



Introduction to Statistical Inference



Giacomo D'Amico

ArQus School 2022, Bergen, Norway

5-9 Sep. 2022



What do we mean by inferring?

Definition of **infer** verb from the Oxford Learner's Dictionary of Academic English



infer *verb*



BrE /ɪn'fɜ:(r)/; NAmE /ɪn'fɜ:r/

+ Verb Forms

to reach an opinion or decide that something is true on the basis of information that is available

What do we mean by inferring?

Definition of **infer** verb from the Oxford Learner's Dictionary of Academic English



infer *verb*

OPAL
written

BrE /ɪn'fɜ:(r)/; NAmE /ɪn'fɜ:r/

+ Verb Forms

to reach an **opinion** or decide that **something** is true on the basis of **information** that is available

An opinion that has to be *quantified* through the instrument of **probability** and **statistics**

A given theoretical model

The data we have collected

The Model



All sheep are white

The data



The opinion



The model is rejected

The Model



1% of the sheep are black

The data



The opinion



?

We will come back
later on this!

Two approaches are used to **quantify** an *opinion* about a **model** given an **observation**

- The **Bayesian approach** tries to answer the question:

*Given our **prior** knowledge and the observed **data**, what is the **probability** that the model is true?*

- The **Frequentist approach** tries to answer the question:

*If I repeat the experiment an infinite time, assuming the model is true, with which **frequency** I would observe a value more **extreme** than the one actually observed?*

The Bayesian approach

The Bayes theorem

- Marginalised probability

$$f(x|I) = \int f(x, y|I) dy$$

- Conditional probability

$$f(x, y|I) = f(x|y, I) \cdot f(y|I)$$

- I represents our prior knowledge
- $f()$ is for a generic probability distribution (or mass) function

The Bayes theorem

- Marginalised probability

$$f(x|I) = \int f(x, y|I) dy$$

- Conditional probability

$$f(x, y|I) = f(x|y, I) \cdot f(y|I)$$

- I represents our prior knowledge
- $f()$ is for a generic probability distribution (or mass) function

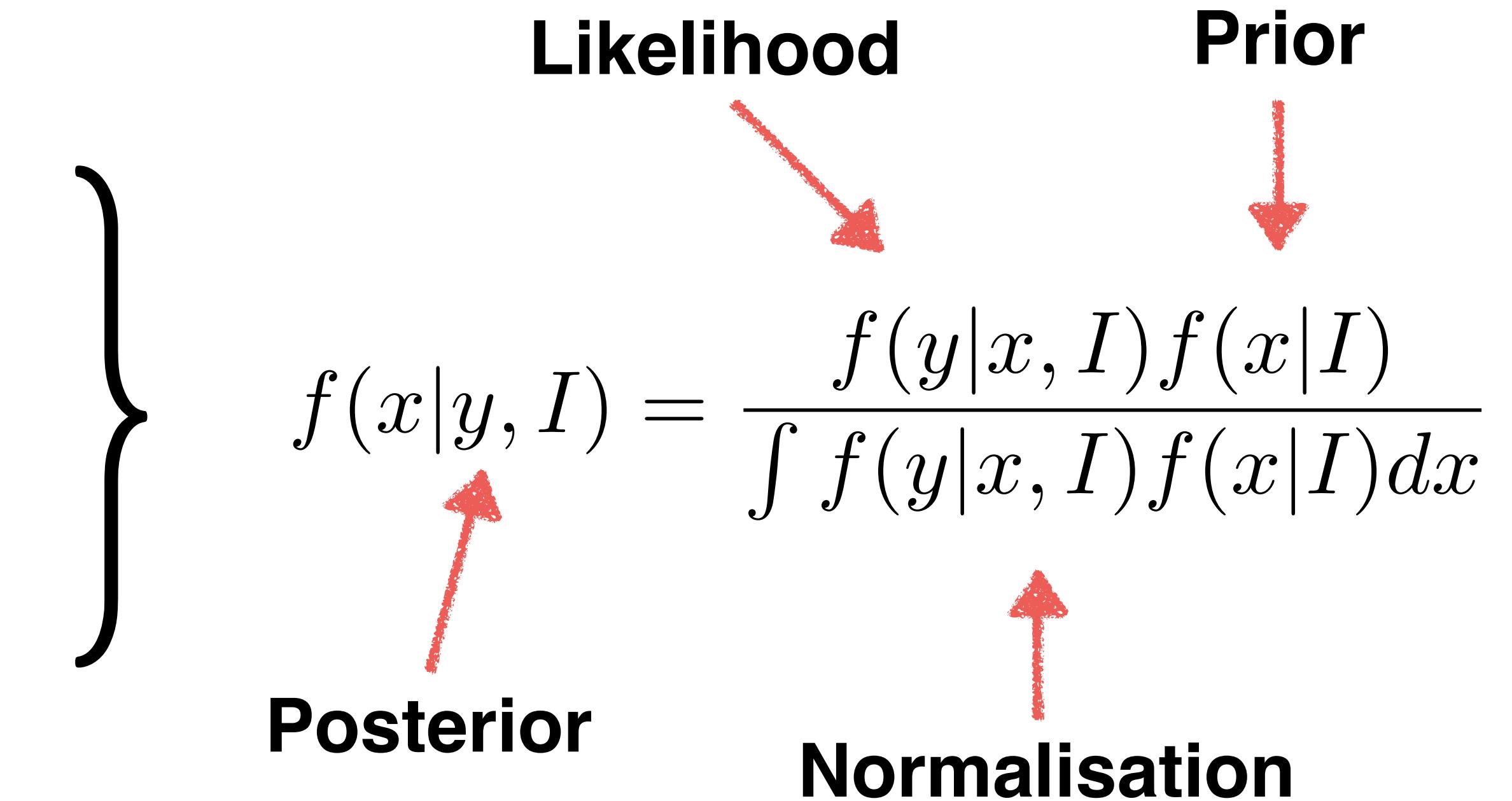
The Bayes theorem

- Marginalised probability

$$f(x|I) = \int f(x, y|I) dy$$

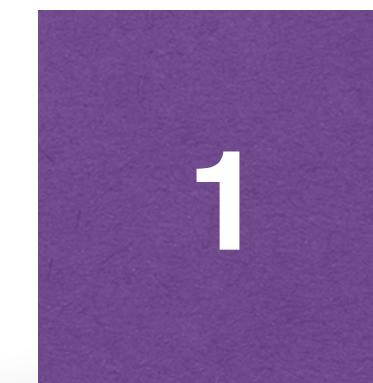
- Conditional probability

$$f(x, y|I) = f(x|y, I) \cdot f(y|I)$$



- I represents our prior knowledge
- $f()$ is for a generic probability distribution (or mass) function

The Monty Hall problem



In two boxes there is a goat and in the other a car

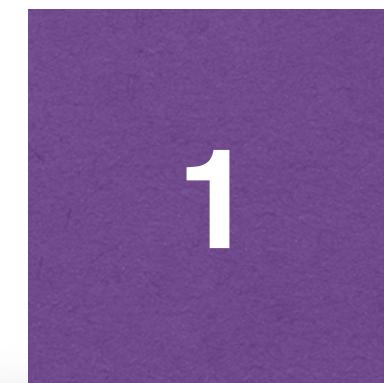
You have to choose one and only one box

The Monty Hall problem



Imagine we randomly pick the first one, but without opening it

The Monty Hall problem



Now the host of the game (who knows where the car is) shows us the content of the third box, which does not contain the car

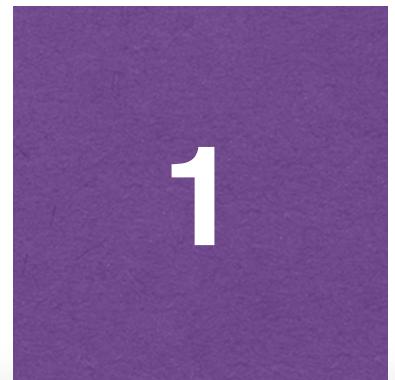
The Monty Hall problem



S/He then give us the opportunity to change our box (n.1) with the other (n. 2)

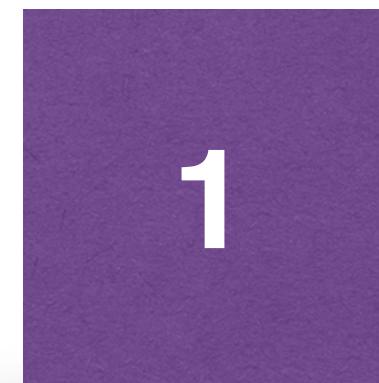
What would you do? Would you accept the opportunity?

The Monty Hall problem



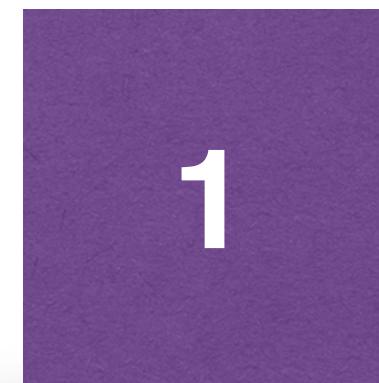
- H_i The hypothesis “the car is in the i-th box”

The Monty Hall problem



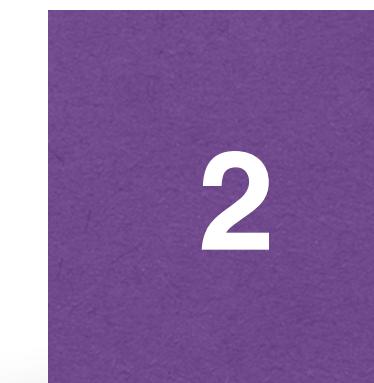
- H_i The hypothesis “the car is in the i-th box”
- E The event “the host shows use the content of the third box”

The Monty Hall problem

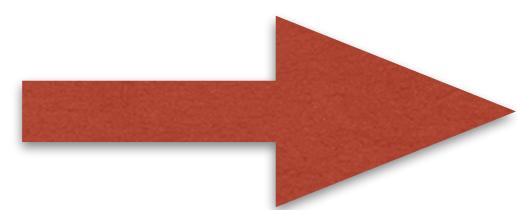


- H_i The hypothesis “the car is in the i-th box”
- E The event “the host shows use the content of the third box”
- I Our prior knowledge
“3 boxes and 1 car” \oplus “the host knows where the car is”

The Monty Hall problem



- H_i The hypothesis “the car is in the i-th box”
- E The event “the host shows use the content of the third box”
- I Our prior knowledge
“3 boxes and 1 car” \oplus “the host knows where the car is”



Posterior

$$f(H_i | E, I)$$

The Monty Hall problem

1

2



$$f(H_1|E, I) = \frac{f(E|H_1, I)f(H_1|I)}{f(E|I)} = \dots$$

$$f(H_2|E, I) = \frac{f(E|H_2, I)f(H_2|I)}{f(E|I)} = \dots$$

$$f(H_3|E, I) = \frac{f(E|H_3, I)f(H_3|I)}{f(E|I)} = \dots$$

The Monty Hall problem

1

2



$$f(H_1|E, I) = \frac{f(E|H_1, I)f(H_1|I)}{f(E|I)} = \underline{\hspace{2cm}} \cdot 1/3$$

$$f(H_2|E, I) = \frac{f(E|H_2, I)f(H_2|I)}{f(E|I)} = \underline{\hspace{2cm}} \cdot 1/3$$

$$f(H_3|E, I) = \frac{f(E|H_3, I)f(H_3|I)}{f(E|I)} = \underline{\hspace{2cm}} \cdot 1/3$$

Priors \rightarrow $f(H_1|I) = f(H_2|I) = f(H_3|I) = \frac{1}{3}$

The Monty Hall problem

1

2



$$f(H_1|E, I) = \frac{f(E|H_1, I)f(H_1|I)}{f(E|I)} = \frac{\cdot 1/3}{1/2}$$

$$f(H_2|E, I) = \frac{f(E|H_2, I)f(H_2|I)}{f(E|I)} = \frac{\cdot 1/3}{1/2}$$

$$f(H_3|E, I) = \frac{f(E|H_3, I)f(H_3|I)}{f(E|I)} = \frac{\cdot 1/3}{1/2}$$

Normalisation $\rightarrow \sum_i f(E|H_i, I)f(H_i|I) = f(E|I) = \frac{1}{2}$

The Monty Hall problem

1

2



$$f(H_1|E, I) = \frac{f(E|H_1, I)f(H_1|I)}{f(E|I)} = \frac{1/2 \cdot 1/3}{1/2} =$$

$$f(H_2|E, I) = \frac{f(E|H_2, I)f(H_2|I)}{f(E|I)} = \frac{1 \cdot 1/3}{1/2} =$$

$$f(H_3|E, I) = \frac{f(E|H_3, I)f(H_3|I)}{f(E|I)} = \frac{0 \cdot 1/3}{1/2} =$$

Likelihoods → $f(E|H_1, I) = \frac{1}{2}$ $f(E|H_2, I) = 1$ $f(E|H_3, I) = 0$

Il “paradosso” di Monty Hall

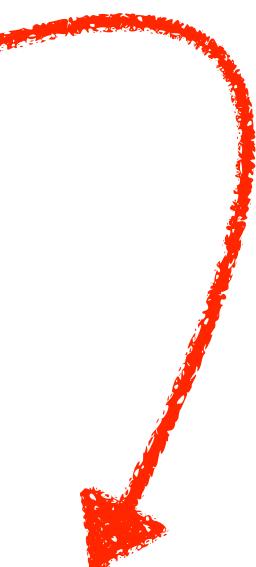


$$f(H_1|E, I) = \frac{f(E|H_1, I)f(H_1|I)}{f(E|I)} = \frac{1/2 \cdot 1/3}{1/2} = \frac{1}{3}$$

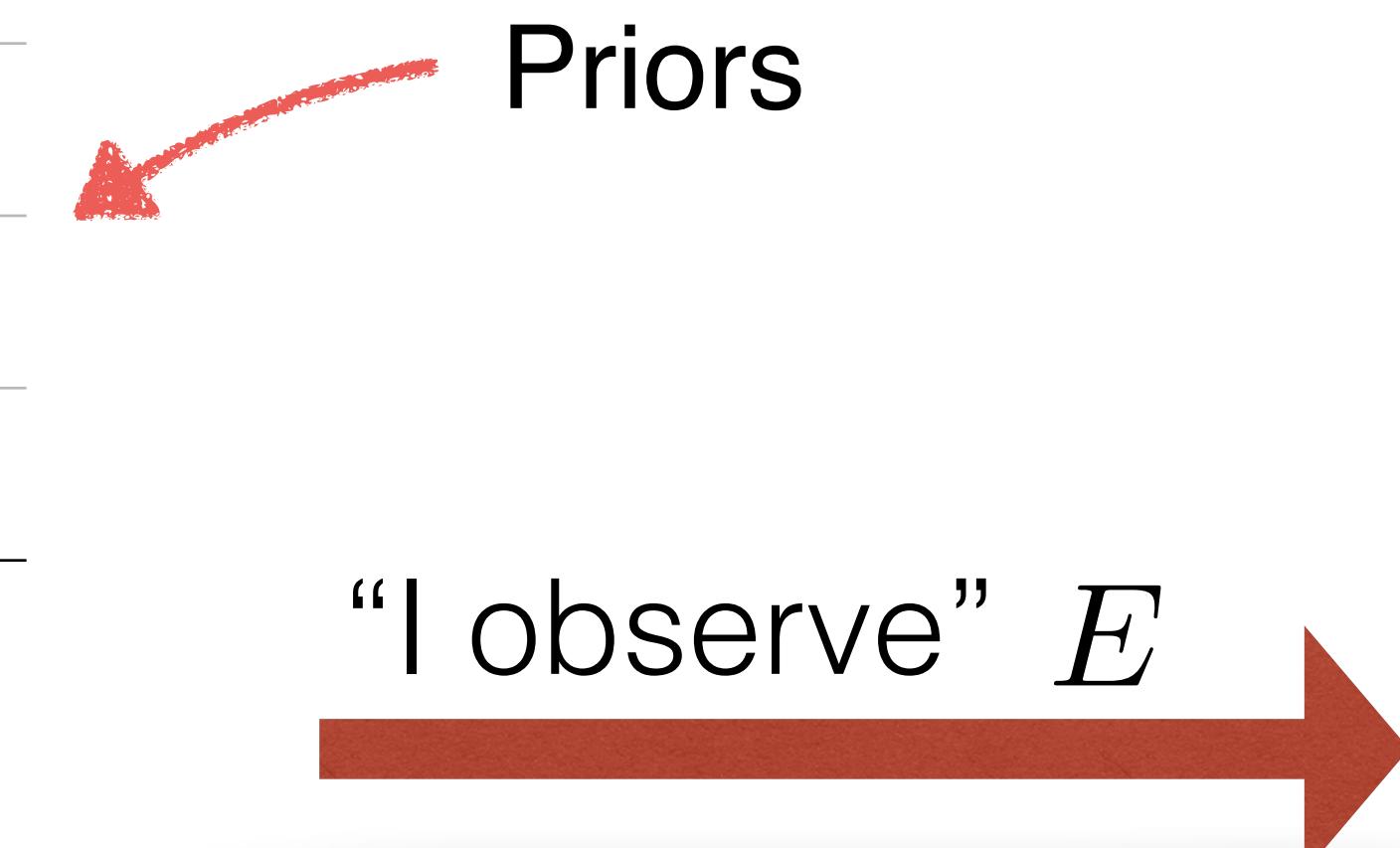
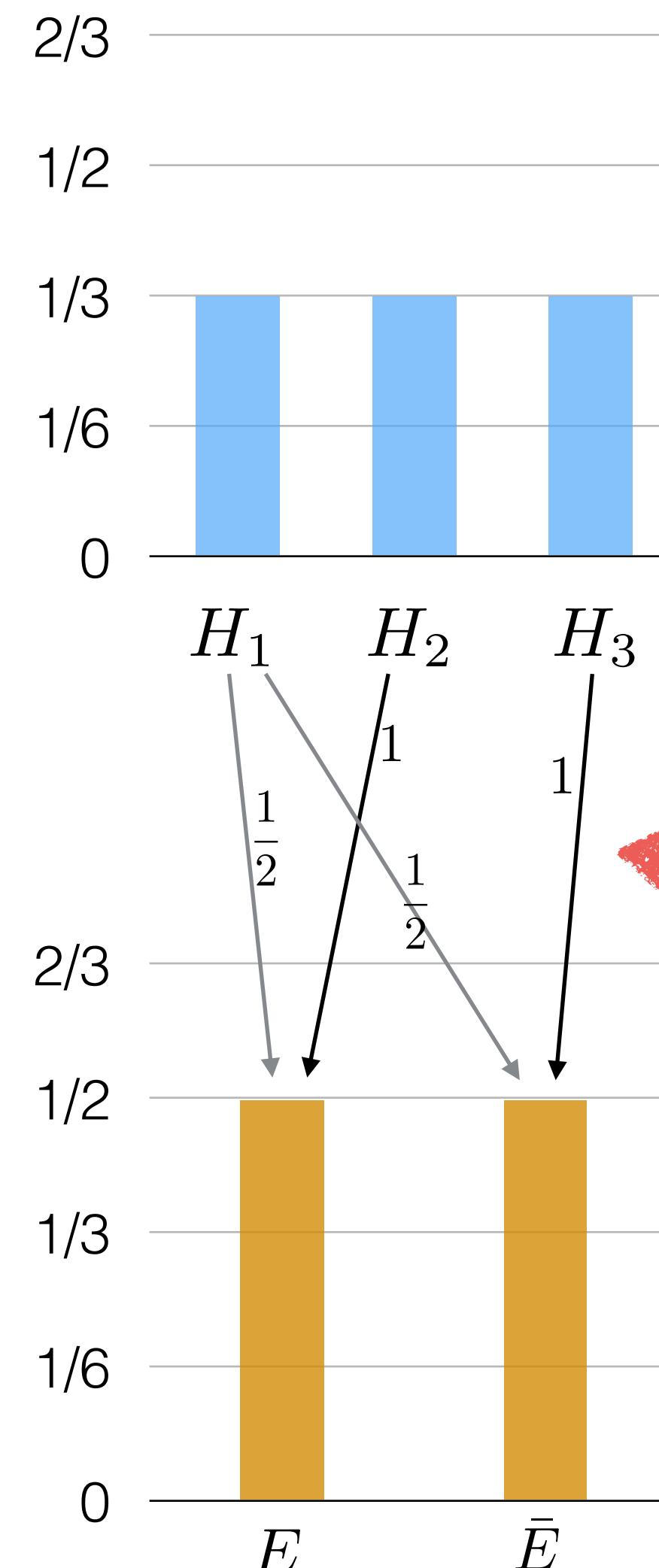
$$f(H_2|E, I) = \frac{f(E|H_2, I)f(H_2|I)}{f(E|I)} = \frac{1 \cdot 1/3}{1/2} = \frac{2}{3}$$

$$f(H_3|E, I) = \frac{f(E|H_3, I)f(H_3|I)}{f(E|I)} = \frac{0 \cdot 1/3}{1/2} = 0$$

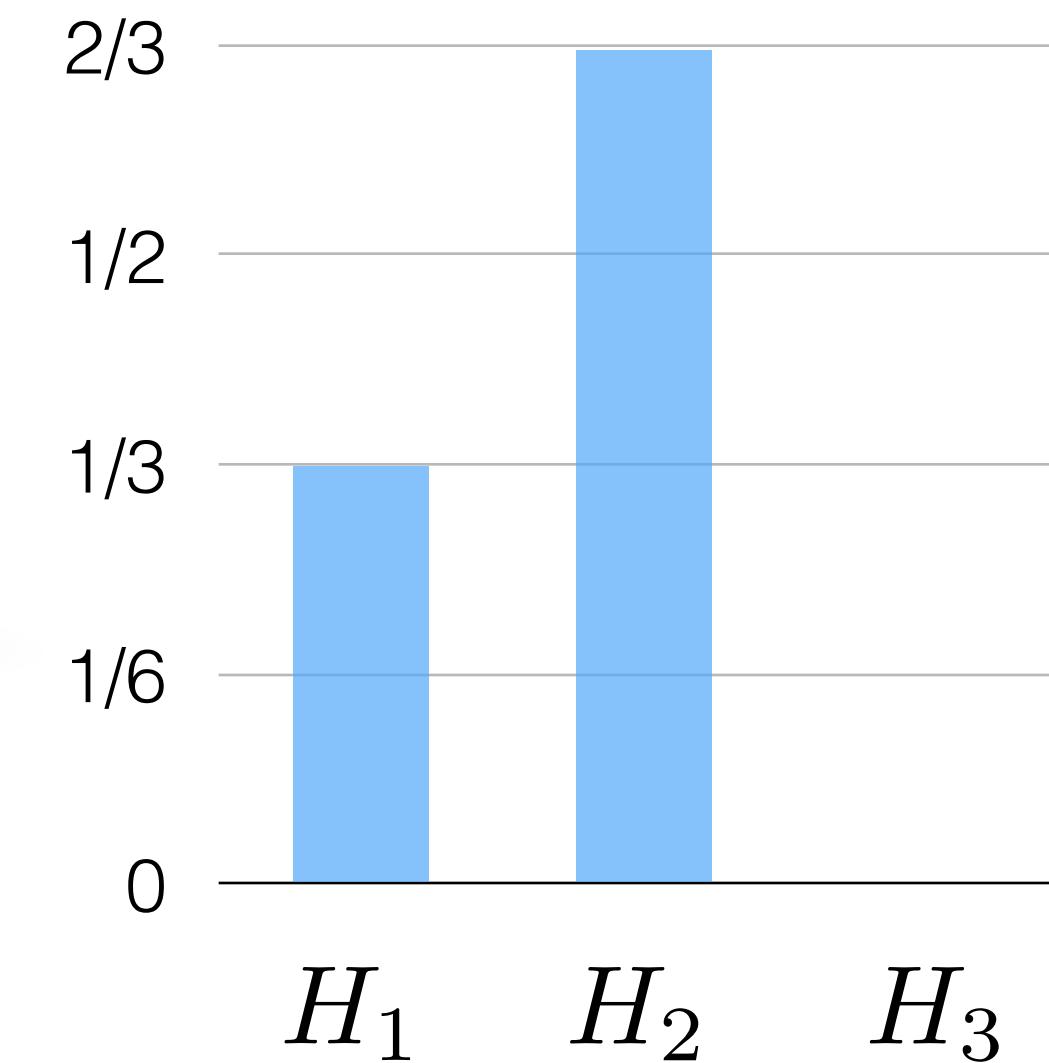
If we want to win the car,
we should change the box!



The Monty Hall problem



Likelihoods



Posteriors

What if the TV-Show host did not know where the car is?

Let's go back to the “sheep” example

The Model



1% of the sheep are black = M

The data



*Out of 1 thousand
sheep 20 are black*



= D

The opinion



$p(M|D)$

Let's go back to the “sheep” example

The Model



1% of the sheep are black = M

The data



*Out of 1 thousand
sheep 20 are black*



= D

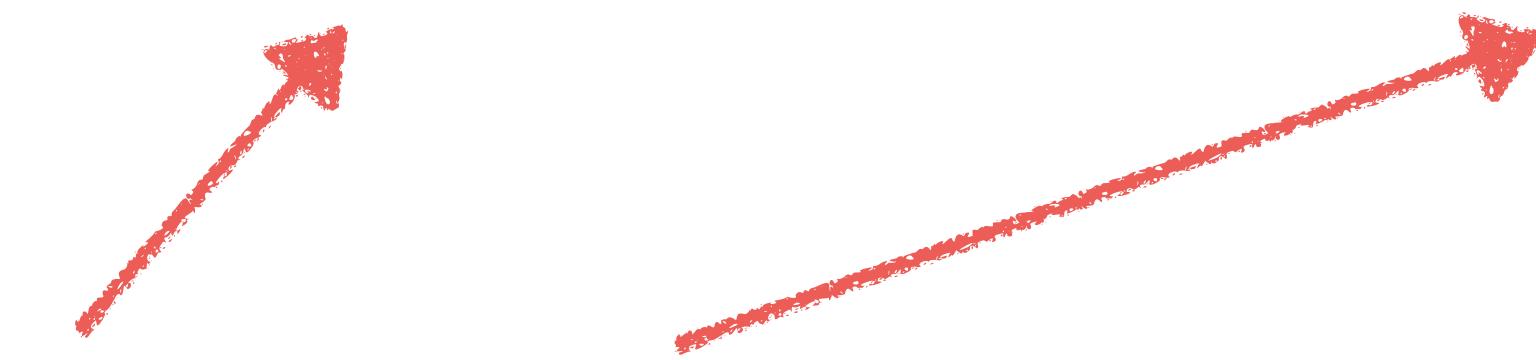
The opinion



$$p(M|D) = \frac{p(D|M)p(M)}{p(D|M)p(M) + p(D|\bar{M})p(\bar{M})}$$

Let's go back to the “sheep” example

$$p(M|D) = \frac{p(D|M)p(M)}{p(D|M)p(M) + p(D|\bar{M})p(\bar{M})}$$

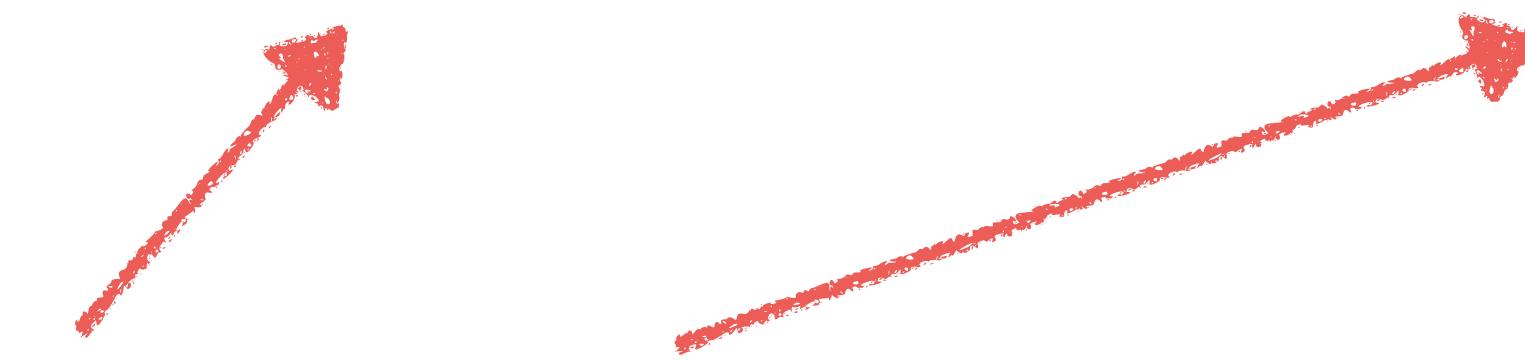


Our prior knowledge:

- How much do you believe in your model before the observation?
- Are there other models/hypotheses that might explain the observation? How likely are they?

Let's go back to the “sheep” example

$$p(M|D) = \frac{p(D|M)p(M)}{p(D|M)p(M) + p(D|\bar{M})p(\bar{M})}$$



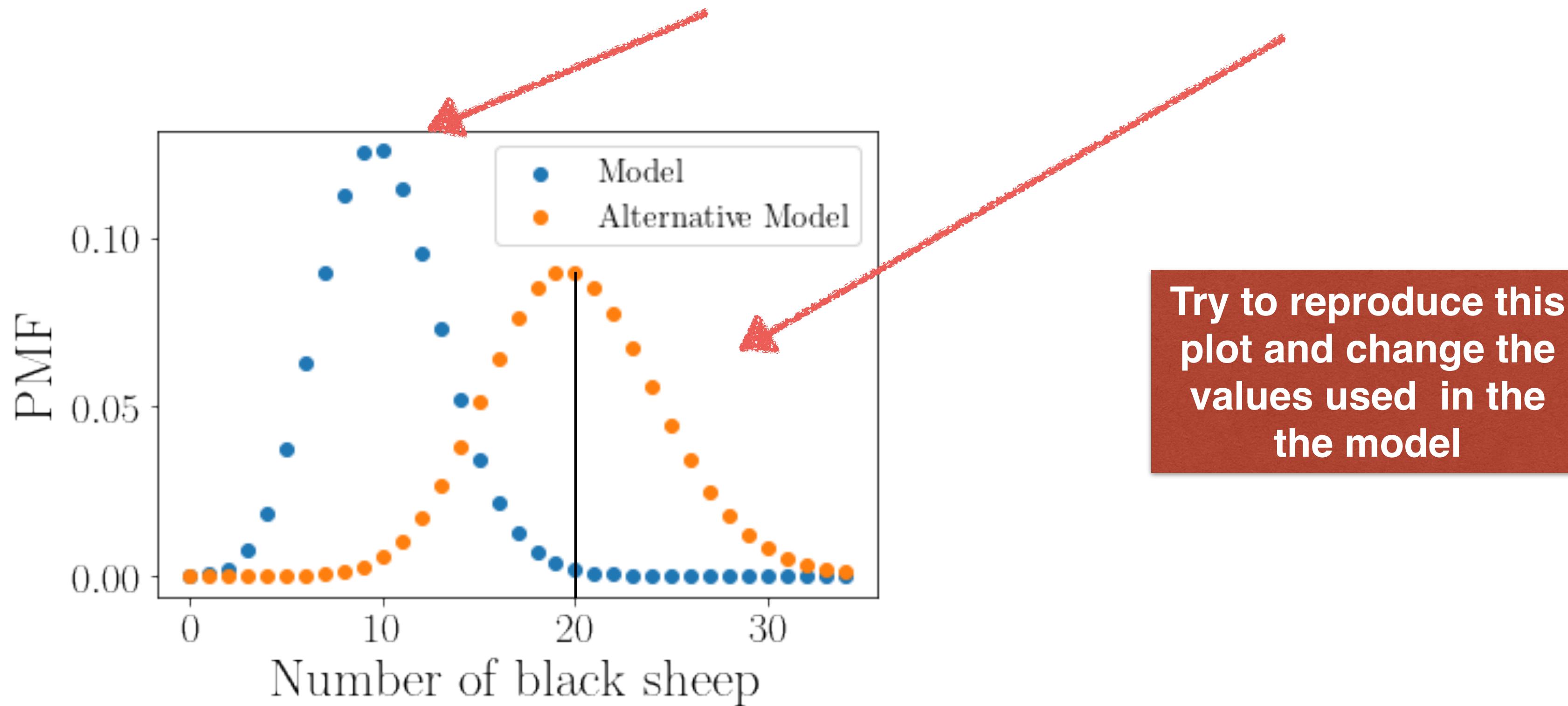
Our prior knowledge:

- We will assume for simplicity that there is only one alternative model “2% of the sheep are black”
- Both models are equally probable

$$p(M) = 1 - p(\bar{M}) = 0.5$$

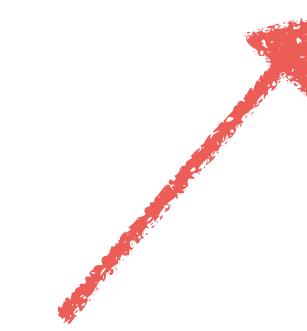
Let's go back to the “sheep” example

$$p(M|D) = \frac{p(D|M)p(M)}{p(D|M)p(M) + p(D|\bar{M})p(\bar{M})}$$



Let's go back to the “sheep” example

$$p(M|D) = \frac{p(D|M)p(M)}{p(D|M)p(M) + p(D|\bar{M})p(\bar{M})}$$



1% of the sheep are black



2% of the sheep are black

$$p(D|M) = \mathcal{B}(20 \mid p = 0.01, N = 10^3)$$

$$\simeq 0.0018$$

$$p(D|\bar{M}) = \mathcal{B}(20 \mid p = 0.02, N = 10^3)$$

$$\simeq 0.090$$

Let's go back to the “sheep” example

$$p(M|D) = \frac{p(D|M)p(M)}{p(D|M)p(M) + p(D|\bar{M})p(\bar{M})} \simeq 2\%$$

$$p(\bar{M}|D) = 1 - p(M|D) \simeq 98\%$$

The alternative model is much more likely of being true and the Bayesian approach let us quantify this “likeliness”

Let's go back to the “sheep” example

The Model



1% of the sheep are black = M

The data



*Out of 1 thousand
sheep 20 are black*



= D

The opinion



$$p(M|D) \approx 2\%$$

... but what if we do not know the priors of the models?

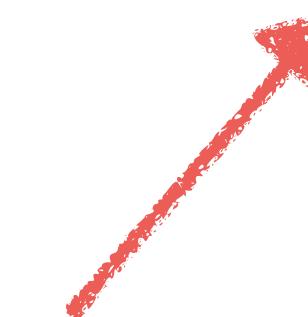
$$\frac{p(M|D)}{p(\bar{M}|D)} = \frac{p(D|M)}{p(D|\bar{M})} \times \frac{p(M)}{p(\bar{M})}$$

... but what if we do not know the priors of the models?

$$\frac{p(M|D)}{p(\bar{M}|D)} = \frac{p(D|M)}{p(D|\bar{M})} \times \frac{p(M)}{p(\bar{M})}$$

Bayes Factor

What is the BF in our example?



Bayes factor BF_{12}		Interpretation
	>	100
30	-	100
10	-	30
3	-	10
1	-	3
		No evidence
1/3	-	1
1/10	-	1/3
1/30	-	1/10
1/100	-	1/30
	<	1/100
		Extreme evidence for M_2

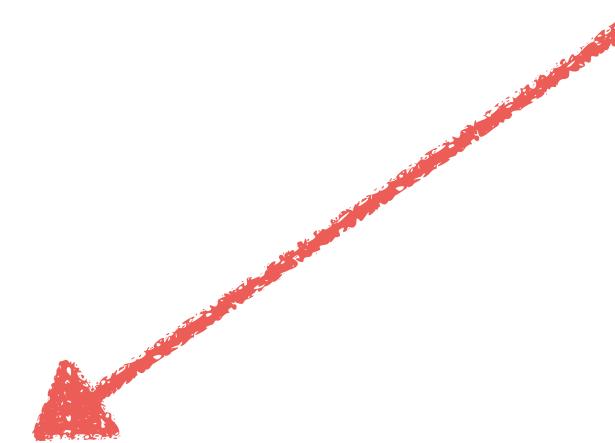
The Frequentist approach

- The **Frequentist approach** tries to answer the question:

*If I repeat the experiment an infinite time, assuming the model is **true**, with which **frequency** I would observe a value more **extreme** than the one actually observed?*

- The **Frequentist approach** tries to answer the question:

*If I repeat the experiment an infinite time, assuming the model is **true**, with which **frequency** I would observe a value more **extreme** than the one actually observed?*



The **data “D”** itself or a function of them known as the **statistic**

$$\mathcal{S} = \mathcal{S}(D)$$

Let's go back to the “sheep” example

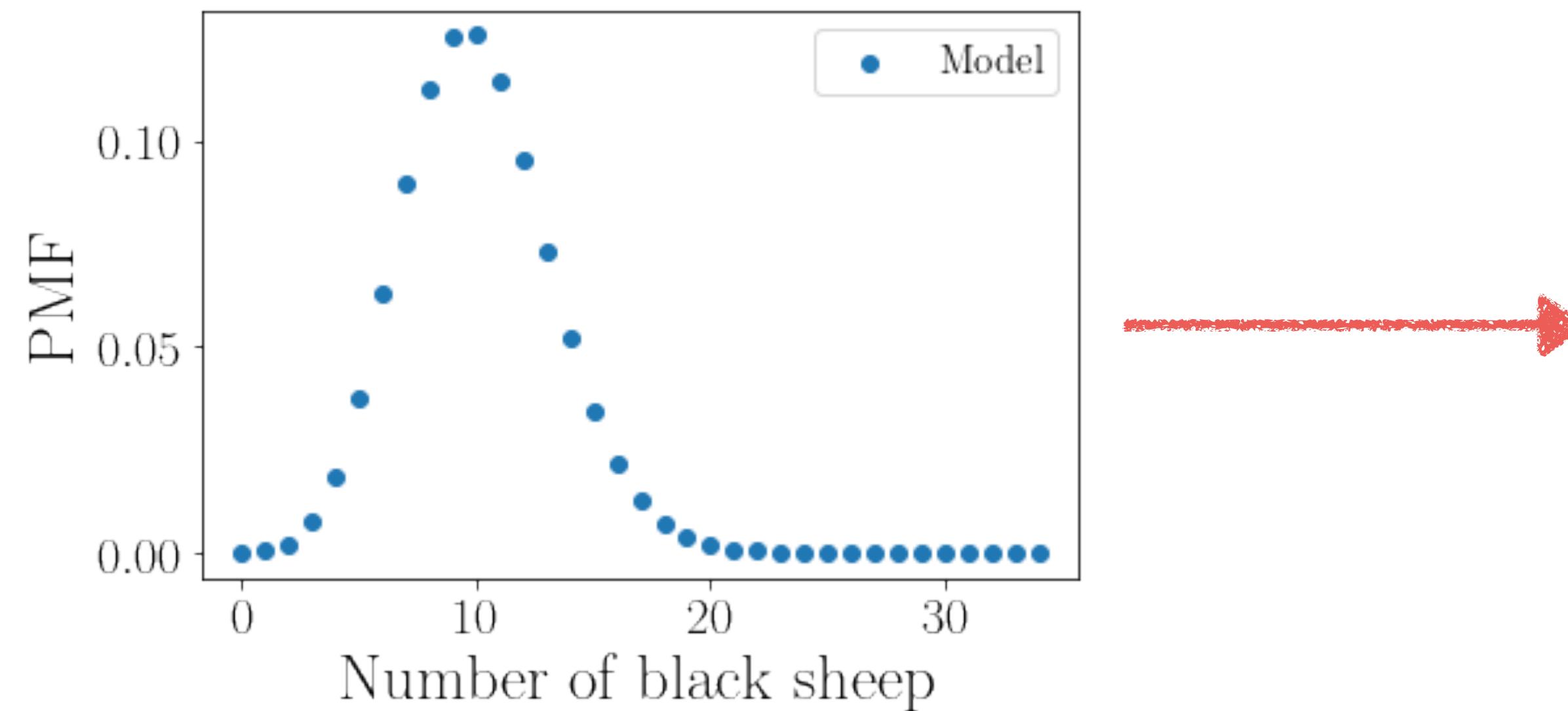
We can use the number of sheep observed as statistics and ask ourselves:

If I repeat the observation an infinity of time, how frequently would I have observed 20 or more sheep?

Let's go back to the “sheep” example

We can use the number of sheep observed as statistics and ask ourselves:

If I repeat the observation an infinity of time, how frequently would I have observed 20 or more sheep?



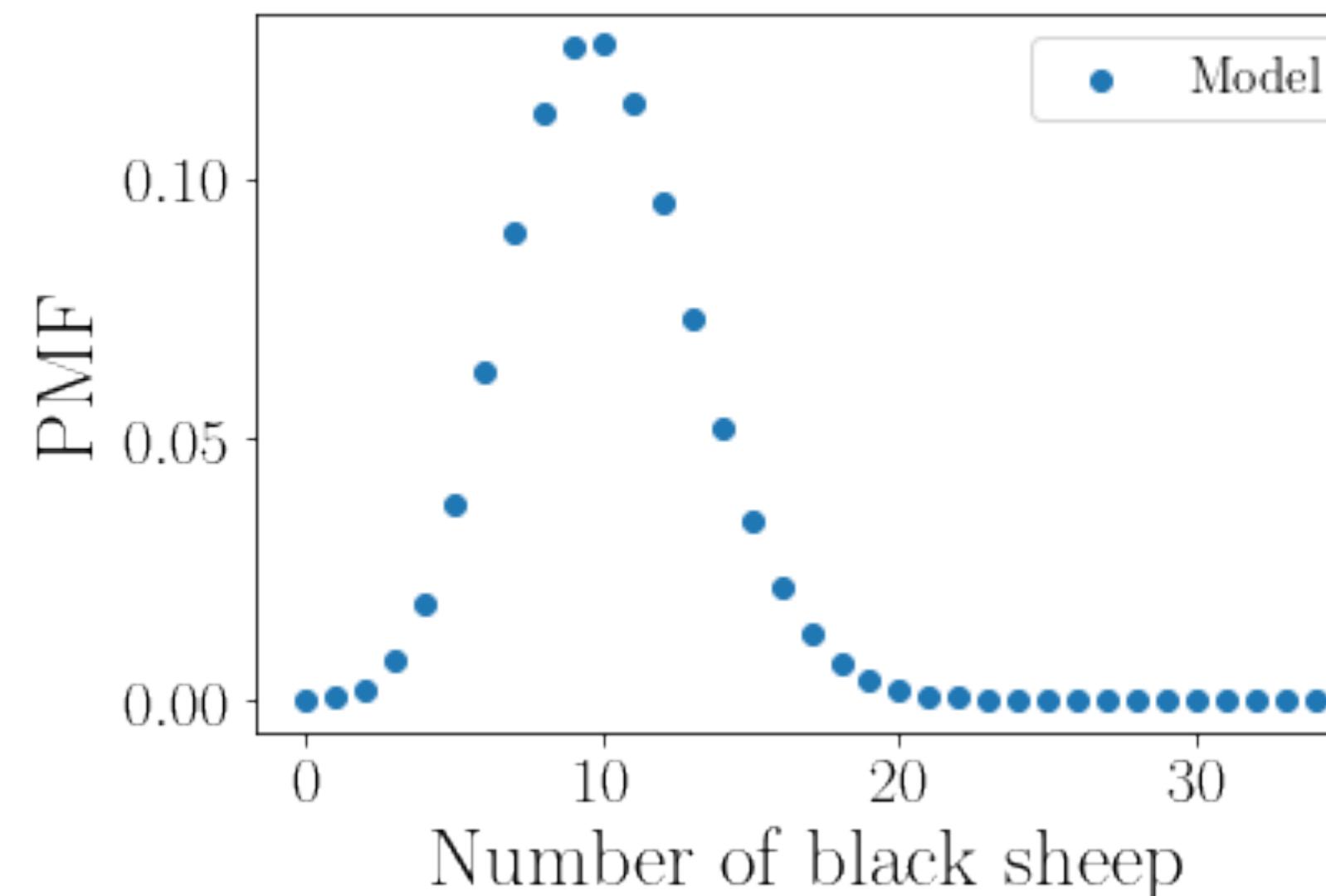
The answer is only 0.1% of the time!

Therefore the frequentist conclusion is that our model is excluded with a 99.9% confidence level.

Let's go back to the “sheep” example

We can use the number of sheep observed as statistics and ask ourselves:

If I repeat the observation an infinity of time, how frequently would I have observed 20 or more sheep?



Try to get this value

P-VALUE

The answer is only 0.1% of the time!



Therefore the frequentist conclusion is that our model is excluded with a 99.9% confidence level.

The **P-VALUE** is the frequency in which we would have observed “something” more extreme assuming the null hypothesis to be true

$$\text{p-value} = p(x \text{ more extreme than } x_{obs} | H_0)$$

The P-VALUE is the frequency in which we would have observed “something” more extreme assuming the null hypothesis to be true

$$\text{p-value} = p(x \text{ more extreme than } x_{obs} | H_0)$$

... but then, what are all these “sigmas”?

We therefore consider emission from the Sgr dSph as an alternative origin for the cocoon. In order to test this possibility, we fit the γ -ray emission observed by *Fermi*-LAT over a region of interest (ROI) containing the cocoon via template analysis. In our *baseline* model these templates include only known point sources and sources of Galactic diffuse γ -ray emission. We contrast the baseline with a *baseline + Sgr dSph* model that invokes these same templates plus an additional template constructed to be spatially coincident with the bright stars of the Sgr dSph (Extended Data (E.D.) Figure 1 and S.I. Figure 1); full details of the fitting procedure are provided in Methods and S.I. sec. 3. Using the best motivated choice of templates, we find that the baseline + Sgr dSph model is preferred at 8.1σ significance over the baseline model. We also repeat the analysis for a wide range of alternative templates for both Galactic diffuse emission and for the Sgr dSph (Table 1) and obtain $> 5\sigma$ detections for all combinations but one. Moreover, even this is an extremely conservative estimate, because our baseline model uses a structured template for the FBs that absorbs some of the signal that is spatially coincident with the Sgr dSph into a structure of unknown origin. If we follow the method recommended by the Fermi collaboration [2] and use a flat FB template in our analysis, the significance of our detection of the Sgr dSph is always $> 14\sigma$. Despite this, for the remainder of our analysis we follow the most conservative choice by using the structured template in our baseline model. In Methods, we also show that our analysis passes a series of validation tests: the residuals between our best-fitting model and the data are consistent with photon counting statistics (E.D. Figure 2 and Figure 3), our pipeline reliably recovers synthetic signals superimposed on a realistic background (E.D. Figure 4), fits using a template tracing the stars of the Sgr dSph yield significantly better results than fits using purely geometric templates (S.I. Table 1), and if we artificially rotate the Sgr dSph template on the sky, the best-fitting position angle is very close to the actual one (E.D. Figure 5). By contrast, if we displace the Sgr dSph template, we find moderate (4.5σ significance) evidence that the best-fitting position is $\sim 4^\circ$ from the true position, in a direction very closely aligned with the dwarf galaxy’s direction of travel (E.D. Figure 5); this plausibly represents a small, but real and expected (as explained below) physical offset between the stars and the γ -ray emission.

PKS 1413+135: Bright GeV γ -ray Flares with Hard-spectrum and Hints for First Detection of TeV γ -rays from a Compact Symmetric Object

YING-YING GAN,¹ JIN ZHANG^{†,1} SU YAO,² HAI-MING ZHANG,³ YUN-FENG LIANG,⁴ AND EN-WEI LIANG⁴

¹School of Physics, Beijing Institute of Technology, Beijing 100081, People’s Republic of China; j.zhang@bit.edu.cn

²Max-Planck-Institute für Radioastronomie, Auf dem Hügel 69, 53121 Bonn, Germany

³School of Astronomy and Space Science, Nanjing University, Nanjing 210023, People’s Republic of China

⁴Guangxi Key Laboratory for Relativistic Astrophysics, School of Physical Science and Technology, Guangxi University, Nanning 530004, People’s Republic of China

ABSTRACT

PKS 1413+135, a typical compact symmetric object (CSO) with a two-side pc-scale structure in its miniature radio morphology, is spatially associated with the *Fermi*-LAT source 4FGL J1416.1+1320 and recently announced to be detected in the TeV γ -ray band with the MAGIC telescopes. We present the analysis of its X-ray and GeV γ -ray observations obtained with *Swift*-XRT, *XMM-Newton*, *Chandra*, and *Fermi*-LAT for revealing its high energy radiation physics. No significant variation trend is observed in the X-ray band. Its GeV γ -ray light curve derived from the *Fermi*-LAT 13.5-year observations shows that it is in a low γ -ray flux stage before MJD 58500 and experiences violent outbursts after MJD 58500. The confidence level of the flux variability is much higher than 5σ , and the flux at 10 GeV varies ~ 3 orders of magnitude. The flux variation is accompanied by the clearly

The P-VALUE is the frequency in which we would have observed “something” more extreme assuming the null hypothesis to be true

$$\text{p-value} = p(x \text{ more extreme than } x_{obs} | H_0)$$

... but then, what are all these “sigmas”?

We therefore consider emission from the Sgr dSph as an alternative origin for the cocoon. In order to test this possibility, we fit the γ -ray emission observed by *Fermi*-LAT over a region of interest (ROI) containing the cocoon via template analysis. In our *baseline* model these templates include only known point sources and sources of Galactic diffuse γ -ray emission. We contrast the baseline with a *baseline + Sgr dSph* model that invokes these same templates plus an additional template constructed to be spatially coincident with the bright stars of the Sgr dSph (Extended Data (E.D.) Figure 1 and S.I. Figure 1); full details of the fitting procedure are provided in Methods and S.I. sec. 3. Using the best motivated choice of templates, we find that the baseline + Sgr dSph model is preferred at 8.1σ significance over the baseline model. We also repeat the analysis for a wide range of alternative templates for both Galactic diffuse emission and for the Sgr dSph (Table 1) and obtain $> 5\sigma$ detections for all combinations but one. Moreover, even this is an extremely conservative estimate, because our baseline model uses a structured template for the FBs that absorbs some of the signal that is spatially coincident with the Sgr dSph into a structure of unknown origin. If we follow the method recommended by the *Fermi* collaboration [2] and use a flat FB template in our analysis, the significance of our detection of the Sgr dSph is always $> 14\sigma$. Despite this, for the remainder of our analysis we follow the most conservative choice by using the structured template in our baseline model. In Methods, we also show that our analysis passes a series of validation tests: the residuals between our best-fitting model and the data are consistent with photon counting statistics (E.D. Figure 2 and Figure 3), our pipeline reliably recovers synthetic signals superimposed on a realistic background (E.D. Figure 4), fits using a template tracing the stars of the Sgr dSph yield significantly better results than fits using purely geometric templates (S.I. Table 1), and if we artificially rotate the Sgr dSph template on the sky, the best-fitting position angle is very close to the actual one (E.D. Figure 5). By contrast, if we displace the Sgr dSph template, we find moderate (4.5σ significance) evidence that the best-fitting position is $\sim 4^\circ$ from the true position, in a direction very closely aligned with the dwarf galaxy’s direction of travel (E.D. Figure 5); this plausibly represents a small, but real and expected (as explained below) physical offset between the stars and the γ -ray emission.

astro-ph.HE] 19 Jun 2022

PKS 1413+135: Bright GRB

YING-YING GAN,¹ JIN

¹School of Physics, Beij

²Max-Planck

³School of Astronom

⁴Guangxi Key Laboratory for Relat

PKS 1413+135, a type I GRB, was first detected by the Neil Gehrels Swift Observatory at 13:09:59 UT and recently announced by the *Fermi*-LAT team. We present the analysis of the X-ray data from the *Chandra*, and *Fermi*-LAT. A clear trend is observed in the X-ray light curve. The year observations show two outbursts after MJD 59651. The flux at 10 GeV varies between 10^{-1} and 10^{-2} photons cm $^{-2}$ s $^{-1}$.

GRB 211211A triggered the Burst Alert Telescope (Barthelmy et al. 2005) onboard The Neil Gehrels Swift Observatory at 13:09:59 UT (D’Ai et al. 2021), the Gamma-ray Burst Monitor (Meegan et al. 2009) onboard The Fermi Gamma-Ray Space Telescope at 13:09:59.651 UT (Mangan et al. 2021) and High energy X-ray Telescope onboard Insight-HXMT (Xiao et al. 2022) at 13:09:59 UT on 11 December 2021. The burst is characterized by a spiky main emission phase lasting \sim 13 seconds, and a longer, weaker extended emission phase lasting \sim 55 seconds (Yang et al. 2022). The prompt emission is suggested to be produced by

the fast-cooling synchrotron emission (Gompertz et al. 2022). The discovery of a kilonova associated with this GRB indicates clearly that the progenitor is a compact object merger (Rastinejad et al. 2022). The event fluence ($10\text{-}1000$ keV) of the prompt emission is $(5.4 \pm 0.01) \times 10^{-4}$ erg cm $^{-2}$, making this GRB an exceptionally bright event. The host galaxy redshift of GRB 211211A is $z = 0.0763 \pm 0.0002$ (corresponding to a distance of \approx 350 Mpc (Rastinejad et al. 2022)). At 350 Mpc, GRB 211211A is one of the closest GRBs, only a bit further than GRB 170817A, which is associated with the gravitational wave (GW)-detected binary neutron star (BNS) merger GW170817. For GRB 170817A, no GeV afterglow was detected by the LAT on timescales of minutes, hours, or days after the LIGO/Virgo detection (Ajello et al. 2018).

As the angle from the *Fermi*-LAT boresight at the GBM trigger time of GRB 211211A is 106.5 degrees (Mangan et al. 2021), LAT cannot place constraints on the existence of high-energy ($E > 100$ MeV) emission associated with the prompt GRB emission. We focus instead on constraining high-energy emission on the longer timescale. We analyze the late-time *Fermi*-LAT data when the GRB enters the field-of-view (FOV) of *Fermi*-LAT. We detect a transient source with a significance of $TS_{\text{max}} \simeq 51$, corresponding to a detection significance over 6σ . The result of the data analysis is shown in §2

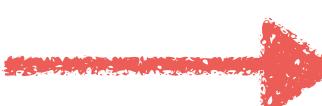
The **P-VALUE** is the frequency in which we would have observed “something” more extreme assuming the null hypothesis to be true

$$\text{p-value} = p(x \text{ more extreme than } x_{obs} | H_0)$$

... but then, what are all these “sigmas”?

It is common to express such probability in multiples S of the standard deviations of a normal distribution via the inverse error function

$$S = \sqrt{2} \operatorname{erf}^{-1} (1 - \text{p-value})$$



Here the (in-)famous number of “sigma”

Let's go back to the “sheep” example

The Model



1% of the sheep are black

The data



*Out of 1 thousand
sheep 20 are black*



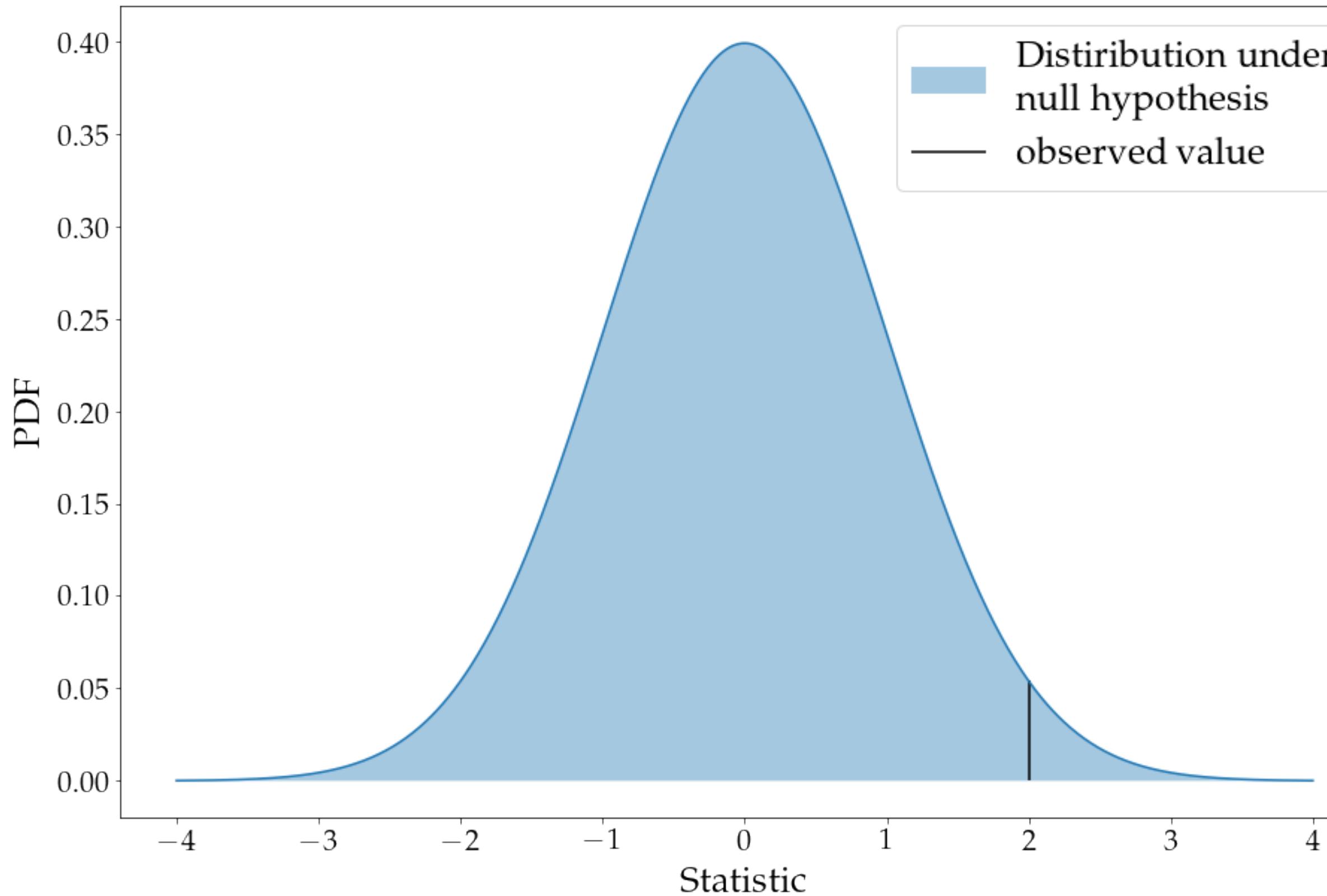
The opinion



*The model is excluded at 3.2
sigma*

Issues of the frequentist approach:

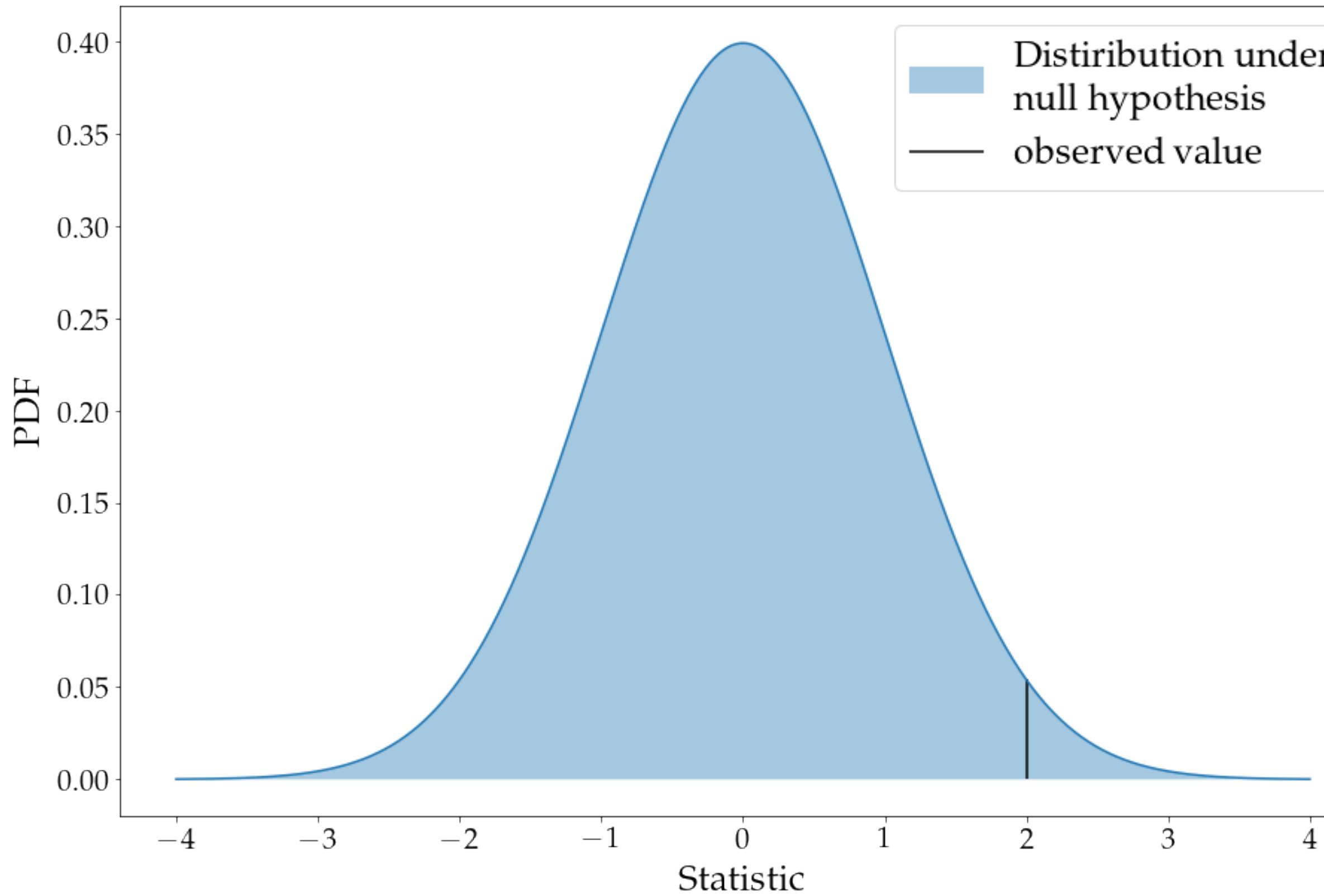
It does not take into account the **alternative hypothesis** that might explain the outcome of an event



Conclusion:
The null hypothesis is rejected with
a **2 sigma** significance

Issues of the frequentist approach:

It does not take into account the **alternative hypothesis** that might explain the outcome of an event



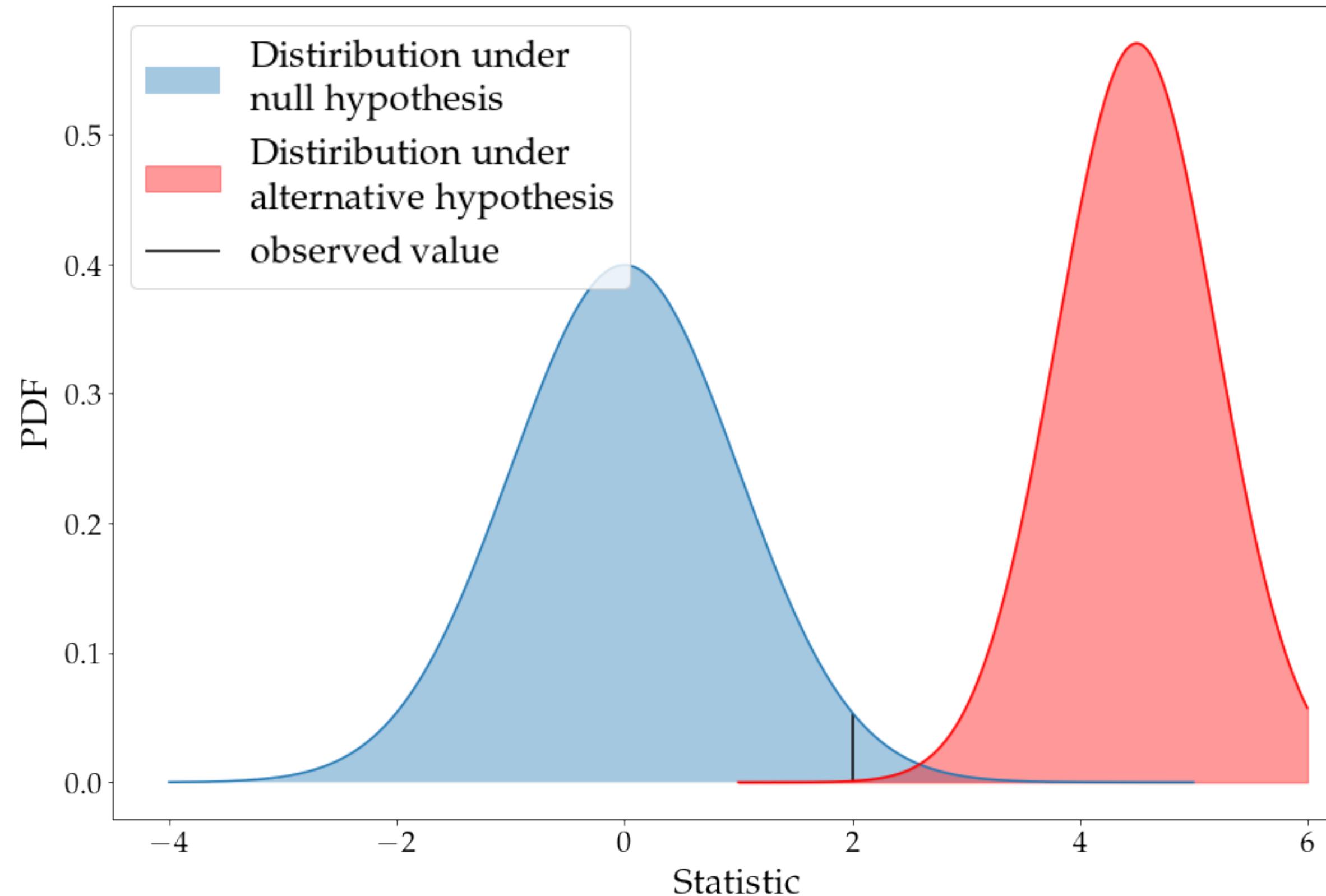
Conclusion:

The null hypothesis is rejected with a **2 sigma** significance

But what about the alternative hypothesis?

Issues of the frequentist approach:

It does not take into account the **alternative hypothesis** that might explain the outcome of an event



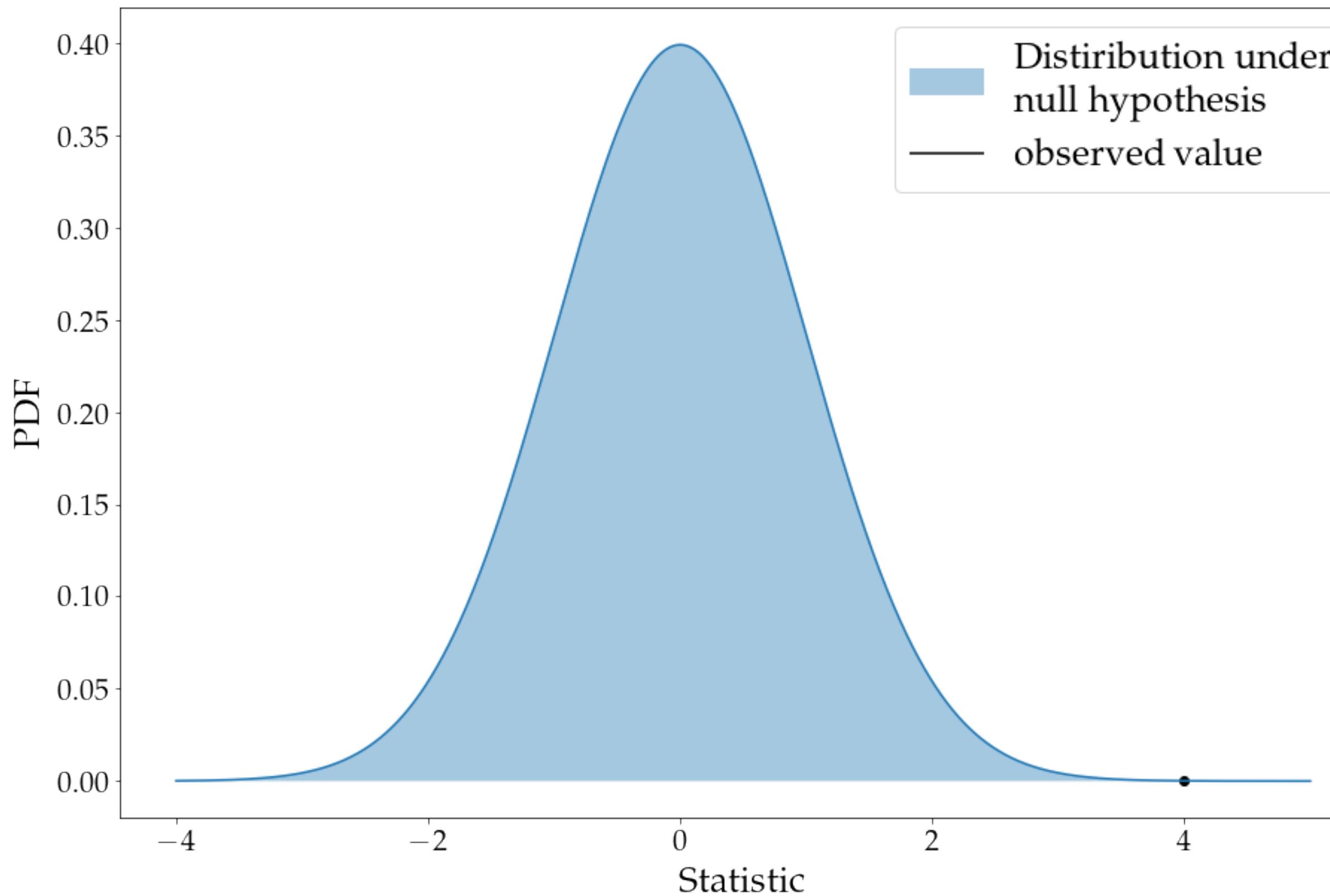
The observed value of 2 is actually more plausible being the outcome of the null hypothesis

By rejecting the null hypothesis we would have done the so-called ***type I error***

This is why a value of sigma bigger than 3 or even 5 is required for making a claim!

Issues of the frequentist approach:

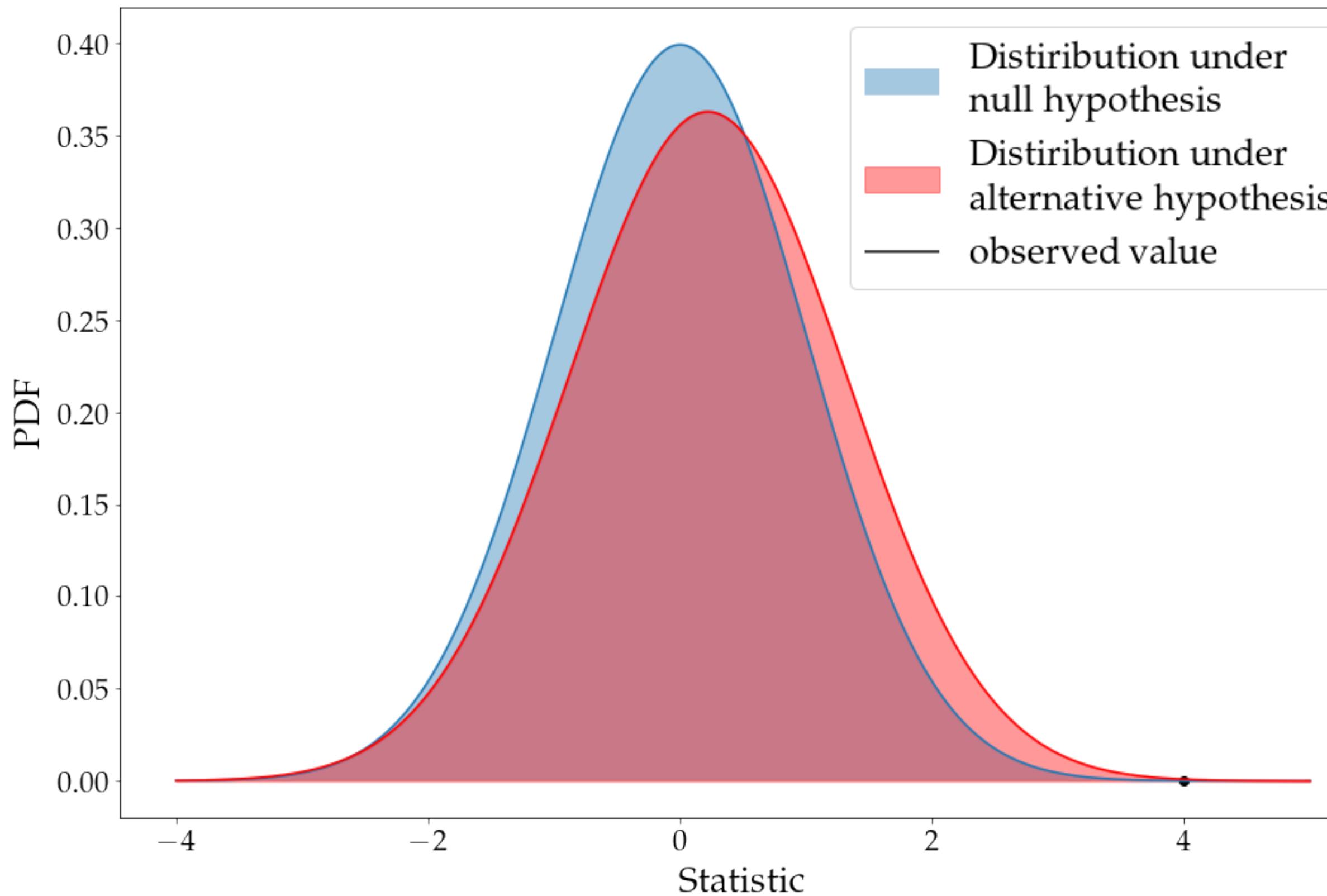
It does not take into account the **alternative hypothesis** that might explain the outcome of an event



So... with a significance of 4 we should be safe?

Issues of the frequentist approach:

It does not take into account the **alternative hypothesis** that might explain the outcome of an event



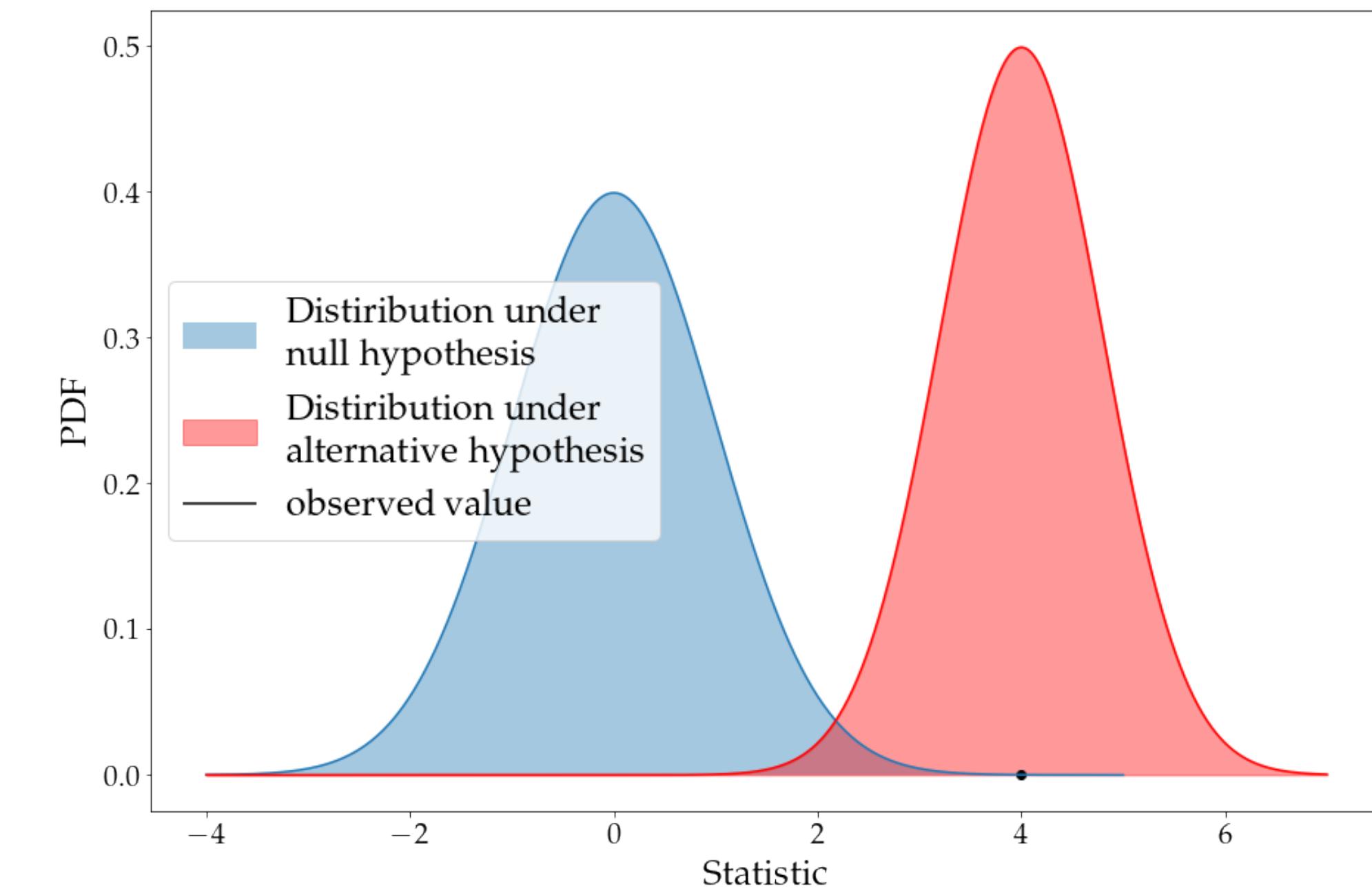
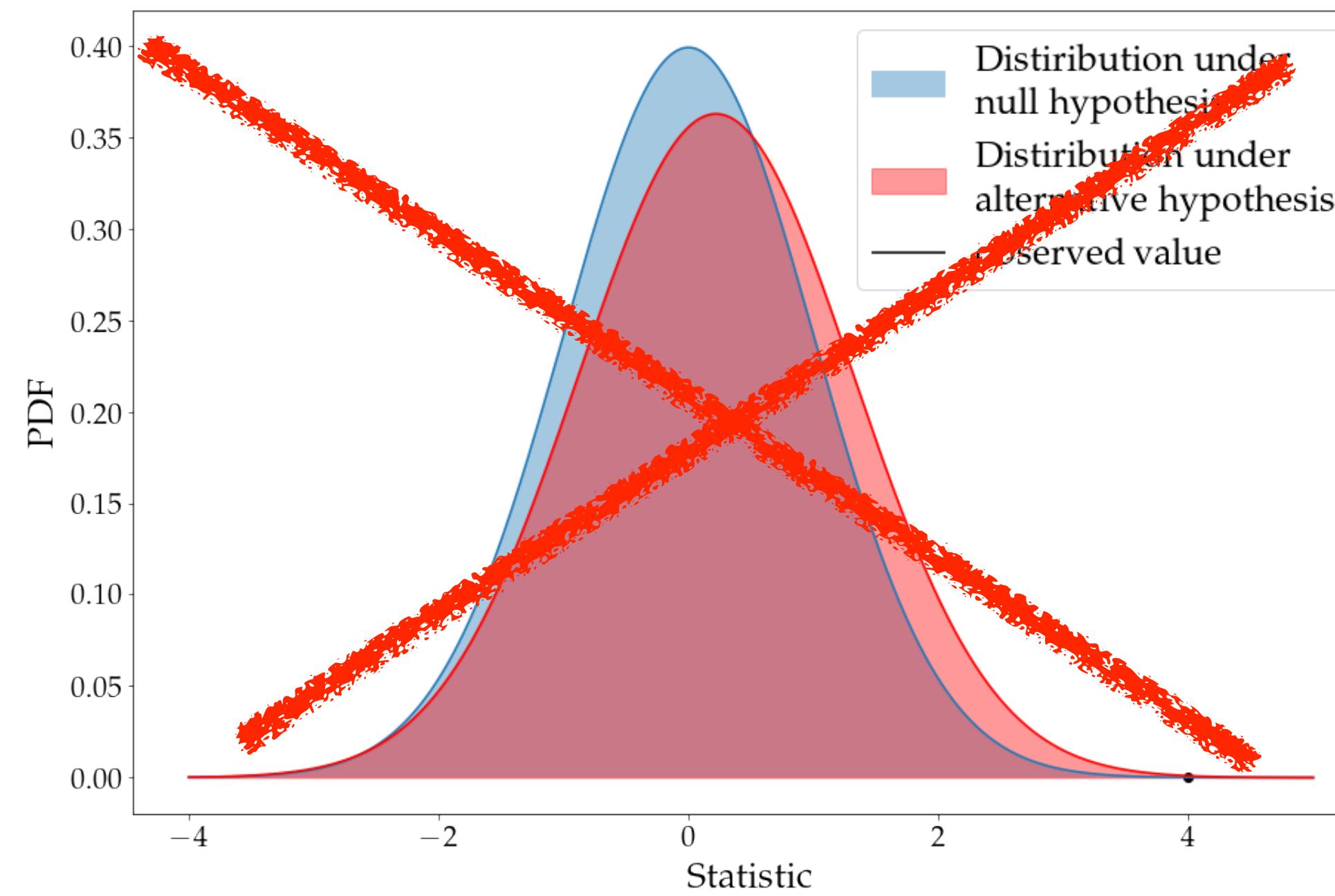
So... with a significance of 4 we should be safe?

The value of 4 is unlikely to be the outcome also of the alternative hypothesis, thus again we could be doing a ***type I error***

Issues of the frequentist approach:

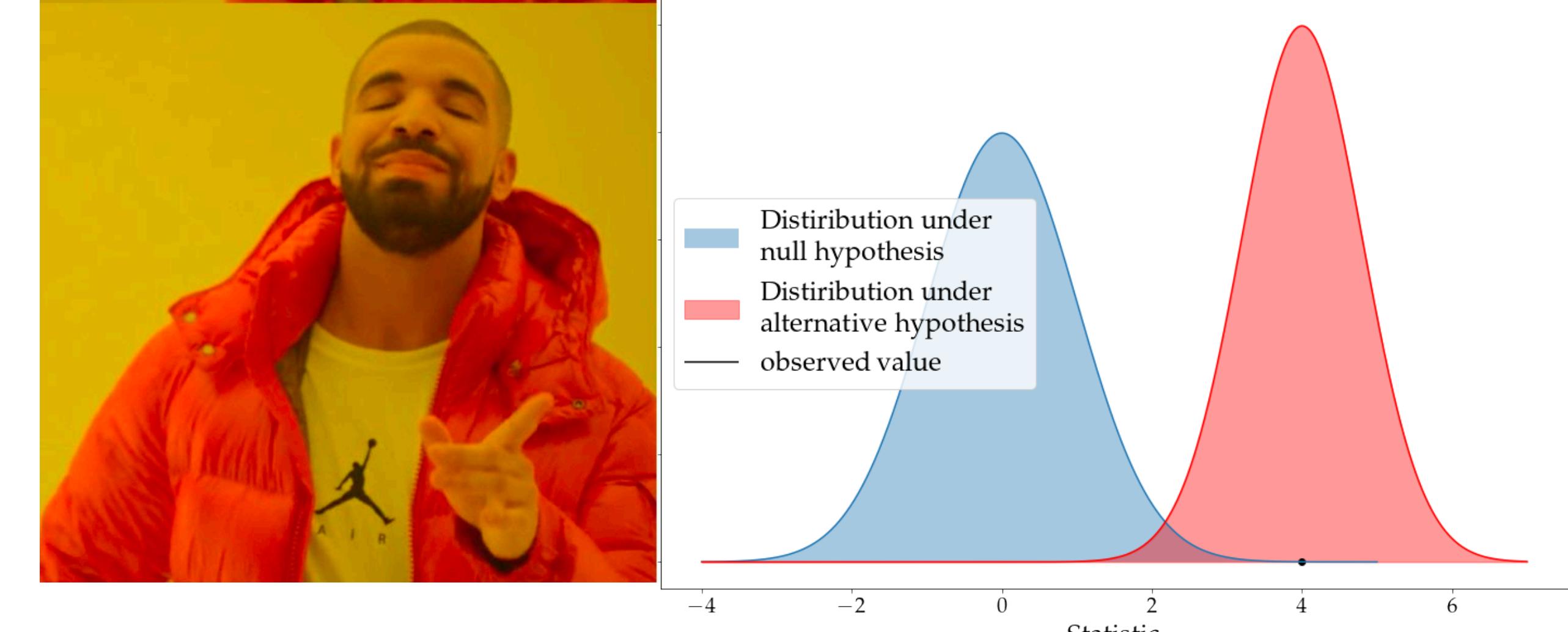
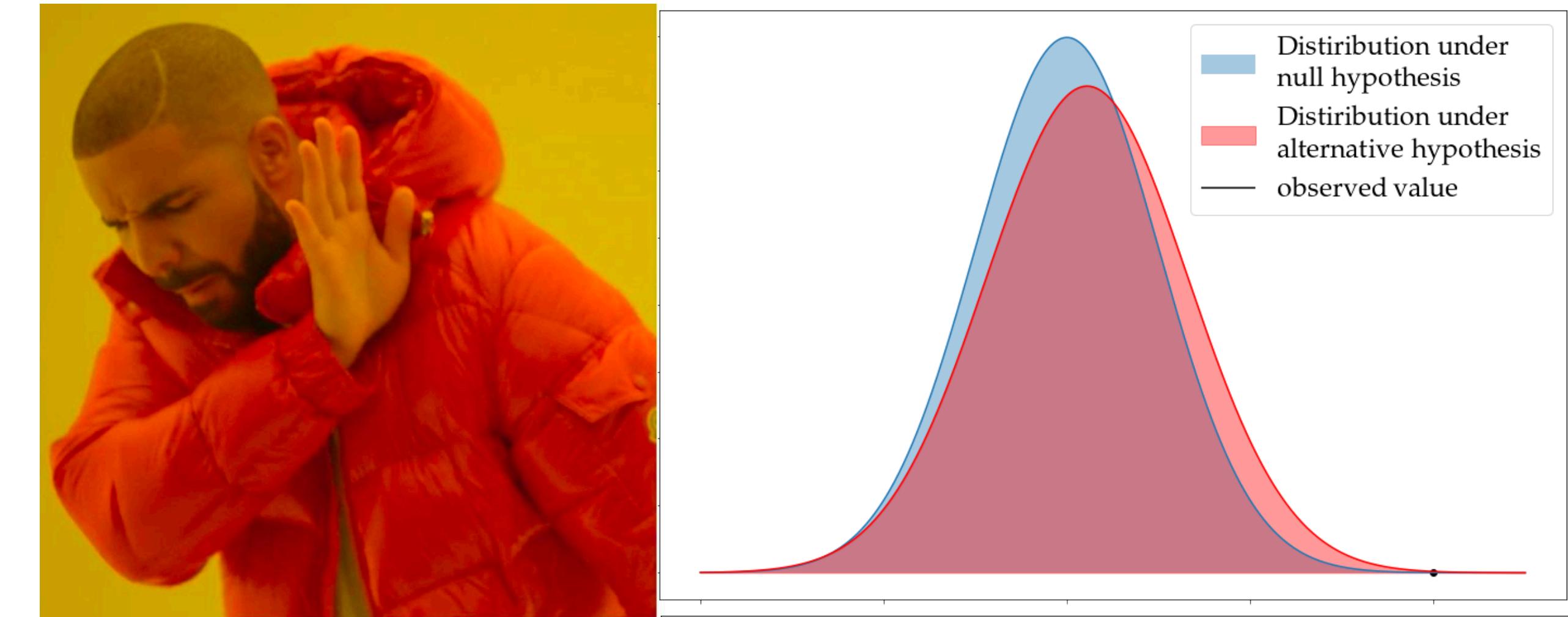
It does not take into account the **alternative hypothesis** that might explain the outcome of an event

The ideal statistic is the one that makes you **reject** a hypothesis that is false!



Issues of the frequentist approach:

The probability of rejecting a hypothesis that is false is called the “**power**” of the statistic



Your statistic must be
POWERFUL!

Thankfully the Neyman-Pearson Lemma tells us that the most “powerful” statistic is the **likelihood ratio**:

Parameter of the
null hypothesis

$$\frac{\mathcal{L}(\theta | D_{obs})}{\mathcal{L}(\hat{\theta} | D_{obs})}$$

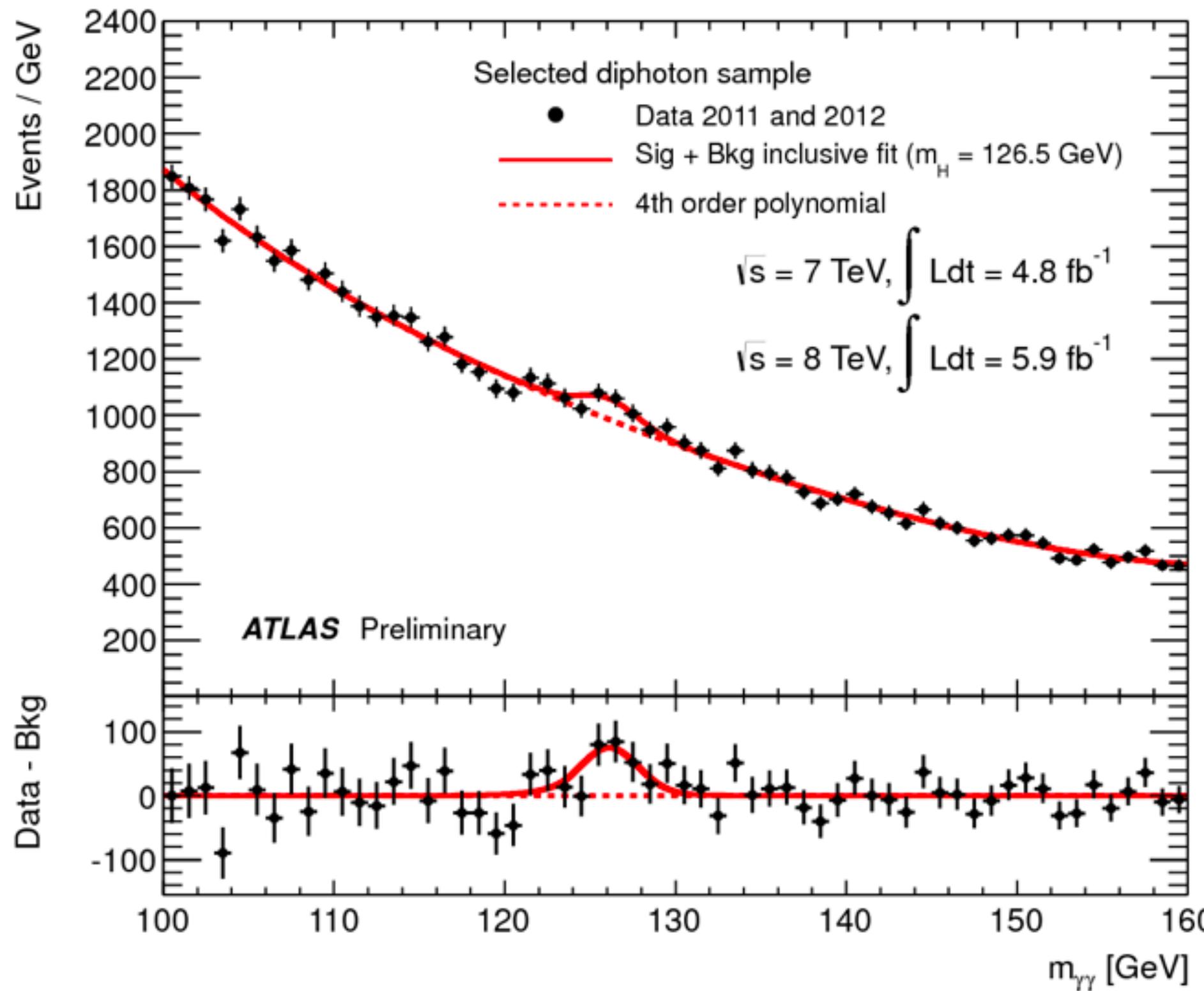
Best fit or value
that maximises
the likelihood

Observed data

Likelihood

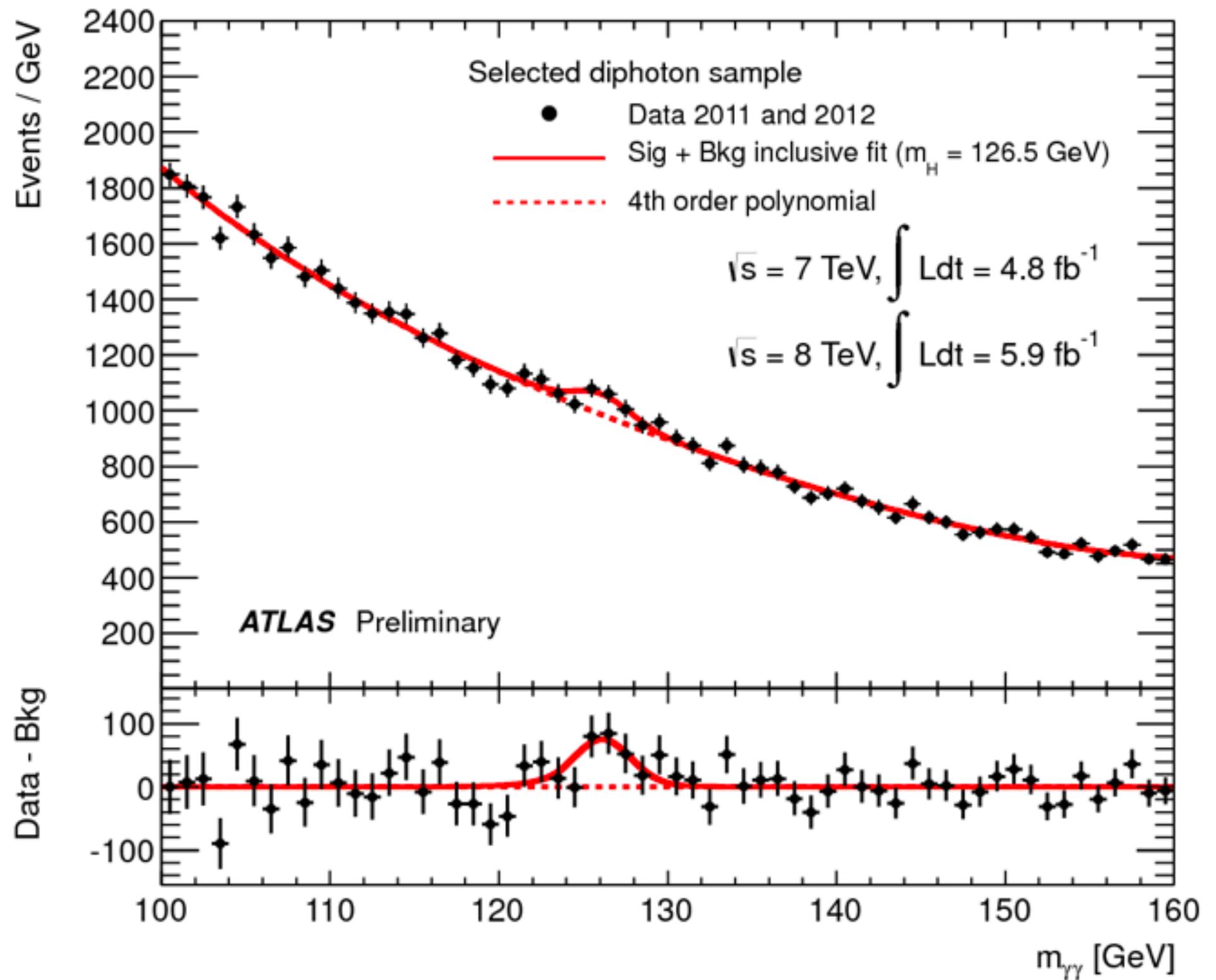
$$\mathcal{L}(\theta | D_{obs}) = p(D_{obs} | \theta)$$

Example:



This is the plot that led ATLAS to claim the **discovery of the HIGGS**.
Let's figure out how they were able to make such a claim with a **Toy Model** and with the **theory** we have learned so far

Example:

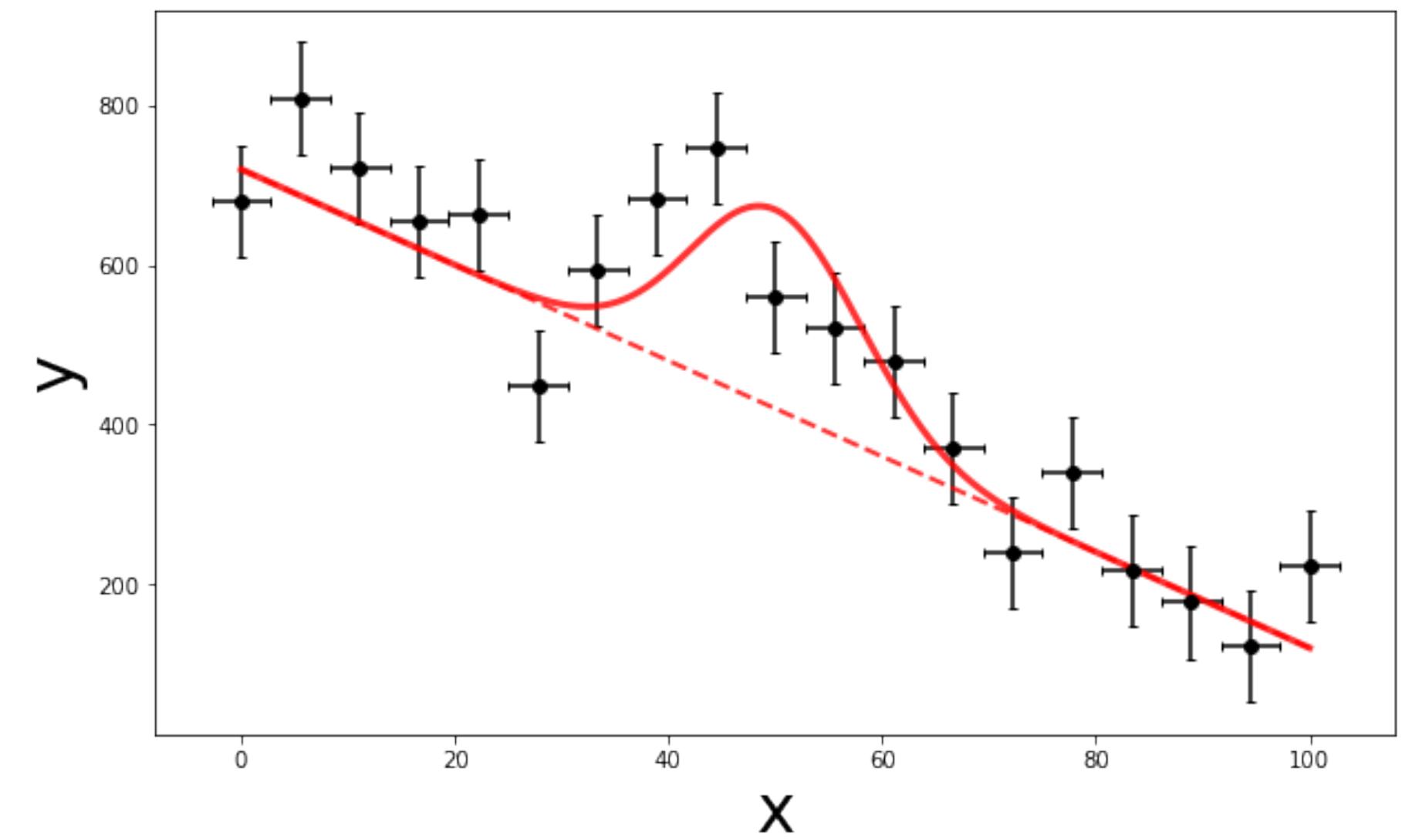


Toy Model



$$y' = mx + q + a \cdot \mathcal{G}(x; \mu = 50, \sigma = 8)$$

$$y \sim \mathcal{N}(\mu = y', \sigma = 70)$$



Null hypothesis H_0

$$a = 0$$

Alternative hyp. H_1

$$a = 5$$

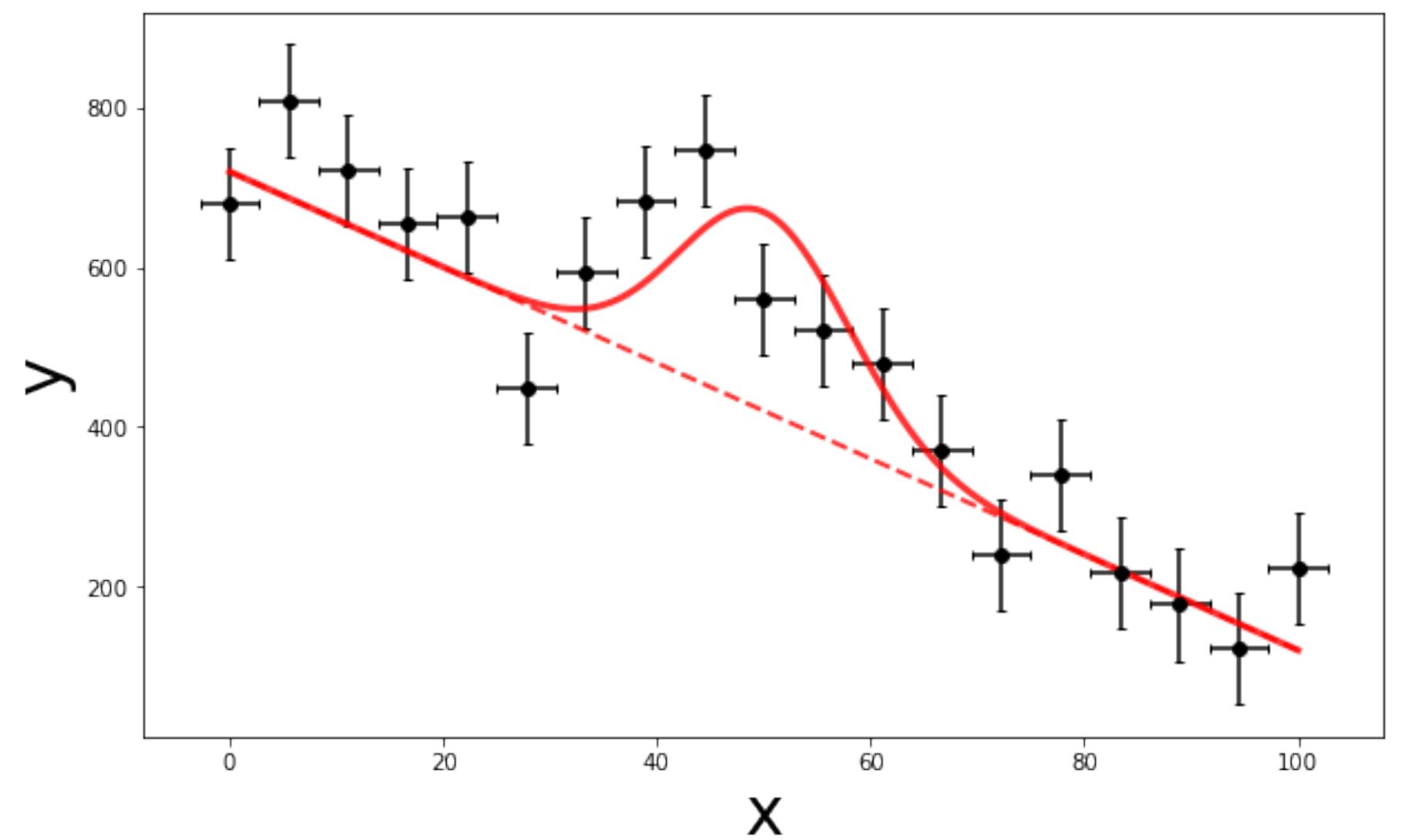
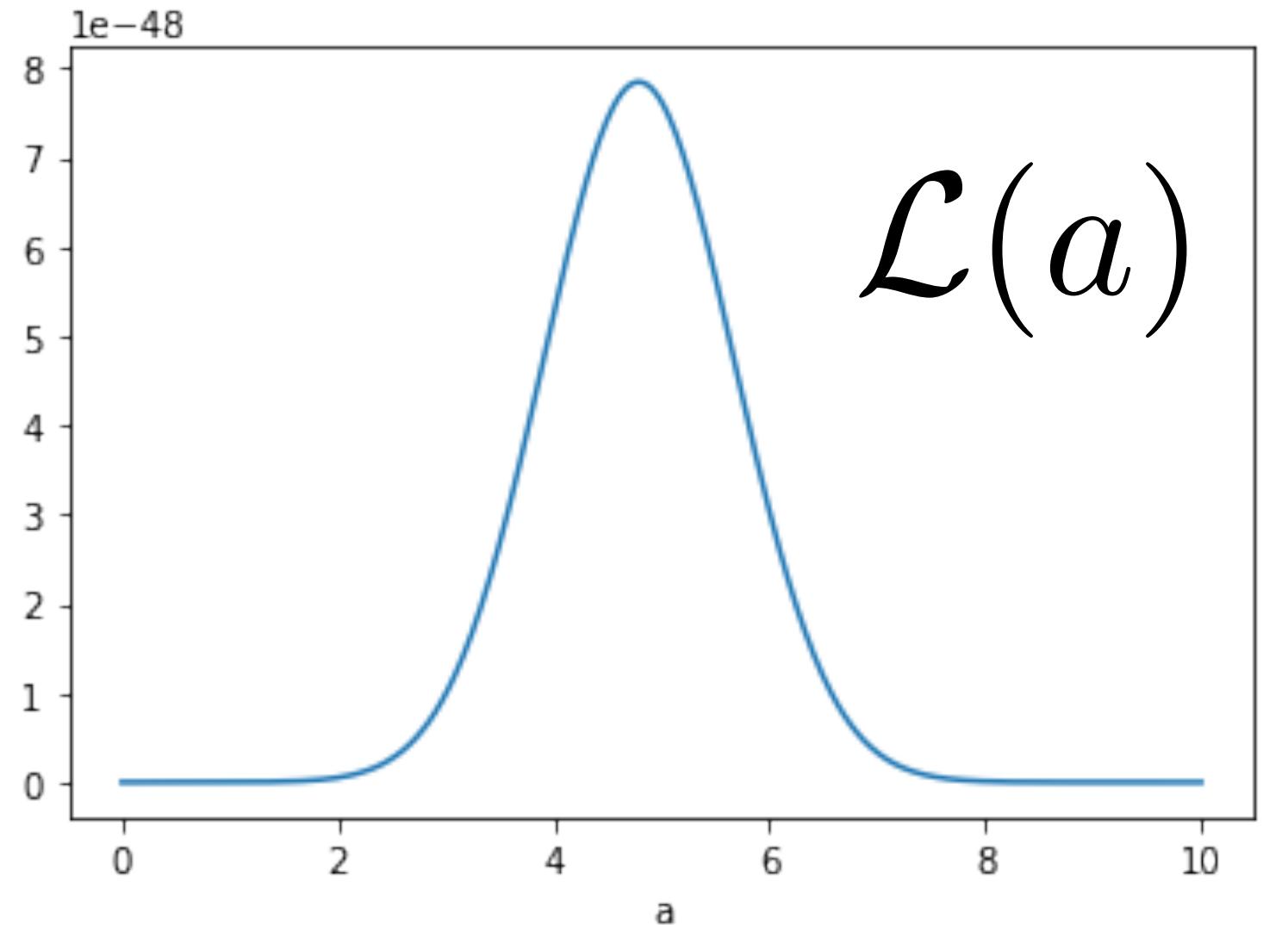
Example:

$$y' = mx + q + a \cdot \mathcal{G}(x; \mu = 50, \sigma = 8)$$

Likelihood

$$\mathcal{L}(a) \equiv p(\vec{x}, \vec{y}|a) = \prod_i p(x_i, y_i|a)$$

$$p(x_i, y_i|a) \propto e^{-\frac{1}{2} \left(\frac{y'_i(a) - y_i}{\sigma} \right)^2}$$



Null hypothesis H_0

$$a = 0$$

Alternative hyp. H_1

$$a = 5$$

Example:

$$y' = mx + q + a \cdot \mathcal{G}(x; \mu = 50, \sigma = 8)$$

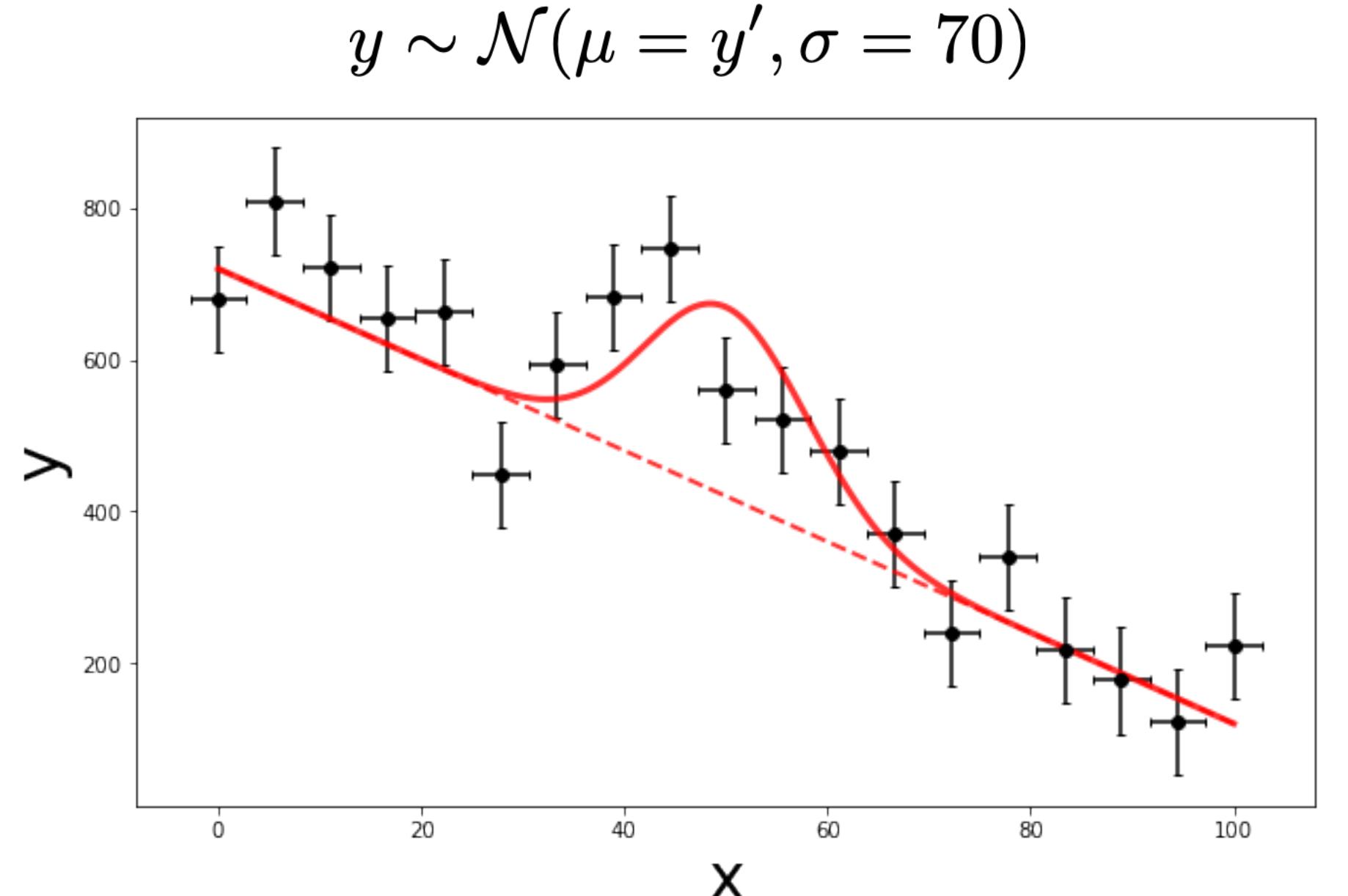
Likelihood

$$\mathcal{L}(a) \equiv p(\vec{x}, \vec{y} | a) = \prod_i p(x_i, y_i | a)$$

$$p(x_i, y_i | a) \propto e^{-\frac{1}{2} \left(\frac{y'_i(a) - y_i}{\sigma} \right)^2}$$

$$\mathcal{S} = \frac{\mathcal{L}(a = 0)}{\mathcal{L}(a = \hat{a})} = 3.52 \cdot 10^{-7}$$

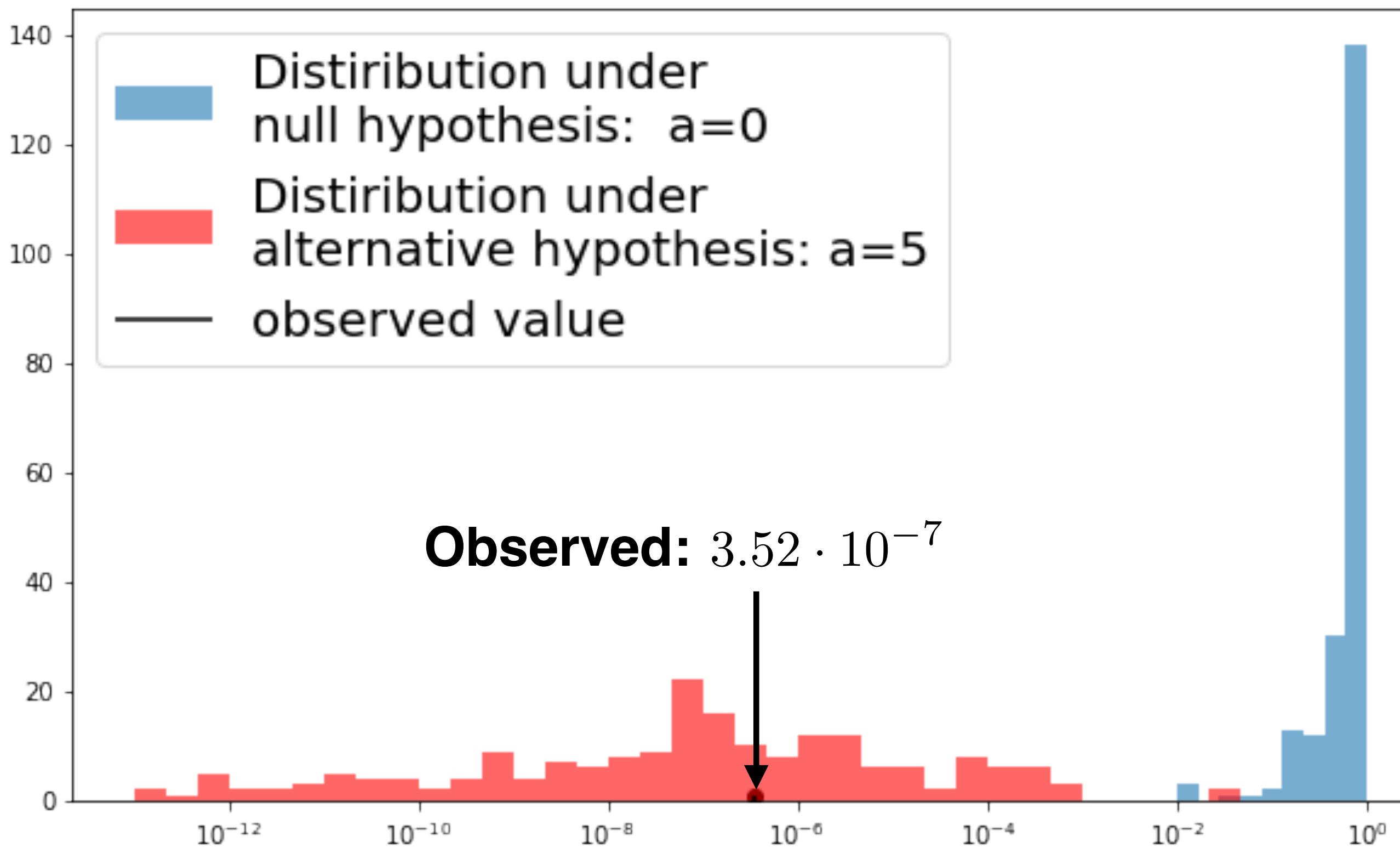
How do we interpret this value of the **statistic**?



Null hypothesis H0
 $a = 0$

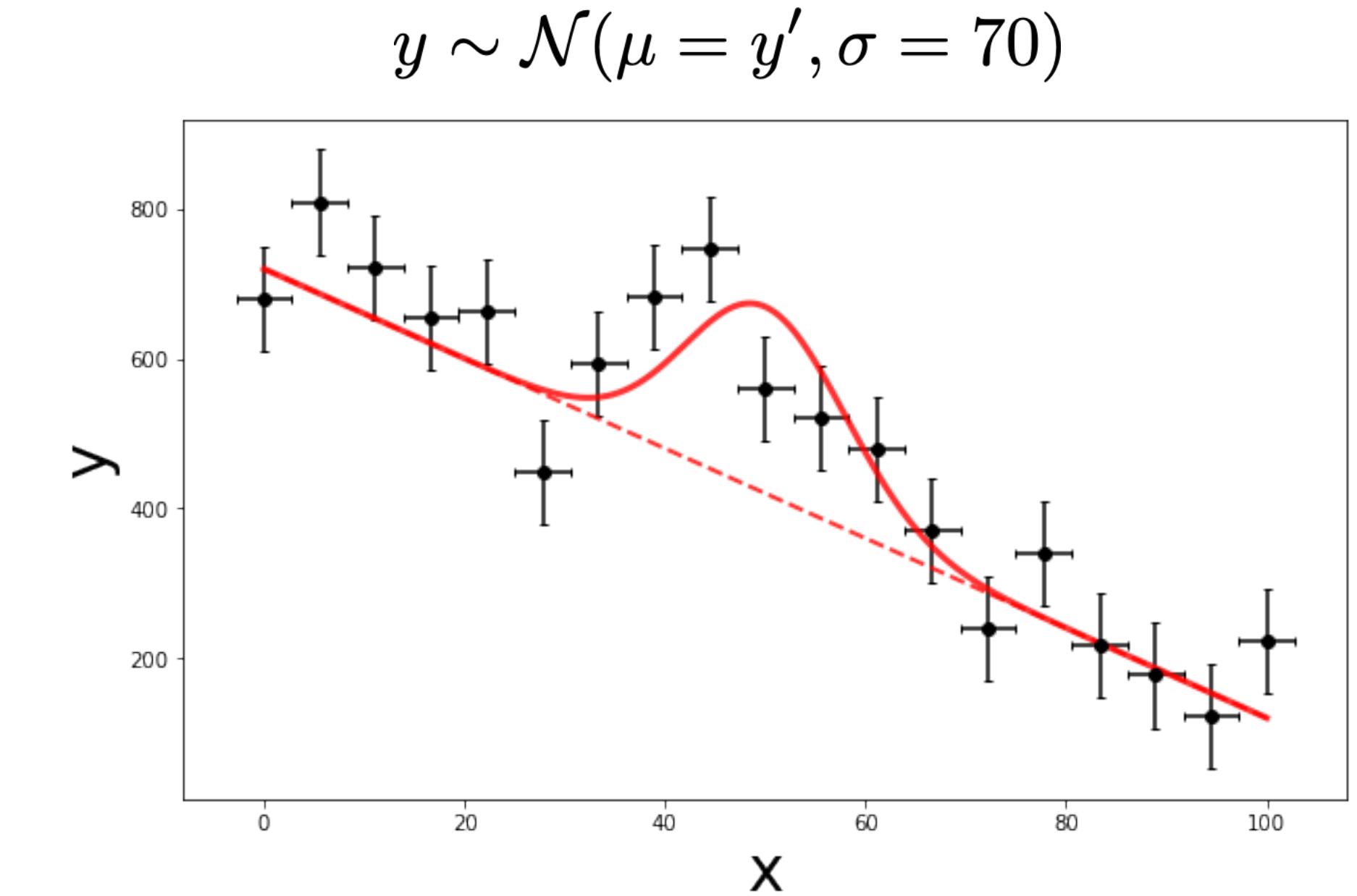
Alternative hyp. H1
 $a = 5$

Example:



$$S = \frac{\mathcal{L}(a = 0)}{\mathcal{L}(a = \hat{a})}$$

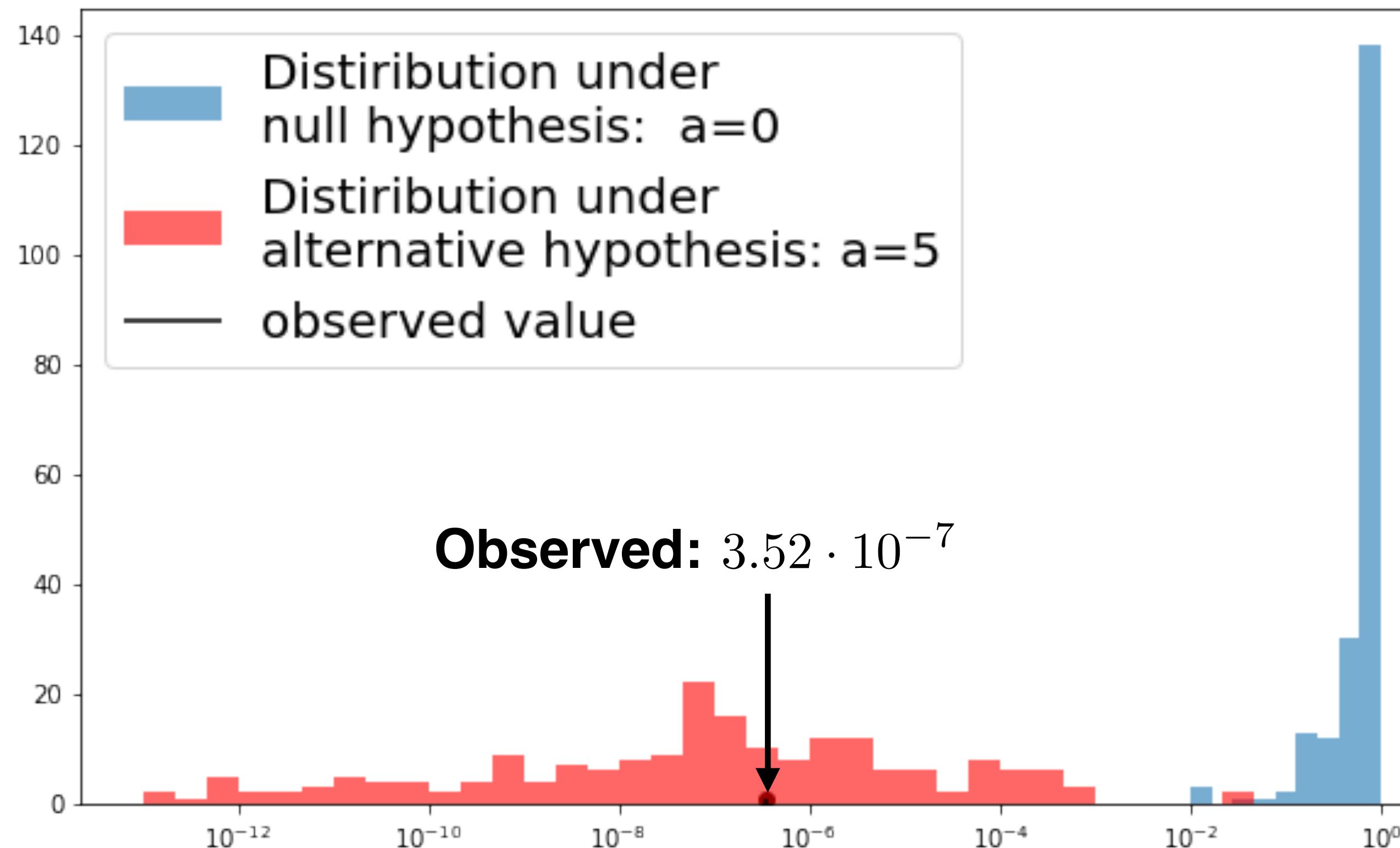
$$y' = mx + q + a \cdot \mathcal{G}(x; \mu = 50, \sigma = 8)$$



Null hypothesis H_0
 $a = 0$

Alternative hyp. H_1
 $a = 5$

Example:



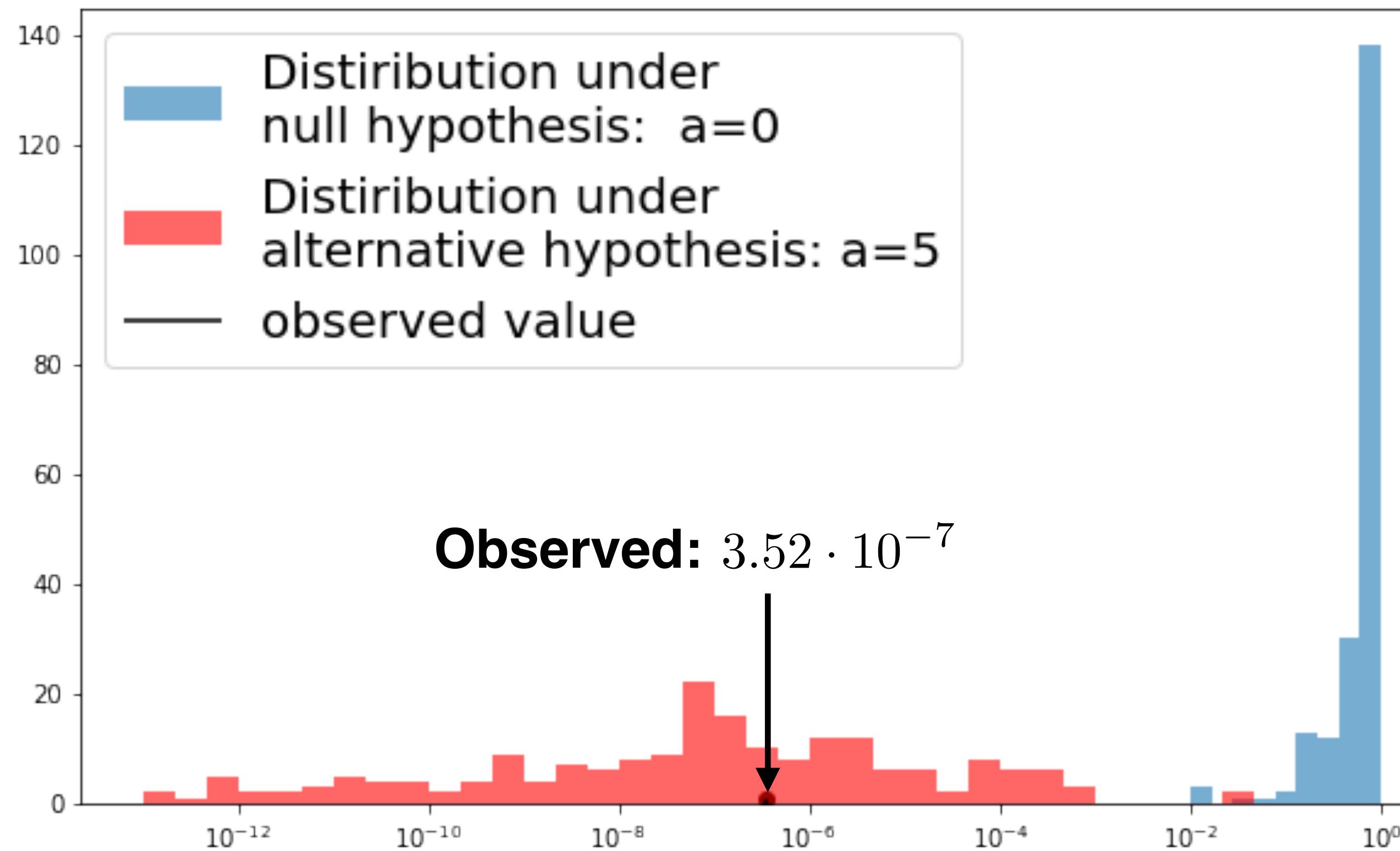
$$S = \frac{\mathcal{L}(a = 0)}{\mathcal{L}(a = \hat{a})}$$

Such a value of the **statistic** is more luckily to have been produced by the **alternative hypothesis** rather than by the **null hypothesis**!

Therefore, we can exclude the null hypothesis and be quite sure of avoiding a type I error.

But with what confidence?

Example:

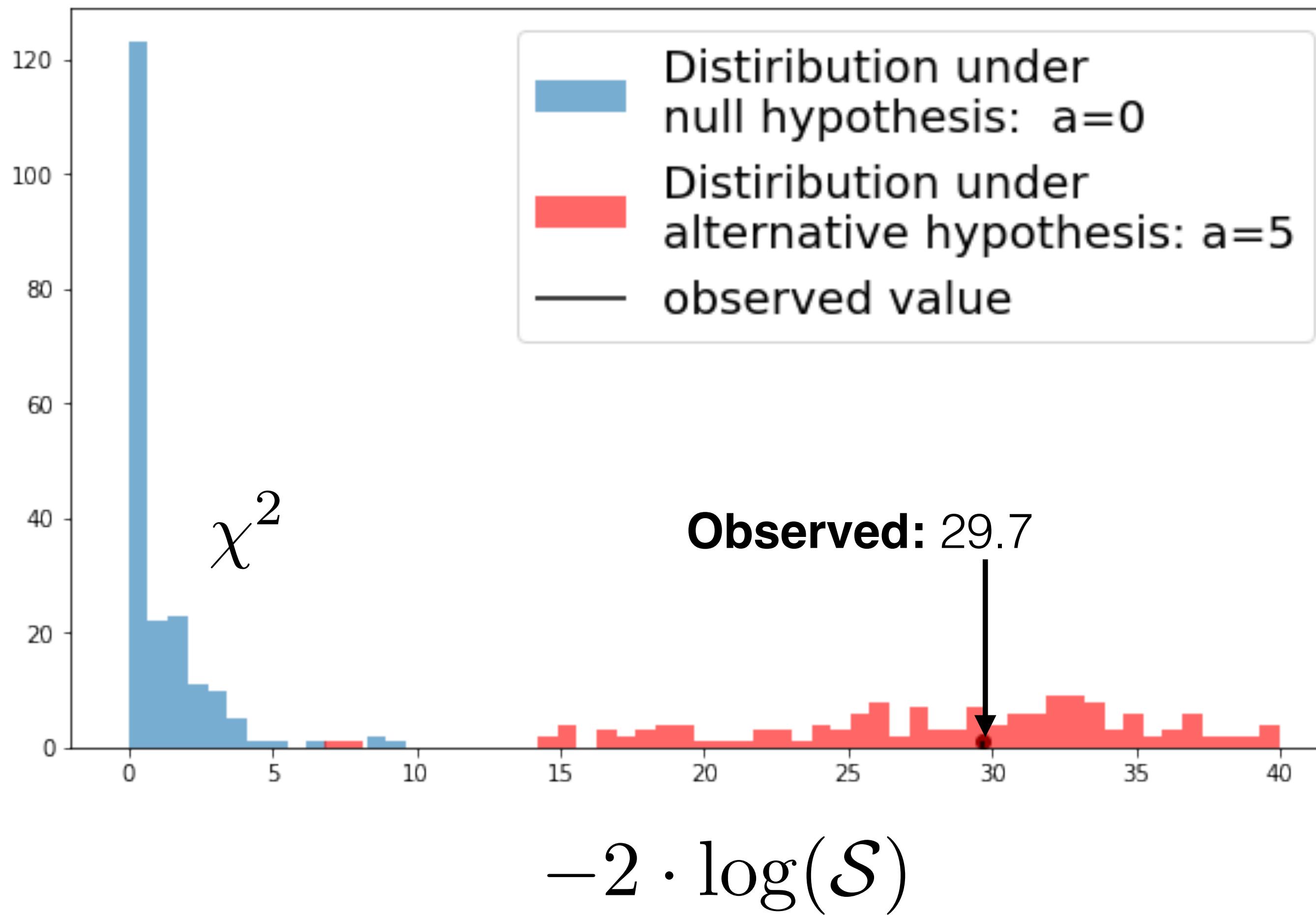


$$\mathcal{S} = \frac{\mathcal{L}(a = 0)}{\mathcal{L}(a = \hat{a})}$$

Taking the $-2 \cdot \log(\mathcal{S})$
the **blue** distribution becomes a
 χ^2 distribution

This is known as the
Wilks' theorem

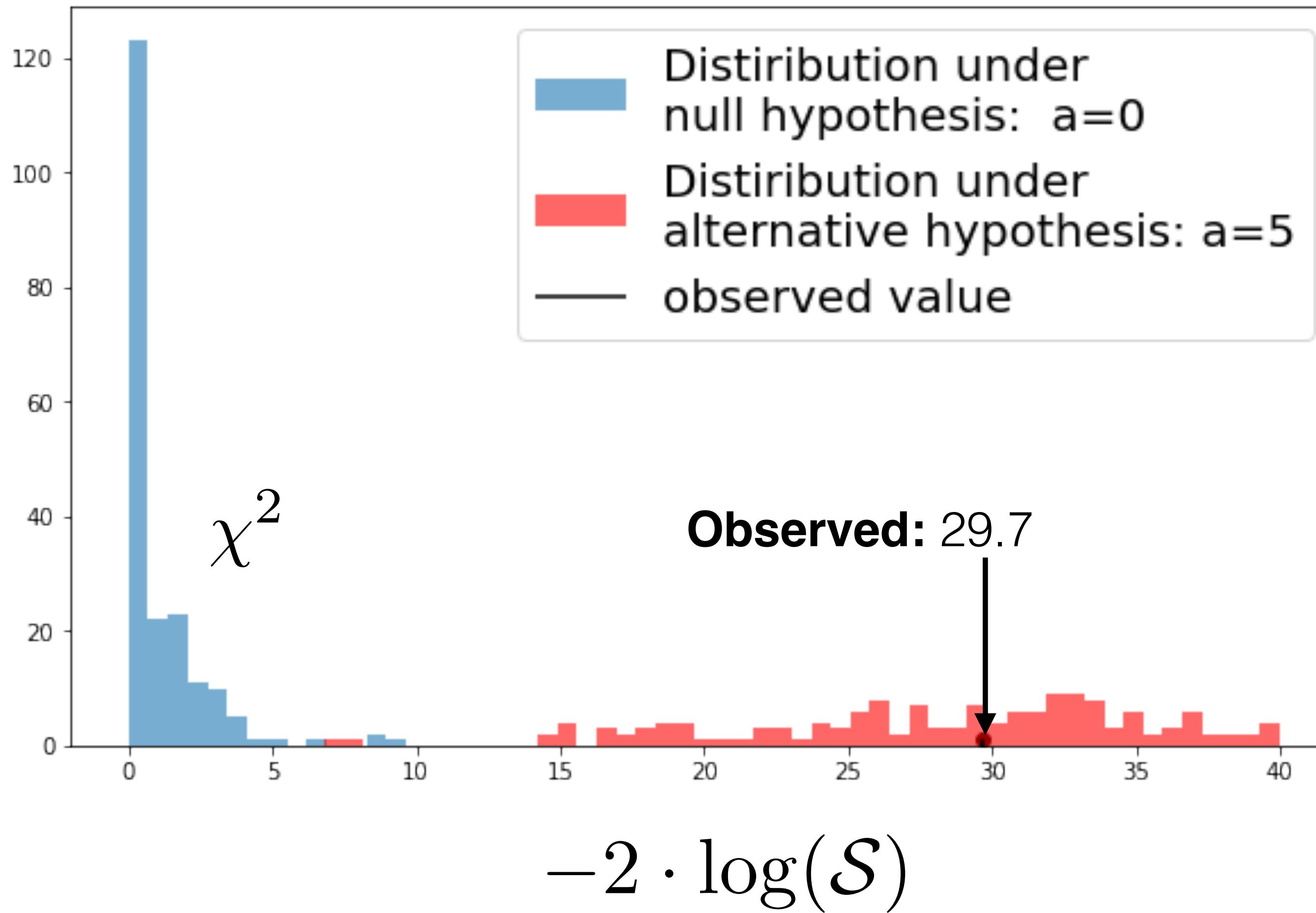
Example:



Taking the $-2 \cdot \log(\mathcal{S})$
the **blue** distribution becomes a
 χ^2 distribution

This is known as the
Wilks' theorem

Example:



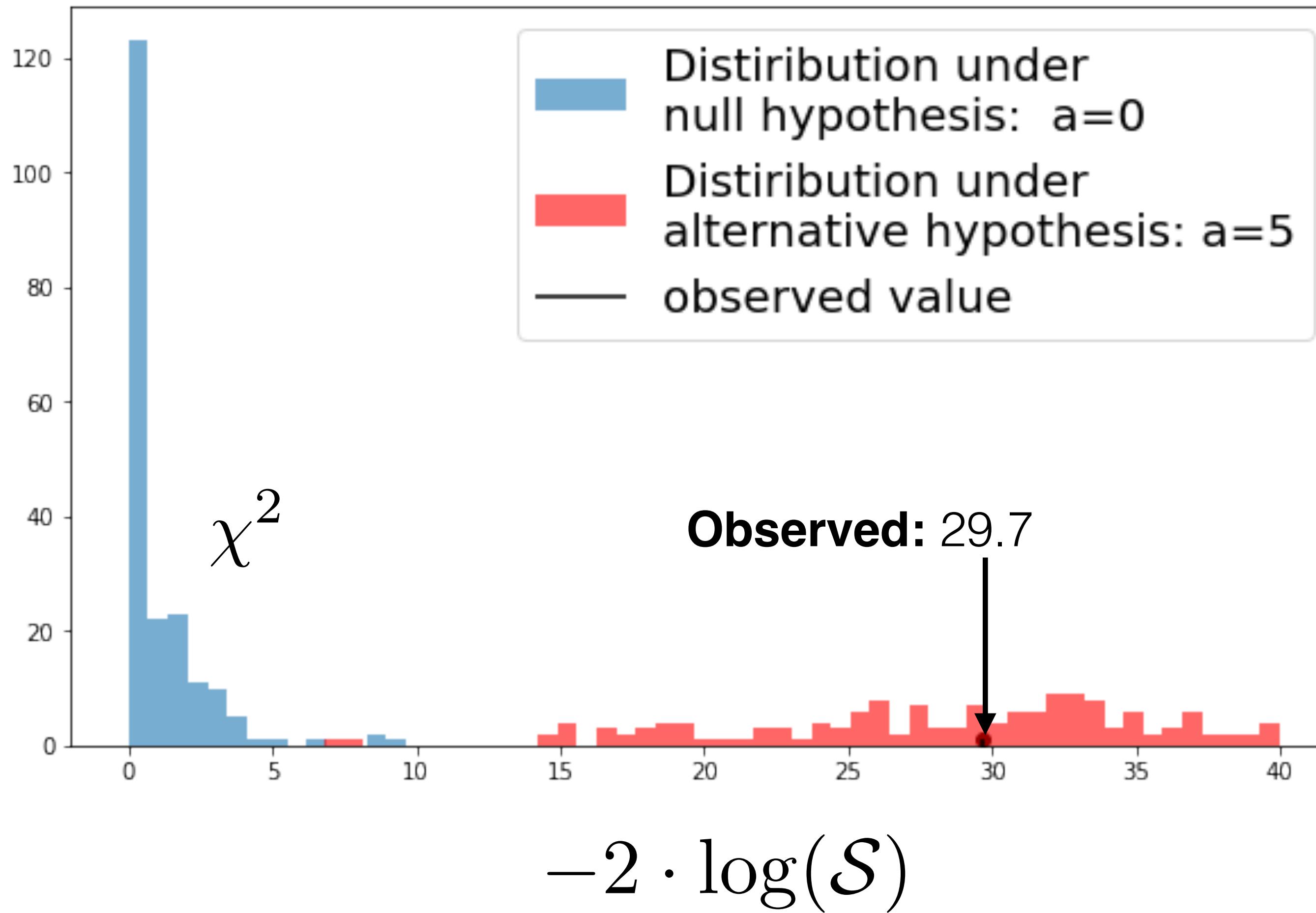
$$p\text{-value} = \int_{29.7}^{\infty} dx \chi^2(x) \simeq 5 \cdot 10^{-8}$$

Converting the p-value to a “sigma”

$$\sqrt{2} \cdot \text{erf}^{-1}(1 - 5 \cdot 10^{-8}) \simeq 5.45$$

We are above the 5 sigmas, we can therefore claim a **discovery!**

Example:



$$p\text{-value} = \int_{29.7}^{\infty} dx \chi^2(x) \simeq 5 \cdot 10^{-8}$$

Converting the p-value to a
“sigma”

$$\sqrt{2} \cdot \text{erf}^{-1}(1 - 5 \cdot 10^{-8}) \simeq 5.45$$

Notice that $\sqrt{29.7} \simeq 5.45$
Why?

Recap:

1. The **Bayesian** approach allows us to quantify our “opinion” on a given model from the observed data using the rules of **probability theory**
 - **Pros:** Alternative hypotheses are taken into account. No need to define a statistic and to know its distribution.
 - **Cons:** One needs a prior distribution.
2. The **frequentist** approach makes us exclude a model with given confidence by looking at infinity repetitions of the experiments in which the model is assumed to be true
 - **Pros:** No need for priors
 - **Cons:** Choice of the statistic is arbitrary. Alternative hypothesis not taken into account. Type I and II errors.