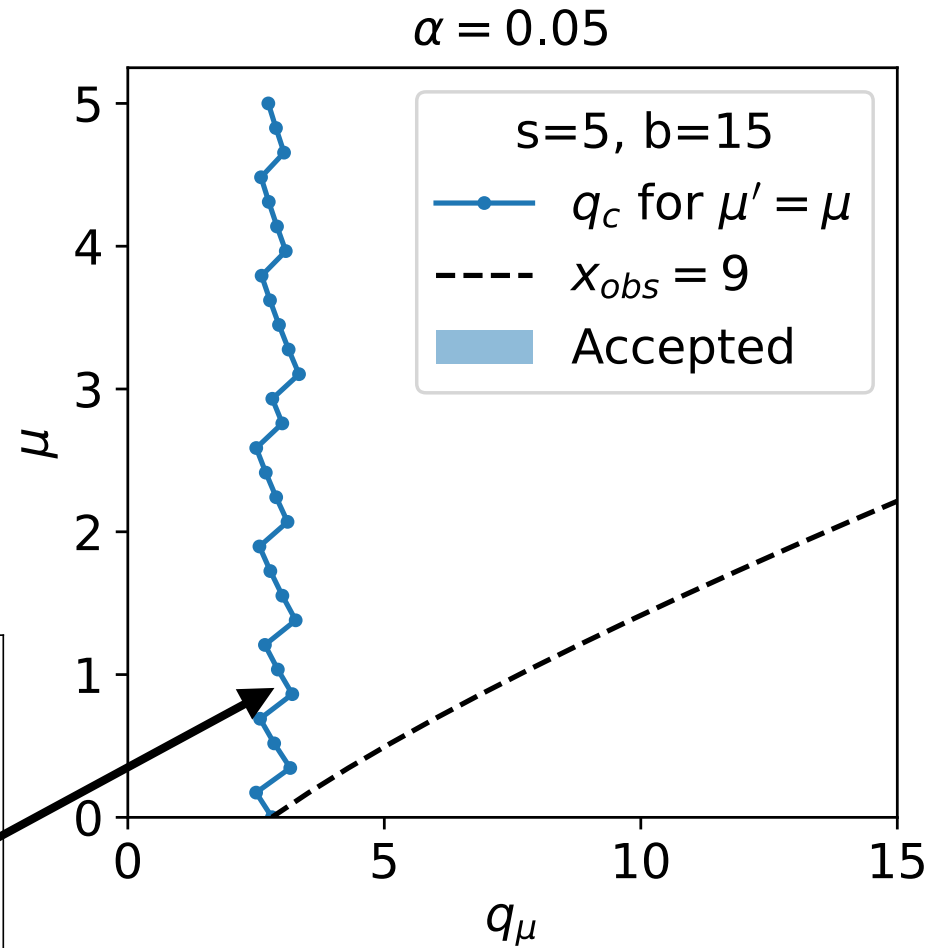
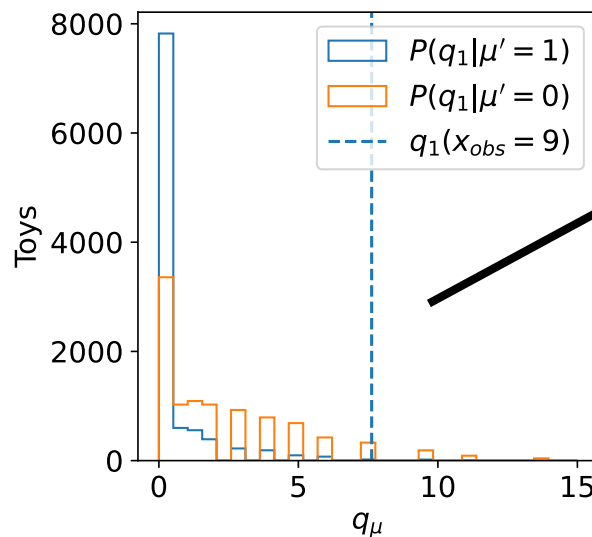
- i.e. restrict $\hat{\mu}$ to be positive



Likelihood ratio examples

- Same example as before, but $b=15$, $x=9$
- Problem: under fluctuation
 - No values accepted!
 - Possible but unsatisfying outcome of frequentist test

The result $x=9$ is rare for both S+B and B-only hypotheses.



$$q_\mu = -2 \ln \frac{\mathcal{L}(\mu)}{\mathcal{L}(\min(\mu, \hat{\mu}))}$$

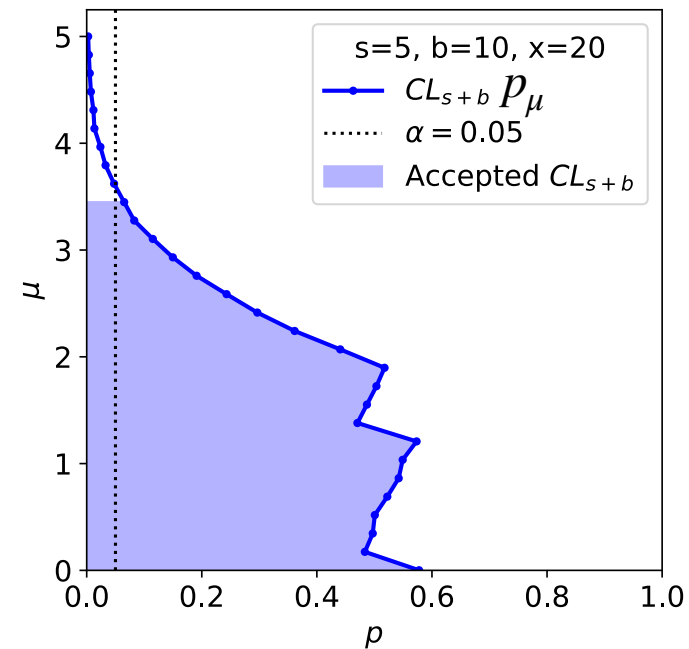
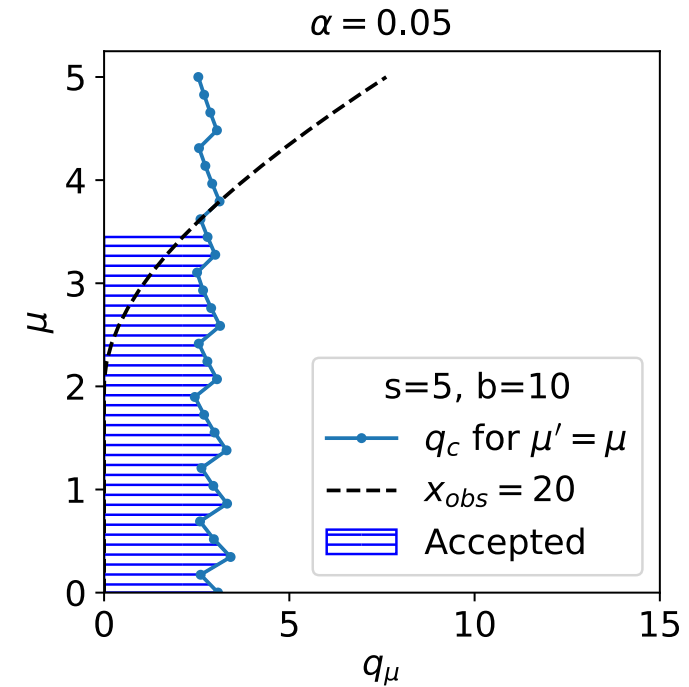
CLs

- CL_s criterion departs from purely frequentist CL to ameliorate the null set problem (among others)
 - Original expositions by [A. Read](#), [T. Junk](#)
 - See also [PDG 40.4.2.4](#)

- First we reformulate our old test:

- Define $p_\mu = \int_{t_\mu(x_{obs})}^{\infty} P(t_\mu | \mu' = \mu) dt_\mu$

- This is a p-value
- Then we accept the region $p_\mu > \alpha$
 - Right: initial S+B example reformulated



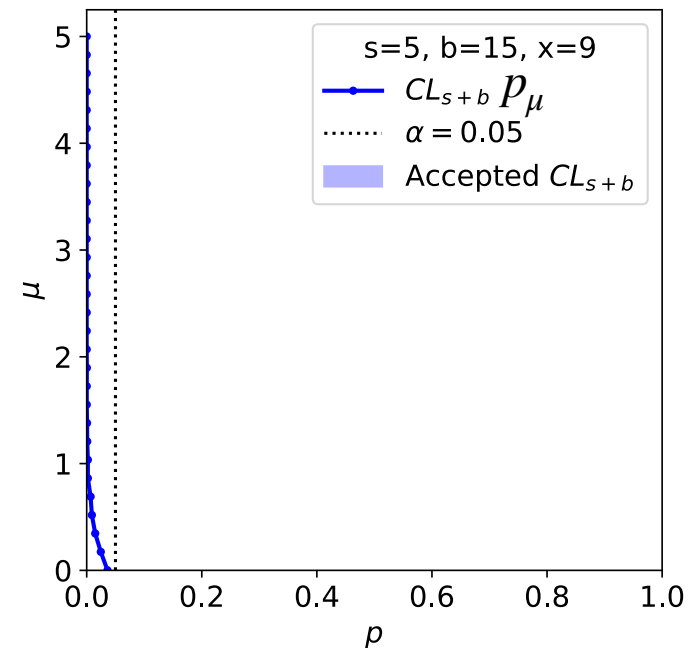
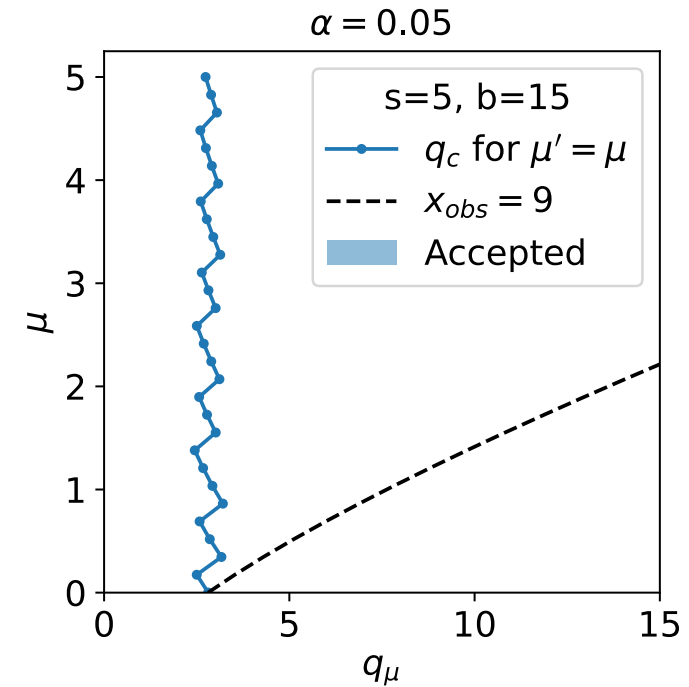
CLs

- CL_s criterion departs from purely frequentist CL to ameliorate the null set problem (among others)
 - Original expositions by [A. Read](#), [T. Junk](#)
 - See also [PDG 40.4.2.4](#)

- First we reformulate our old test:

- Define $p_\mu = \int_{t_\mu(x_{obs})}^{\infty} P(t_\mu | \mu' = \mu) dt_\mu$

- This is a p-value
- Then we accept the region $p_\mu > \alpha$
 - Right: under-fluctuation S+B example reformulated



CLs

- CL_s criterion departs from purely frequentist CL to ameliorate the null set problem (among others)
 - Original expositions by [A. Read](#), [T. Junk](#)
 - See also [PDG 40.4.2.4](#)

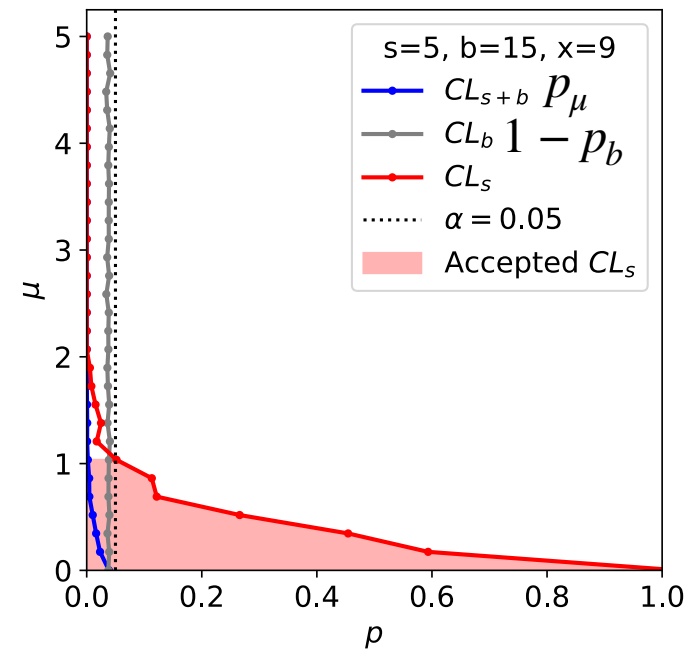
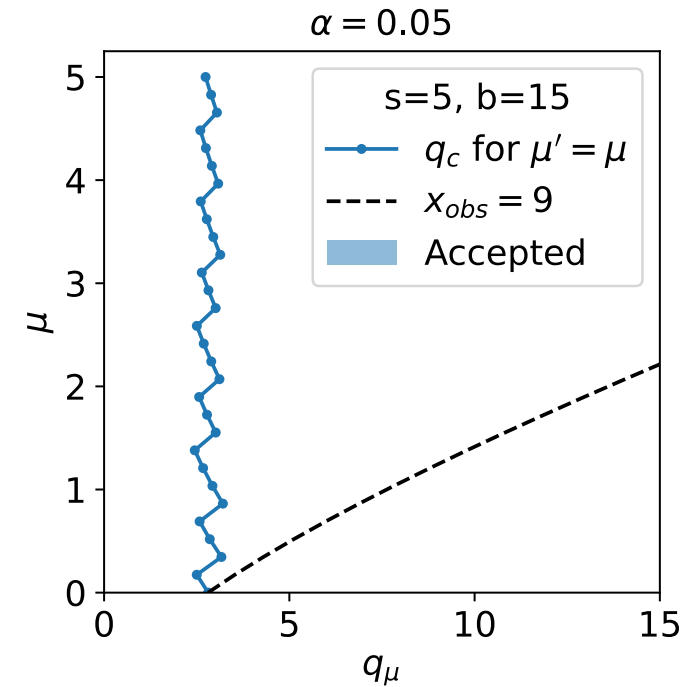
- First we reformulate our old test:

- Define $p_\mu = \int_{t_\mu(x_{obs})}^{\infty} P(t_\mu | \mu' = \mu) dt_\mu$

- Now define background-only p-value

- $1 - p_b = \int_{t_\mu(x_{obs})}^{\infty} P(t_\mu | \mu' = 0) dt_\mu$

- Accept instead $CL_s = \frac{p_\mu}{1 - p_b} > \alpha$



CLs

- CL_s criterion departs from purely frequentist CL to ameliorate the null set problem (among others)
 - Original expositions by [A. Read](#), [T. Junk](#)
 - See also [PDG 40.4.2.4](#)

- First we reformulate our old test:

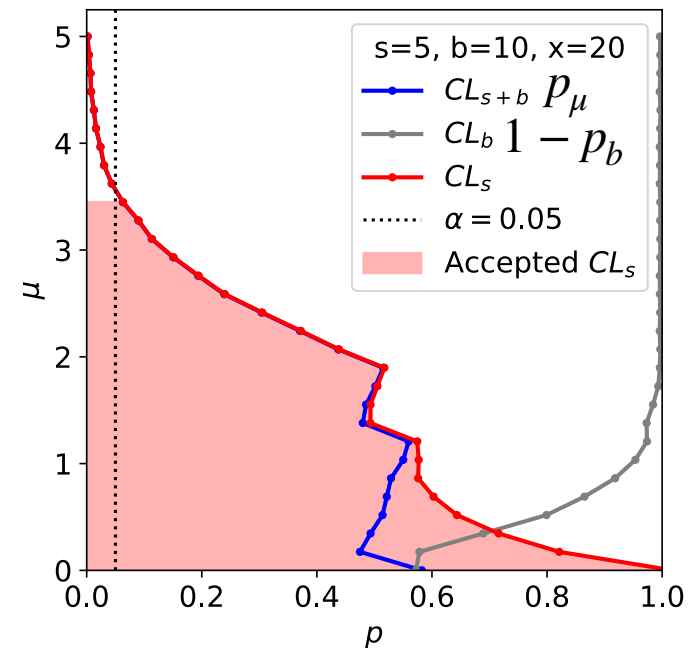
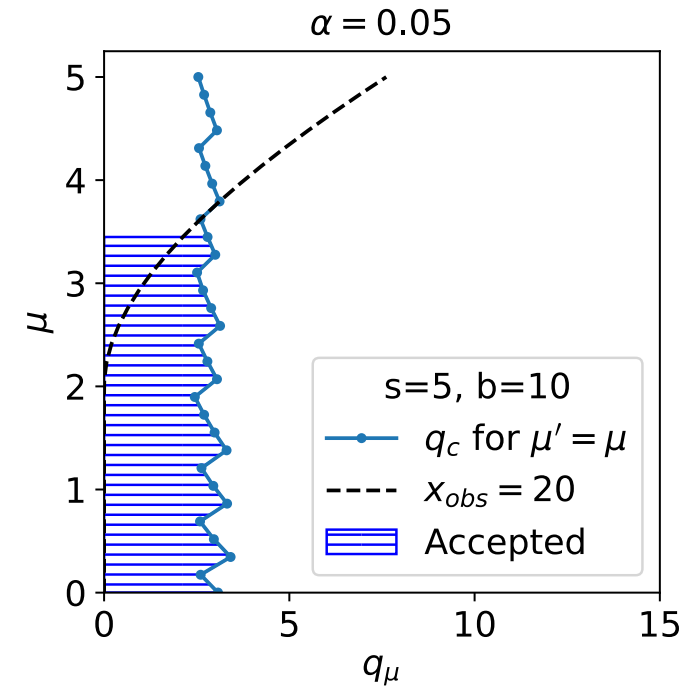
- Define $p_\mu = \int_{t_\mu(x_{obs})}^{\infty} P(t_\mu | \mu' = \mu) dt_\mu$

- Now define background-only p-value

- $1 - p_b = \int_{t_\mu(x_{obs})}^{\infty} P(t_\mu | \mu' = 0) dt_\mu$

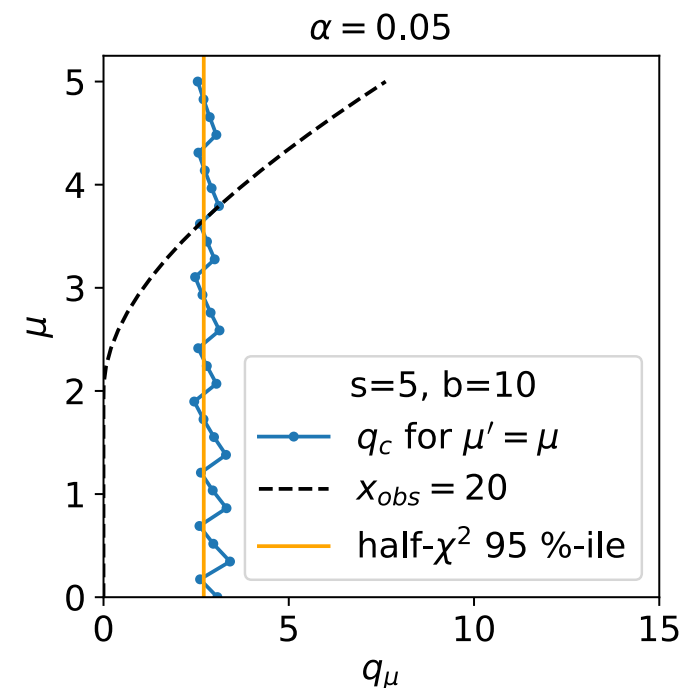
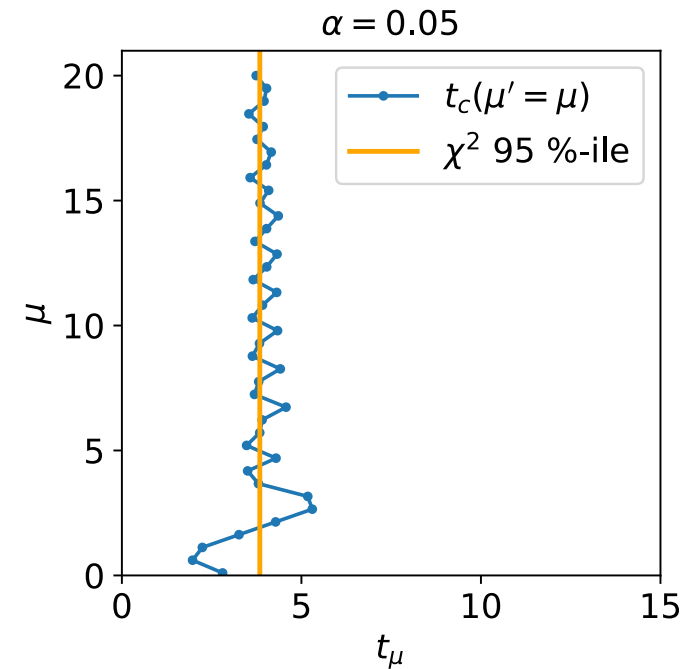
- Accept instead $CL_s = \frac{p_\mu}{1 - p_b} > \alpha$

- No effect in the first example



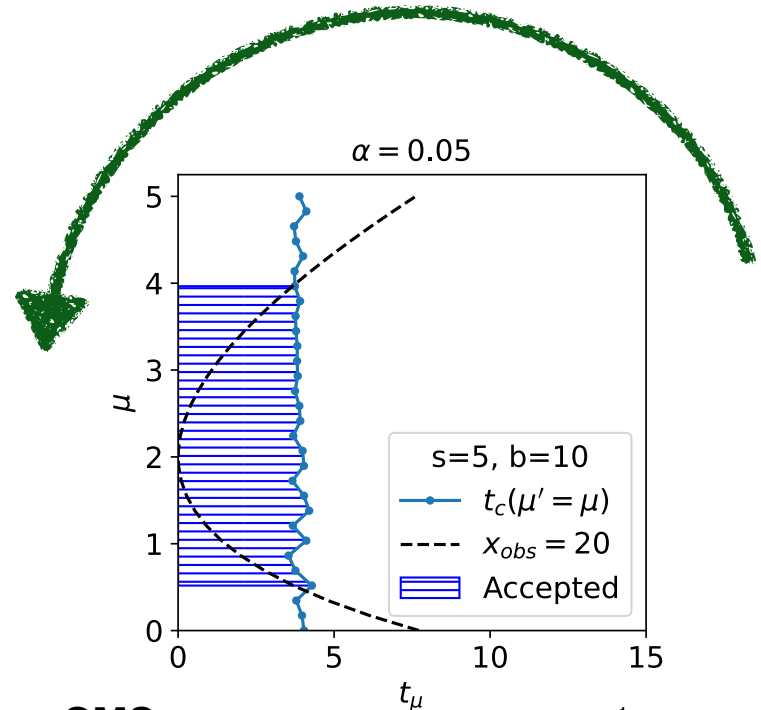
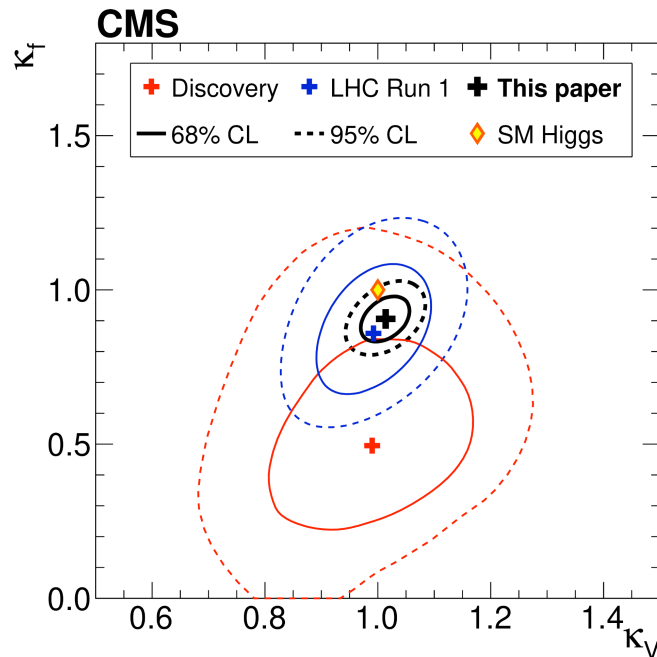
Asymptotic behavior

- Notice how $t_{\mu,c}$ and $q_{\mu,c}$ tend towards a constant?
- This is [Wilk's theorem](#) in action
 - Statement: as sample size grows, the distribution of the likelihood ratio $P(t_{\theta|\theta'})$ approaches a χ^2 distribution
 - With $df = \dim(\theta)$
 - Hence we can approximate by just evaluating $t_{\theta}(x_{obs})$!
- For q , formulas slightly more complex
 - [CCGV](#) provide the recipe: non-central half- χ^2
 - The non-centrality is found using the *Asimov* dataset
 - A special x_{μ} for a given μ such that $\hat{\mu}(x) = \mu$
 - Note for Poisson data, it may be non-integral!
 - This dataset produces the median expected limit



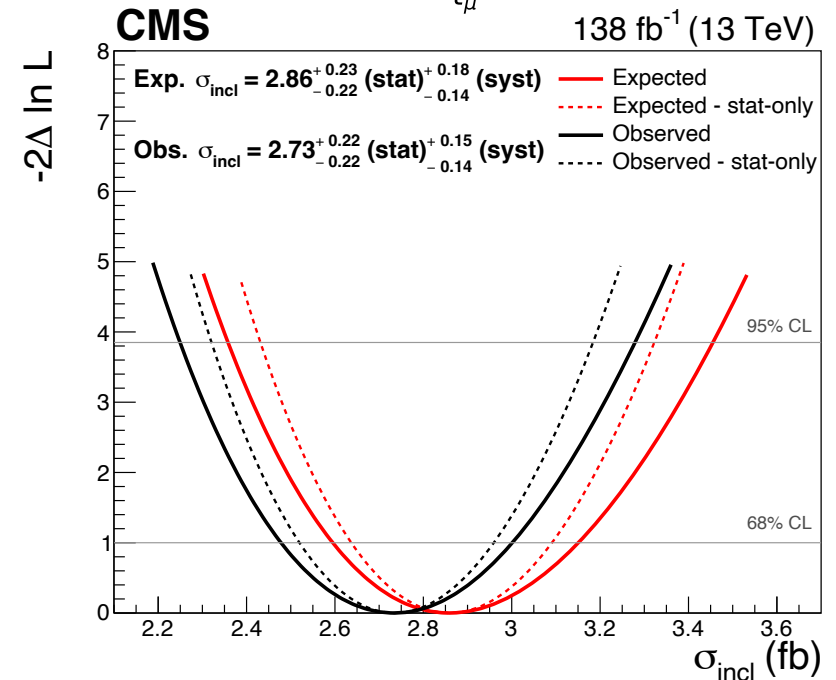
Asymptotic behavior

- This is how we make deltaNLL contours

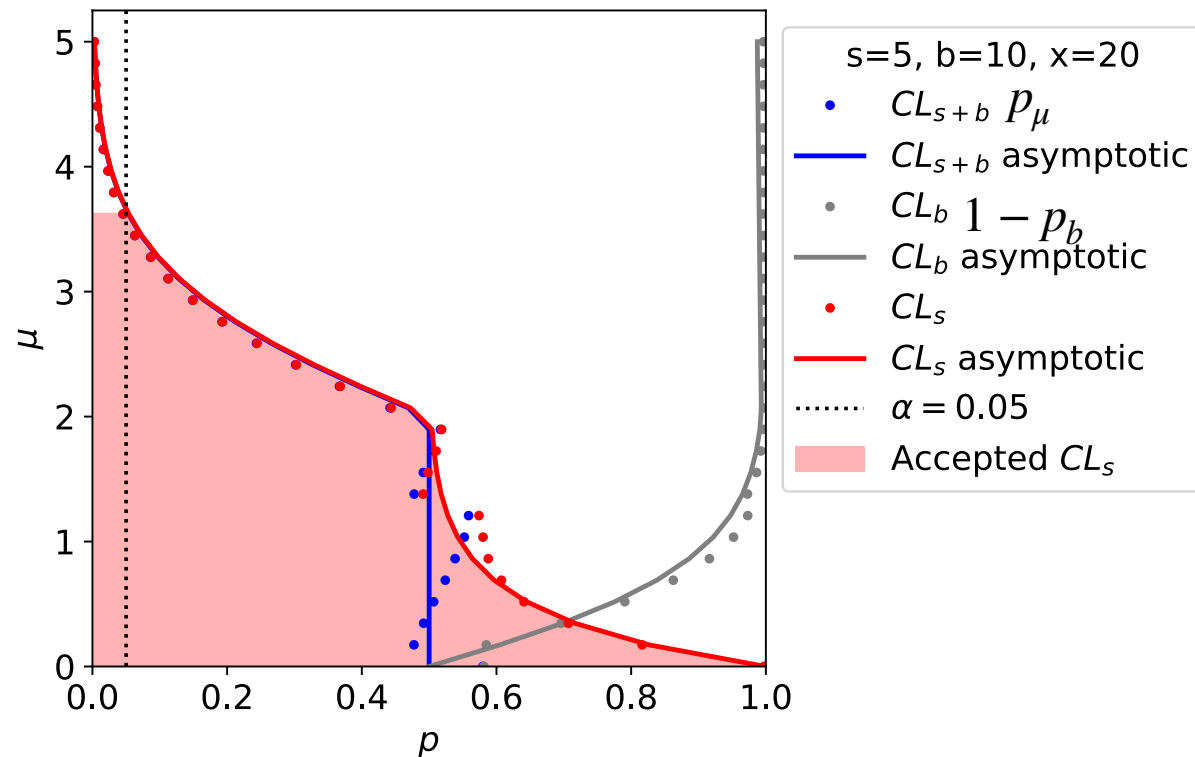
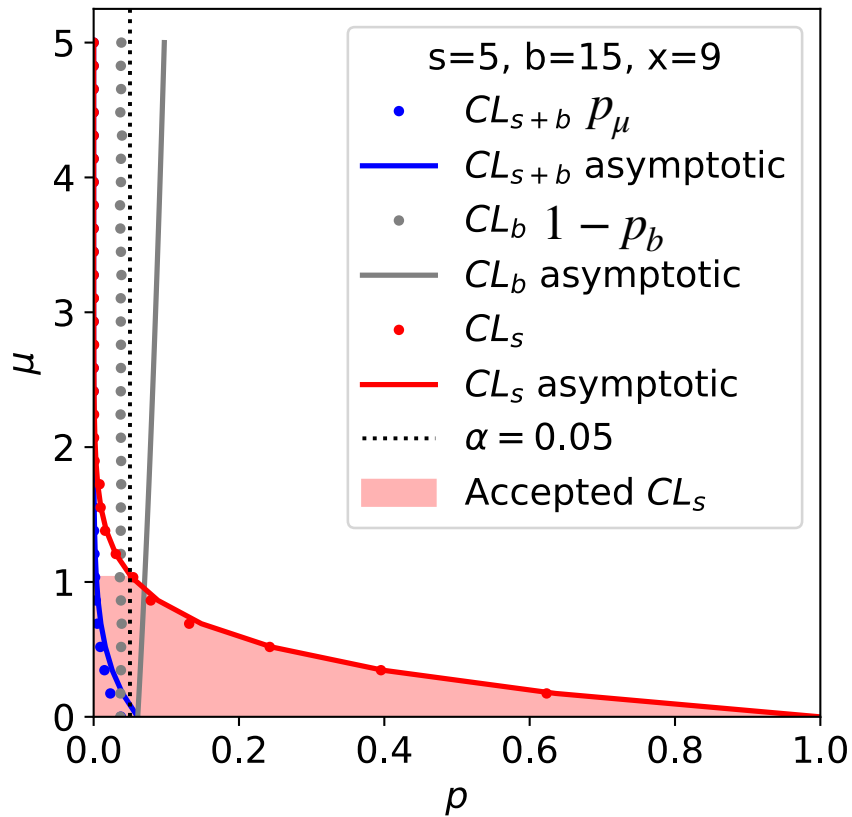


`scipy.stats.chi2.ppf(q, df)`

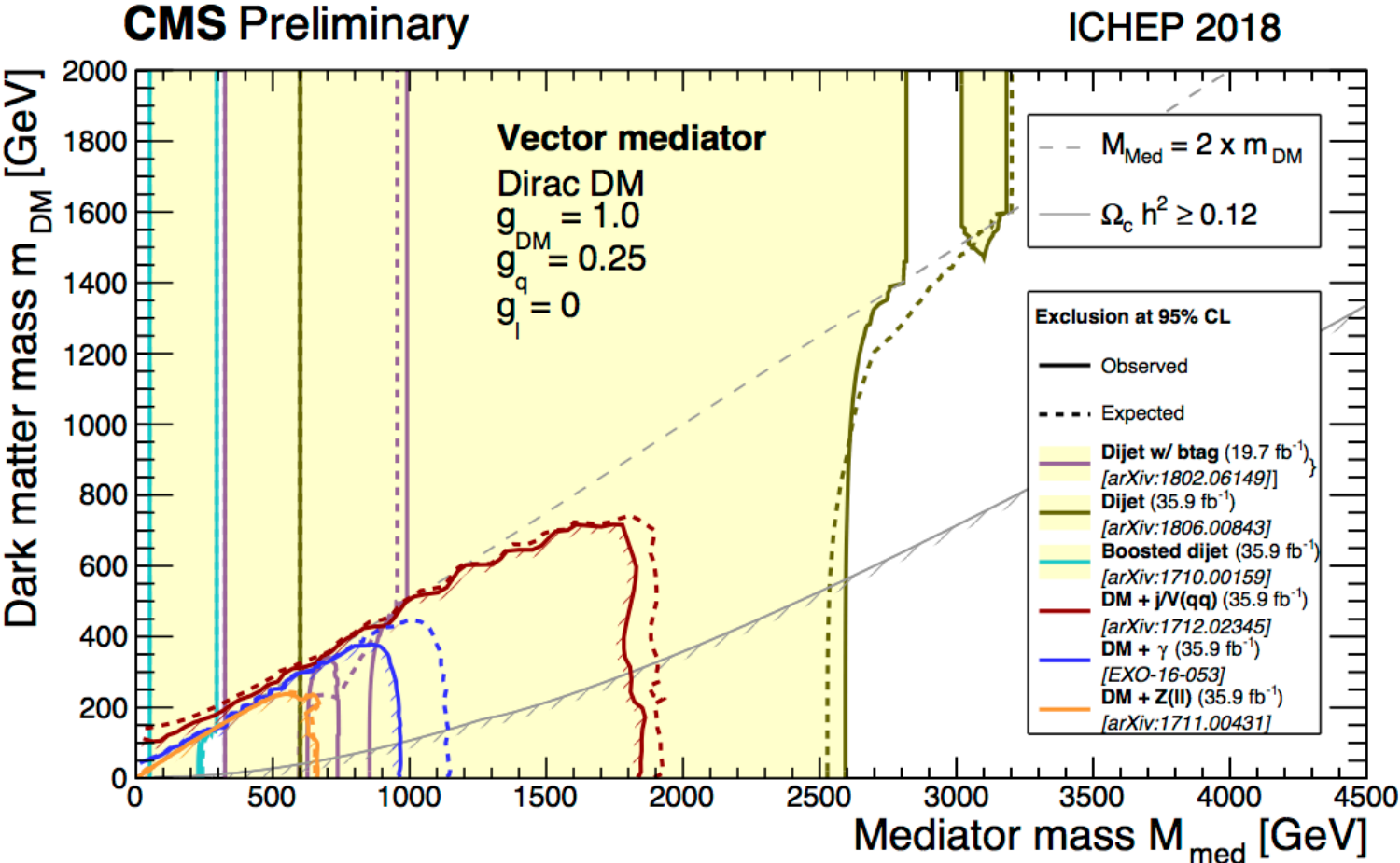
Quantile	t_c	
	$df=1$	$df=2$
0.68	0.989	2.279
1σ (0.6827...)	1	2.296
0.95	3.841	5.991
2σ (0.9545...)	4	6.180



Asymptotic CLs



Examples



Test statistic for discovery

- Poisson with background example

$$f_P(x; \mu s + b) = \frac{(\mu s + b)^x e^{-(\mu s + b)}}{x!}$$

- $s=20$, $b=15$ fixed, $x=39$

- Cannot use t_μ :

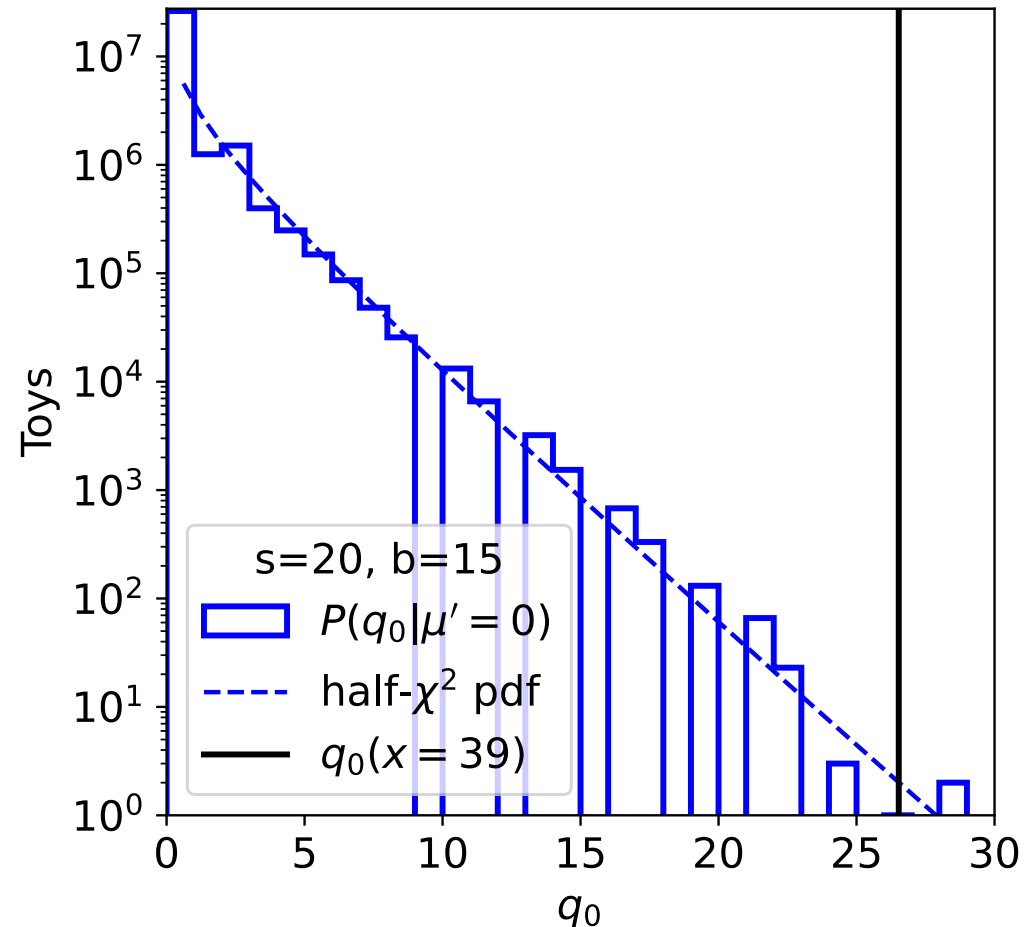
- Severe under-fluctuation would count as discovery! Certainly something was discovered, but not an excess over background. Disallow in test statistic:

- Define $q_0 = -2 \ln \frac{\mathcal{L}(0)}{\mathcal{L}(\max(0, \hat{\mu}))}$

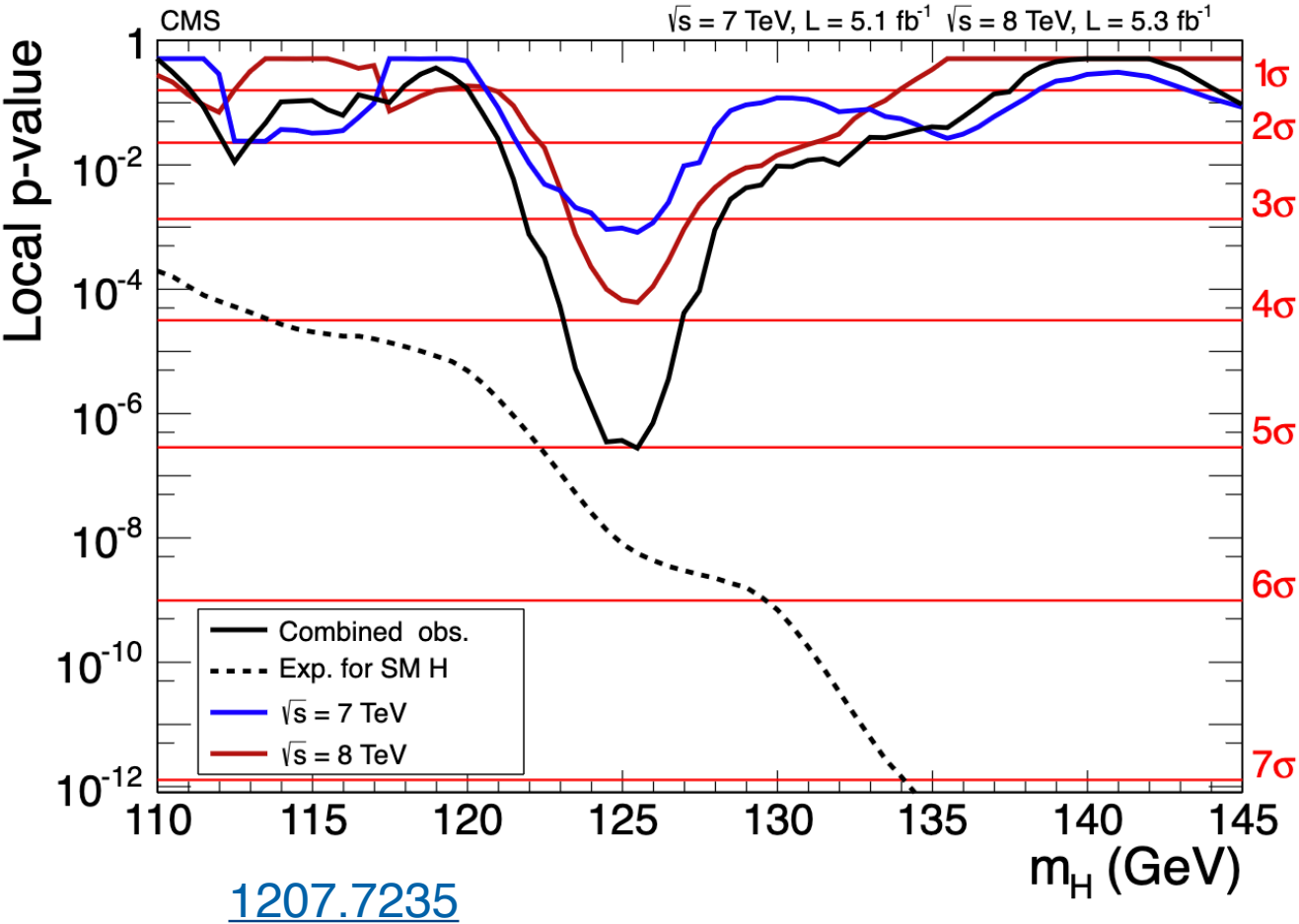
- i.e. under-fluctuations are “not extreme”

- Deceptively simple result: $Z = \sqrt{q_0(x_{obs})}$

- Only true if one POI

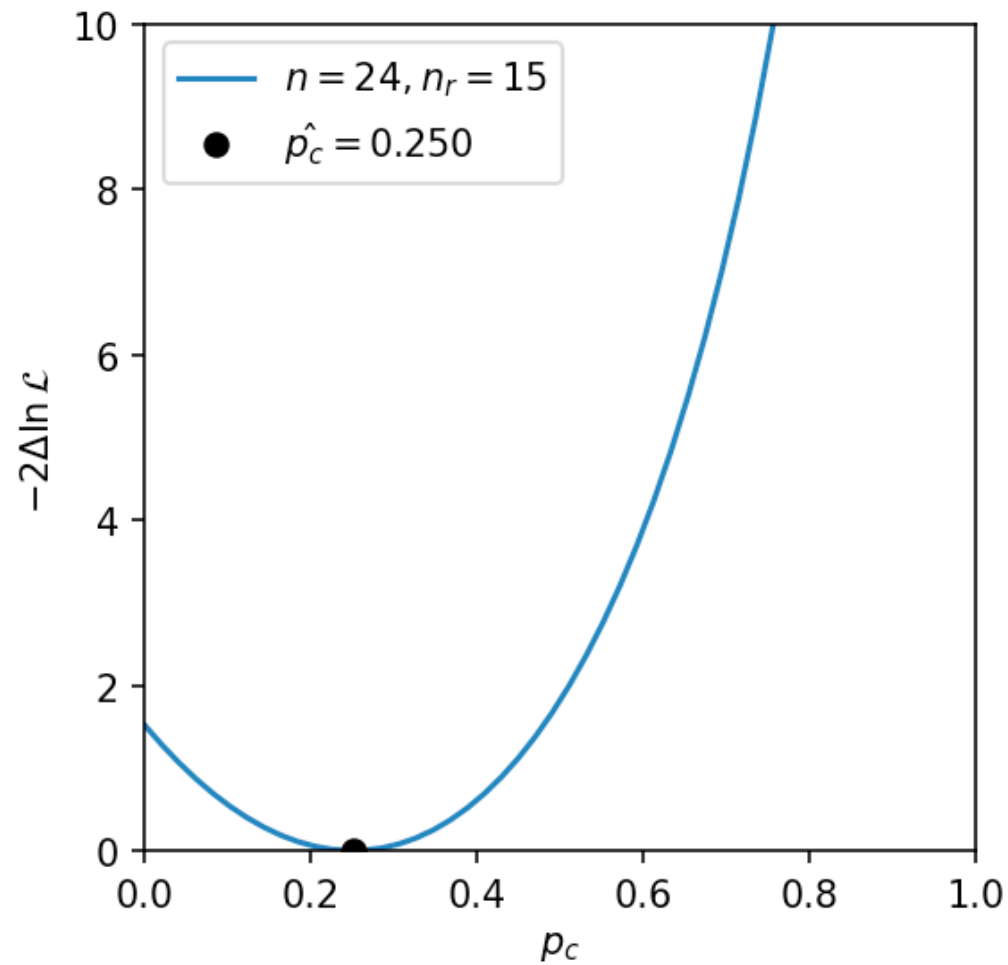


Discovery example



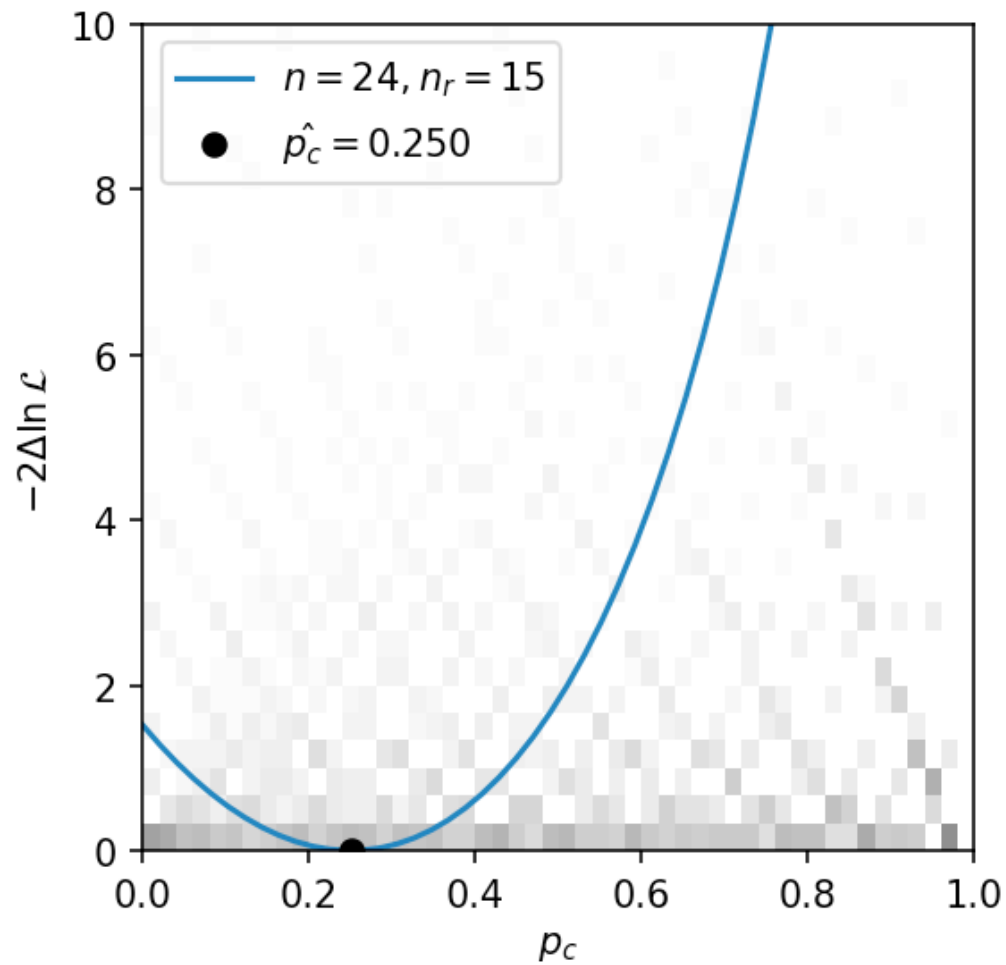
Intervals for $P(\text{cheat})$

- Coming back to our favorite problem



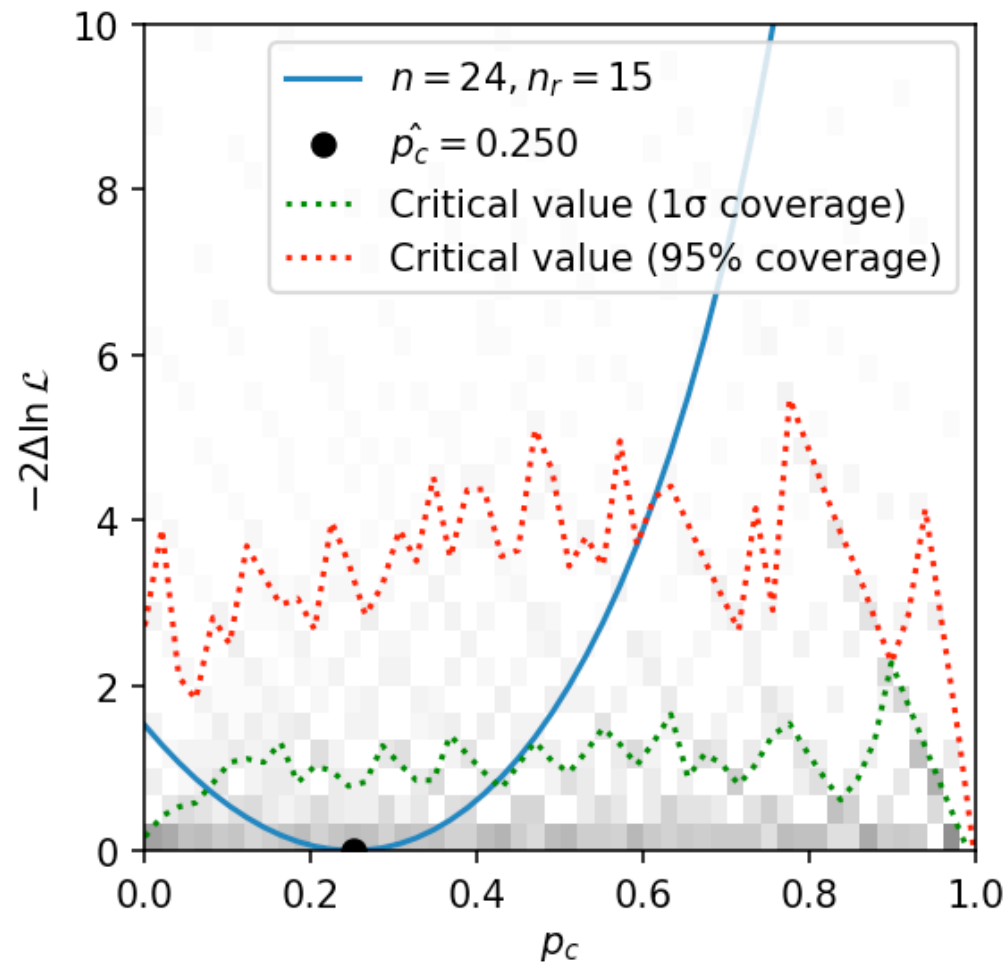
Intervals for P(cheat)

- Coming back to our favorite problem
 - For no-cheat null hypothesis, $p_0 \approx 0.2$, Z score: 0.88σ



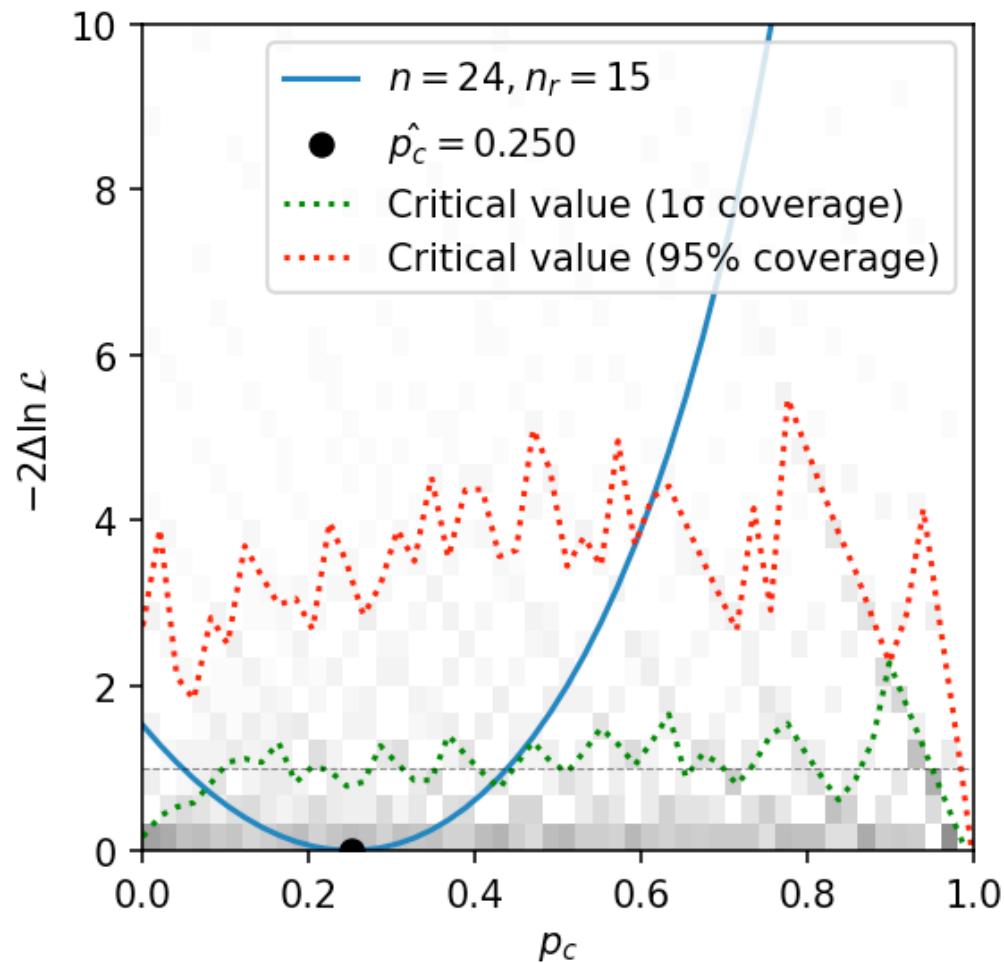
Intervals for $P(\text{cheat})$

- Coming back to our favorite problem
 - For no-cheat null hypothesis, $p_0 \approx 0.2$, Z score: 0.88σ
 - 1σ central interval: $0.08 < p_c < 0.41$



Intervals for $P(\text{cheat})$

- Coming back to our favorite problem
 - For no-cheat null hypothesis, $p_0 \approx 0.2$, Z score: 0.88σ
 - 1σ central interval: $0.08 < p_c < 0.41$



Intervals for $P(\text{cheat})$

- Coming back to our favorite problem
 - For no-cheat null hypothesis, $p_0 \approx 0.2$, Z score: 0.88σ
 - 1σ central interval: $0.08 < p_c < 0.41$

