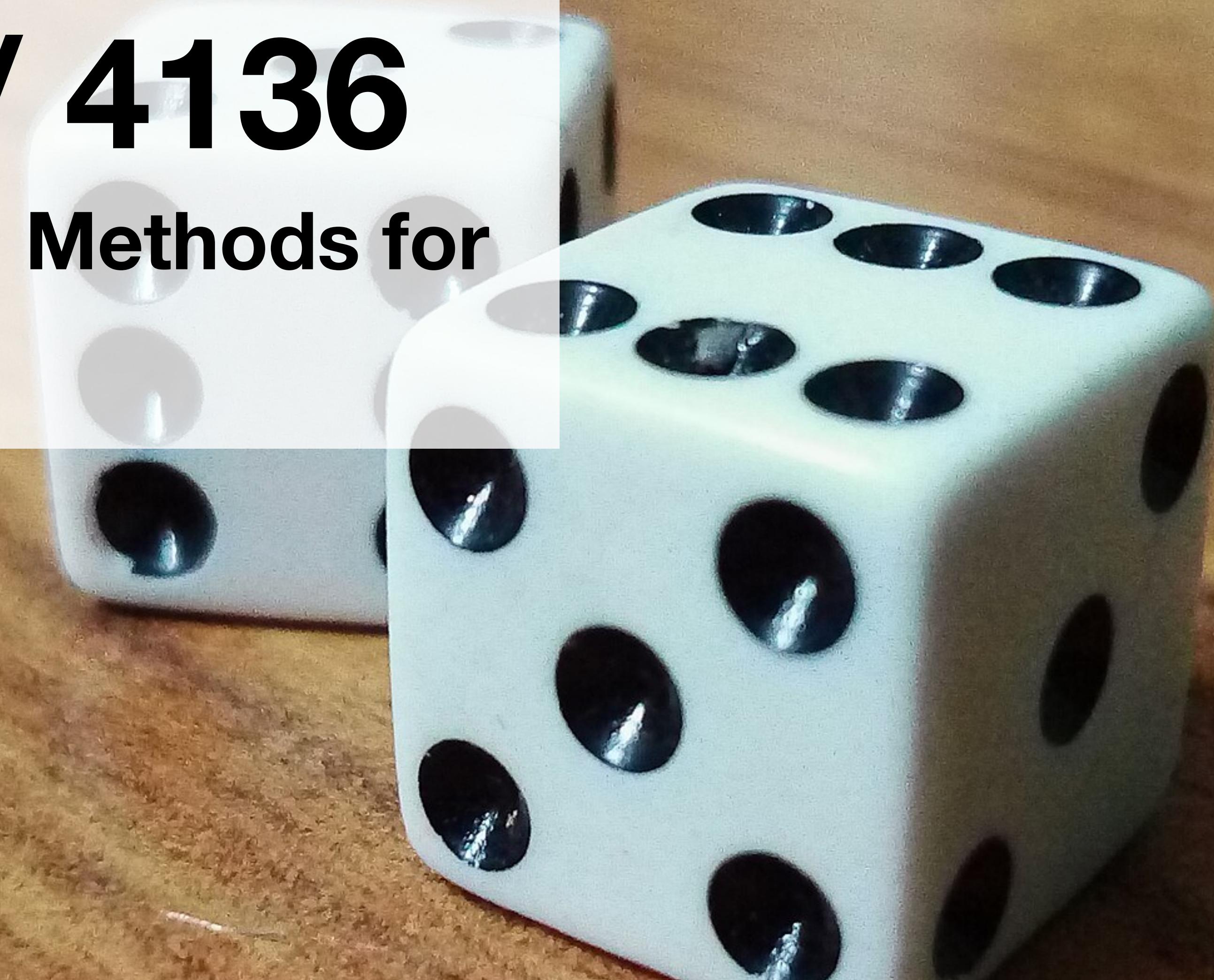


# **LING2136 / 4136**

## **Advanced Statistical Methods for Language Students**



### **Session-05: Bayes**

Lecturer: Timo Roettger

# REASONING ABOUT PROBABILITY







$P(\text{"Rain"}) = 0$

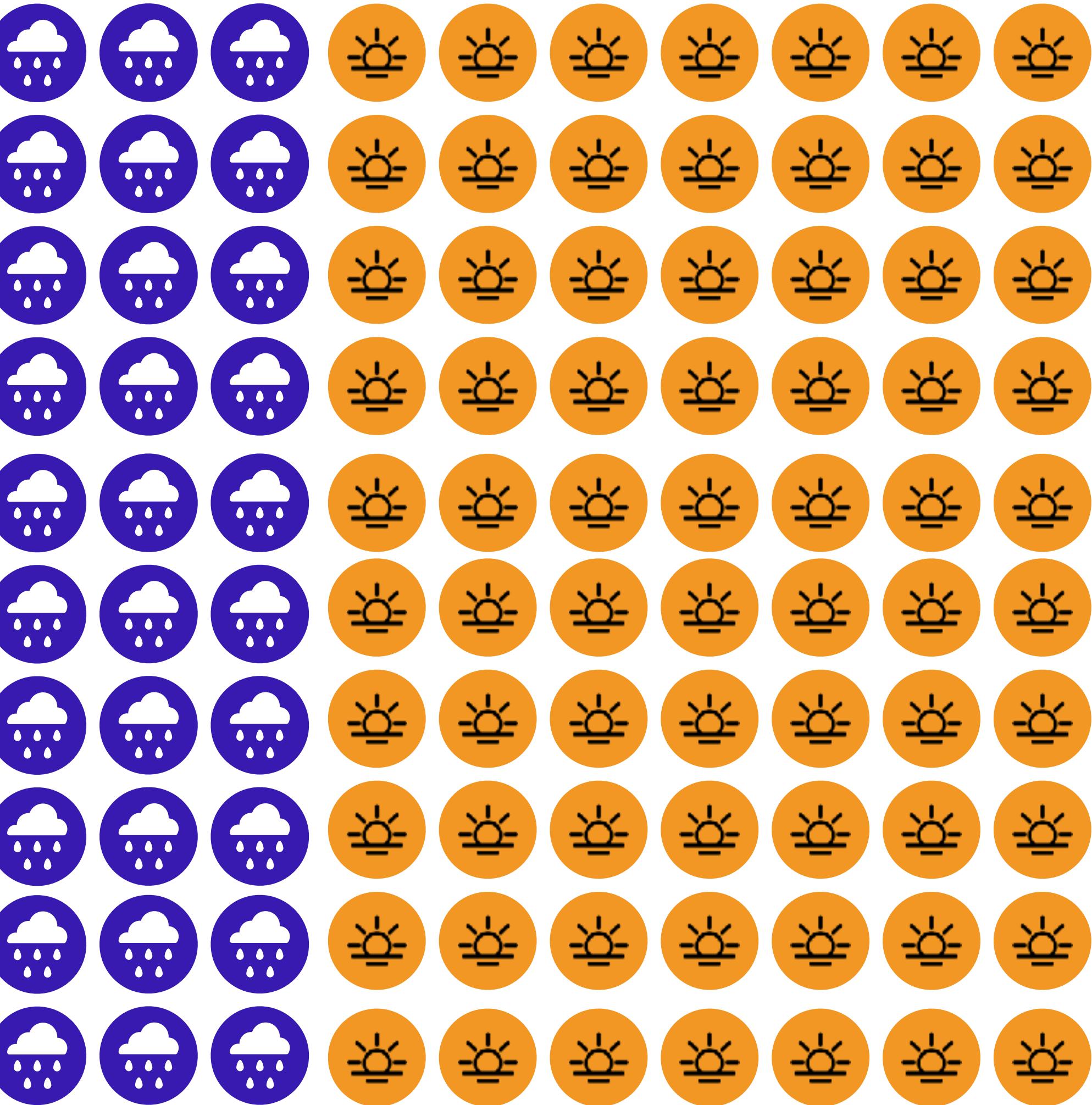
absolutely **certain** that it won't rain

$P(\text{"Rain"}) = 1$

absolutely **certain** that it rains

$P(\text{Rain}) = 0.3$

$P(\neg \text{Rain}) = 0.7$



I observe evidence: 

$$P(\text{Rain}) = 0.3$$



$$P(\neg\text{Rain}) = 0.7$$

$$P(\text{Clouds} \mid \text{Rain}) = 0.8$$

$$P(\text{Clouds} \mid \neg\text{Rain}) = 0.2$$

**DO I CHANGE MY BELIEF?**



I observe evidence: 

$$P(\text{Rain} \mid \text{Clouds}) = \frac{24}{24 + 14} \approx 0.63$$

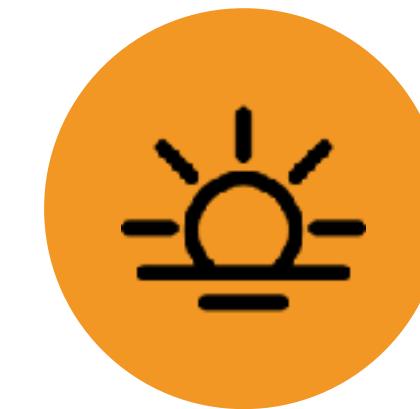
$$P(\text{Clouds} \mid \text{Rain}) = 0.8$$

$$P(\text{Rain}) = 0.3$$

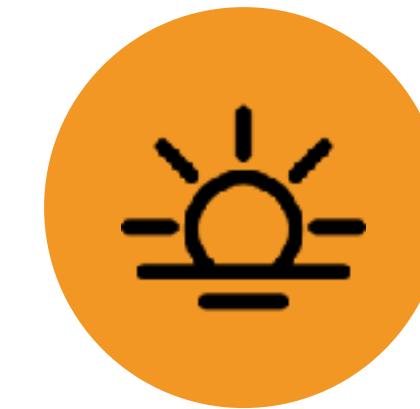
$$P(\neg \text{Rain}) = 0.7$$



$$P(\text{Clouds} \mid \neg \text{Rain}) = 0.2$$



EVIDENCE  
↓



# When to use Bayes' theorem

You have a  
**hypothesis:**



It's gonna rain

You have observed  
evidence:



There are dark clouds

You want to  
know:

$$P(\text{H} | \text{E})$$

$$P(\text{hypothesis} \text{ given } \text{evidence})$$

**“Prior”** →  $P(\text{ Hypothesis }) = 30 / 100 = 0.3$



**“Prior”** →  $P(\text{Hypothesis}) = 30 / 100 = 0.3$

**“Likelihood”**  
 $P(\text{Evidence} | \text{H}) = 0.8$



**“Prior”** →  $P(\text{Hypothesis}) = 30 / 100 = 0.3$

**“Likelihood”**  
 $P(\text{Evidence} | \text{H}) = 0.8$



$\neg$  means “not”  
 $P(\text{E} | \neg \text{H}) = 0.2$

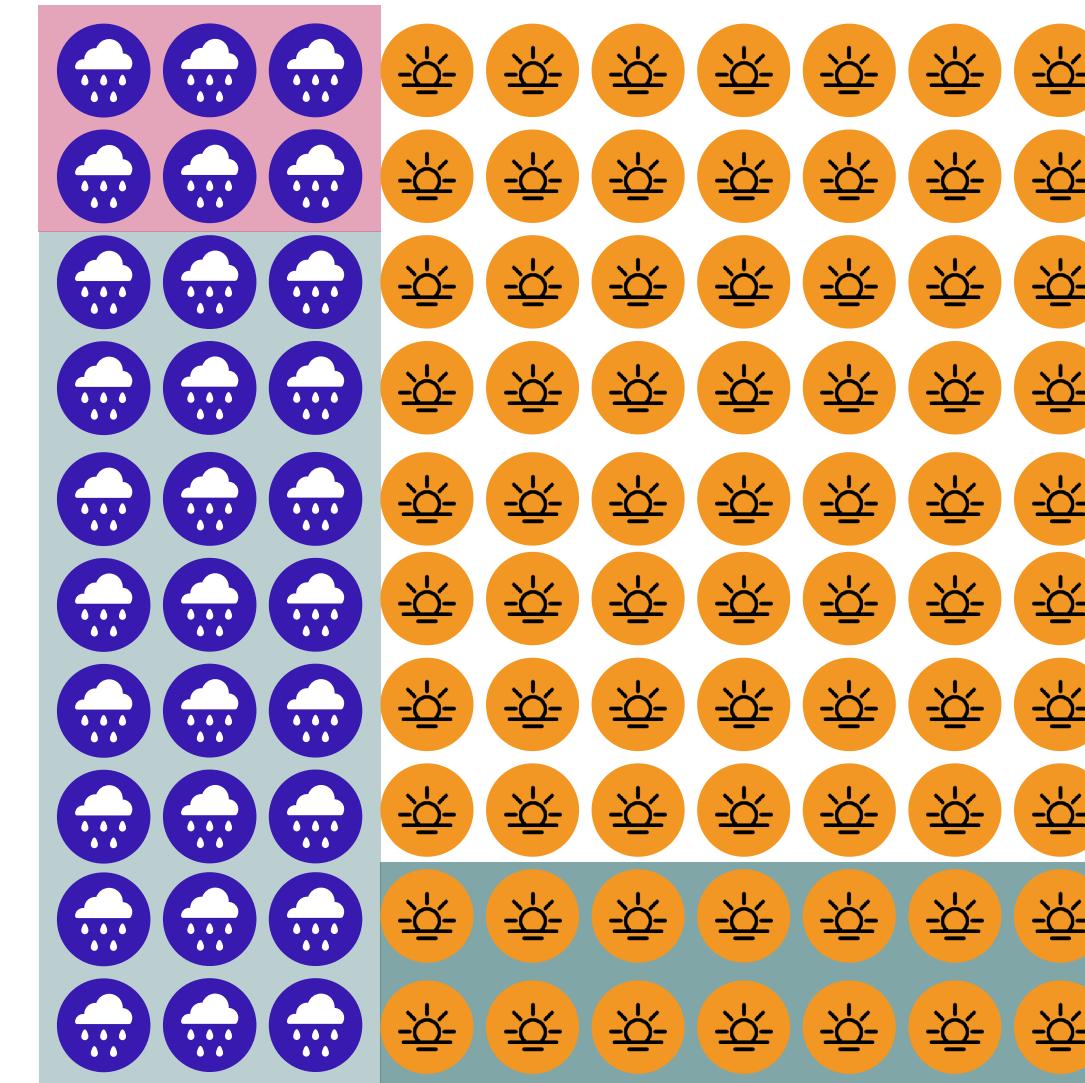
$$P(H|E) = \frac{0.3 * 0.8}{0.3 * 0.8 + 0.7 * 0.2} = 0.38$$

**“Prior” “Likelihood”**

$$= \frac{P(H) P(E|H)}{P(H) P(E|H) + P(\neg H) P(E|\neg H)}$$

$$P(H) = 0.3$$

$$P(E|H) = 0.8$$



$$P(\neg H) = 0.7$$

$$P(E|\neg H) = 0.2$$

$$P(H|E) = \frac{P(H) P(E|H)}{P(E)} = \frac{P(H) P(E|H)}{P(H)P(E|H) + P(\neg H)P(E|\neg H)}$$

“Prior”   “Likelihood”

“Posterior”

$$P(H) = 0.3$$

$$P(E|H) = 0.8$$



$$P(\neg H) = 0.7$$

$$P(E|\neg H) = 0.2$$

# This is Bayes Theorem

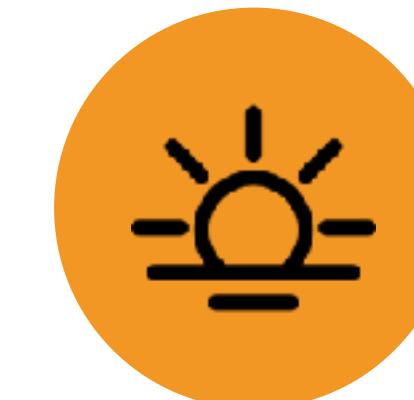
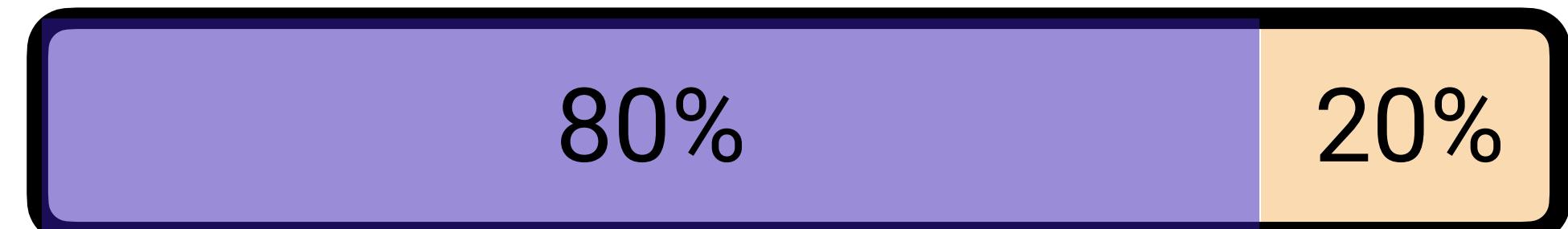
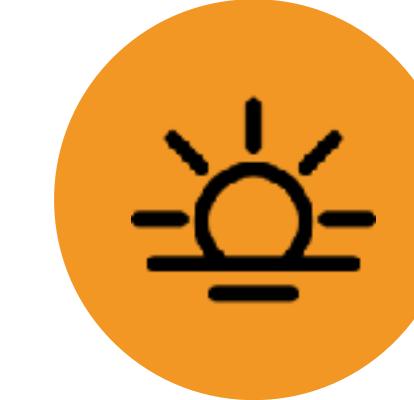
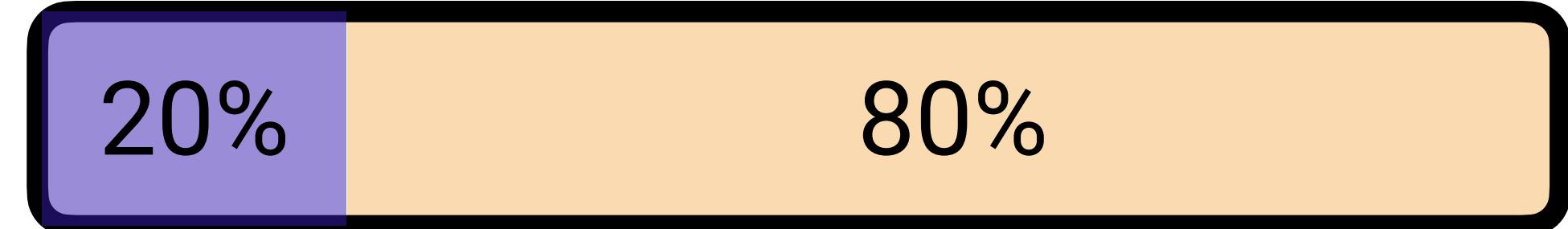
$$P(\mathbf{H} | \mathbf{E}) = \frac{P(\mathbf{H}) P(\mathbf{E} | \mathbf{H})}{P(\mathbf{E})}$$

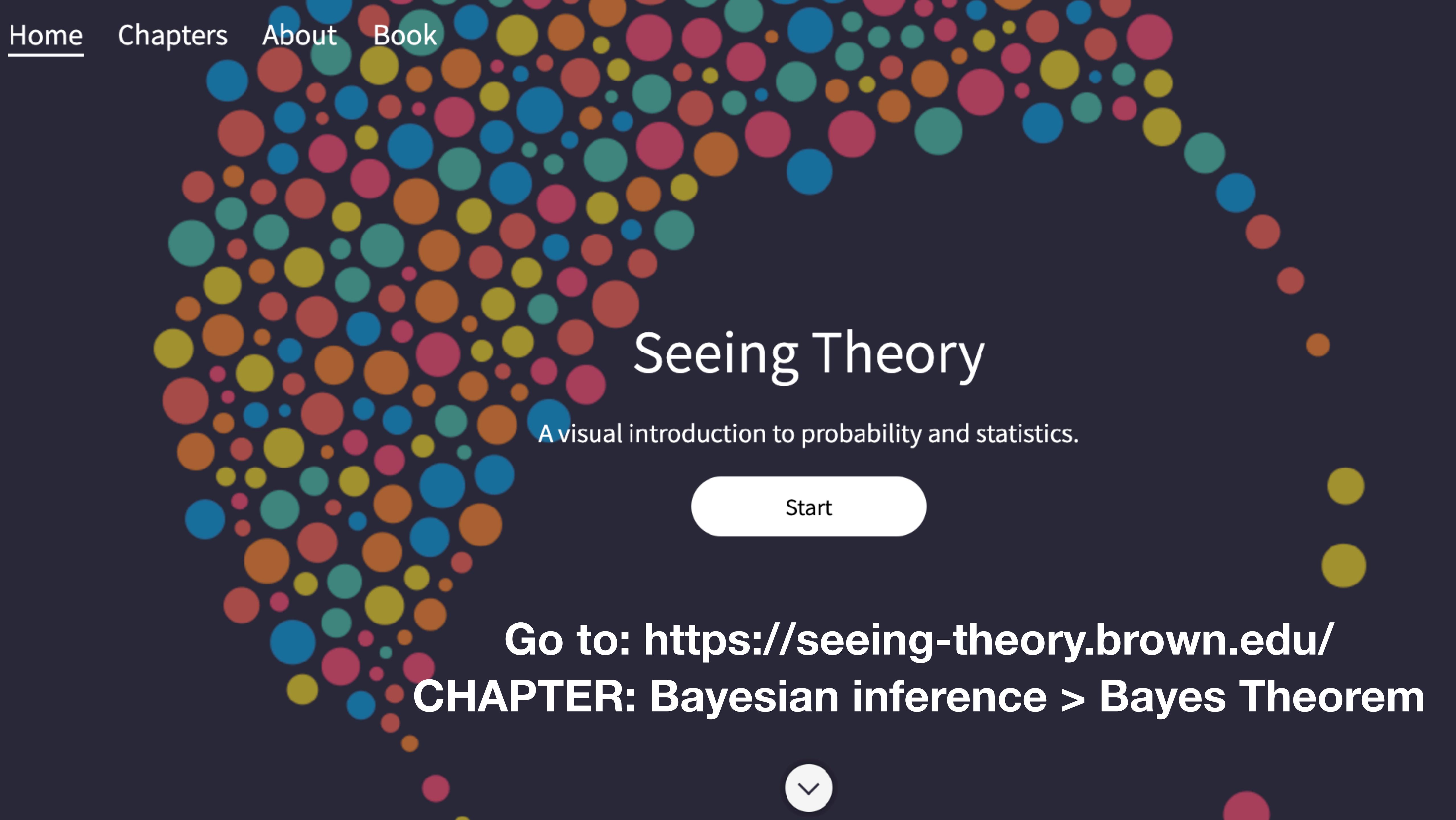
**prior:** initial degree of belief in hypothesis

**likelihood:** the probability of the evidence given the hypothesis

**posterior:** degree of belief in hypothesis, after seeing evidence

The diagram illustrates the components of Bayes' Theorem. It shows the formula  $P(\mathbf{H} | \mathbf{E}) = \frac{P(\mathbf{H}) P(\mathbf{E} | \mathbf{H})}{P(\mathbf{E})}$ . Three curved arrows point from text definitions to specific terms in the formula: one arrow points from the 'prior' definition to  $P(\mathbf{H})$ ; another arrow points from the 'likelihood' definition to  $P(\mathbf{E} | \mathbf{H})$ ; and a third arrow points from the 'posterior' definition to  $P(\mathbf{H} | \mathbf{E})$ .





# Seeing Theory

A visual introduction to probability and statistics.

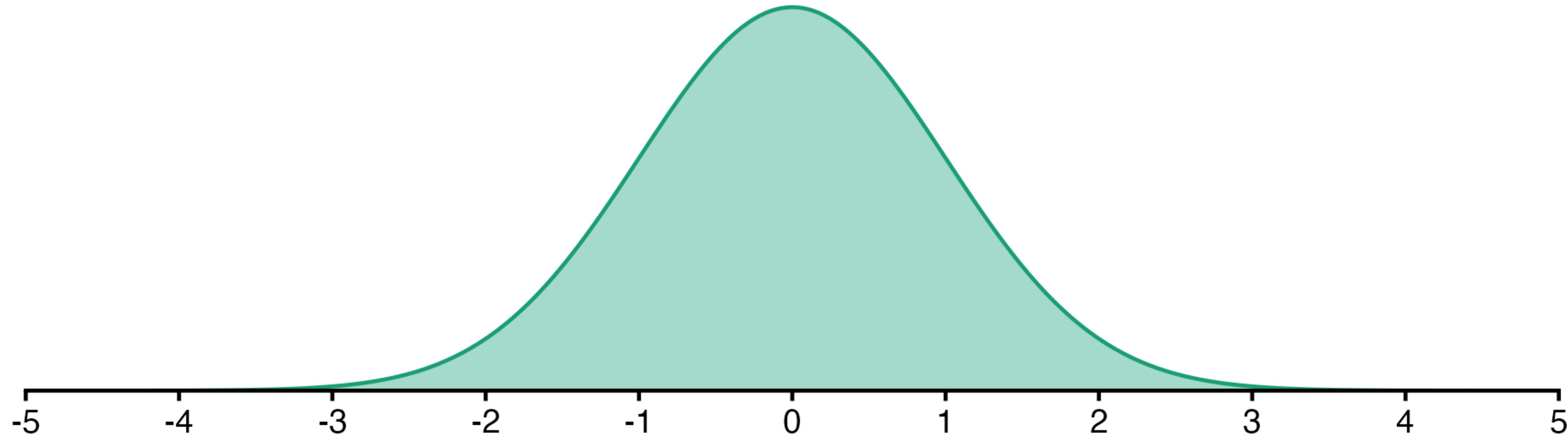
Start

Go to: <https://seeing-theory.brown.edu/>

CHAPTER: Bayesian inference > Bayes Theorem

# The Gaussian (normal) distribution

has parameters: mean ( $\mu$ ) & standard deviation ( $\sigma$ )



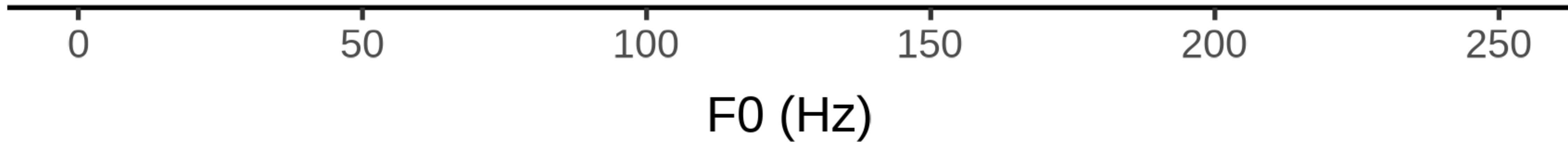
## **Demo with fundamental frequency (F0) data** (corresponds to ‘voice pitch’ in perception)

**data:** average F0 from 32 male speakers

Winter et al. 2021

$f_0 \sim 1$

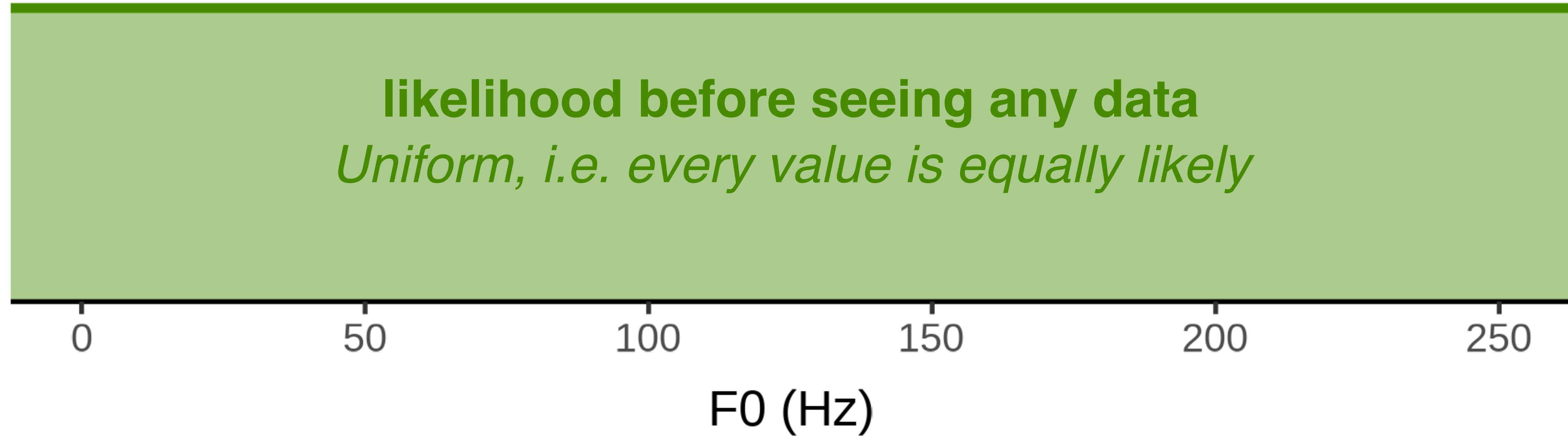
**What we want to know:**  
given some F0 measurements,  
what is the average F0 of a male speaker



## Likelihood functions

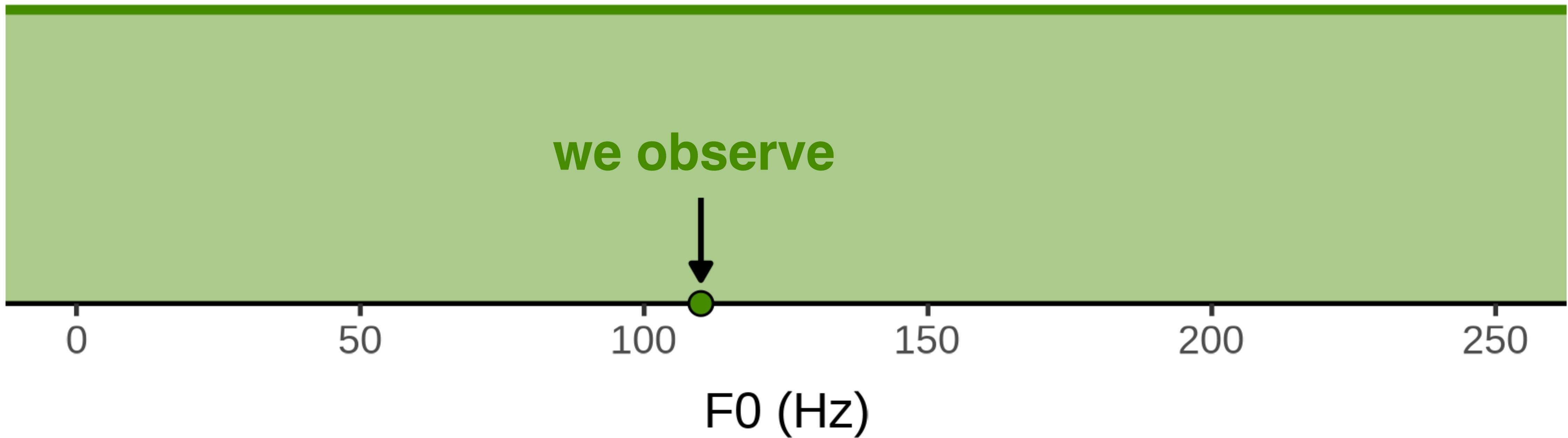
Knowing the data,  
what is the **likelihood** of  
the  $\mu$  parameter?

**likelihood before seeing any data**  
*Uniform, i.e. every value is equally likely*



## Likelihood functions

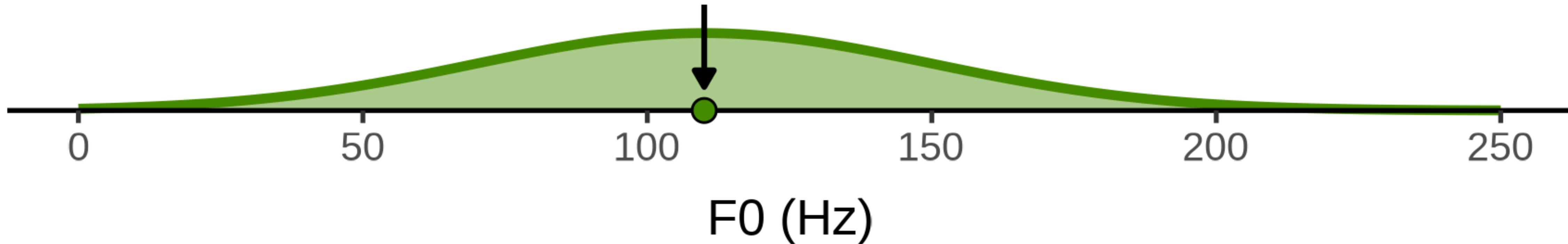
Knowing the data,  
what is the **likelihood** of  
the  $\mu$  parameter?



## Likelihood functions

Knowing the data,  
what is the **likelihood** of  
the  $\mu$  parameter?

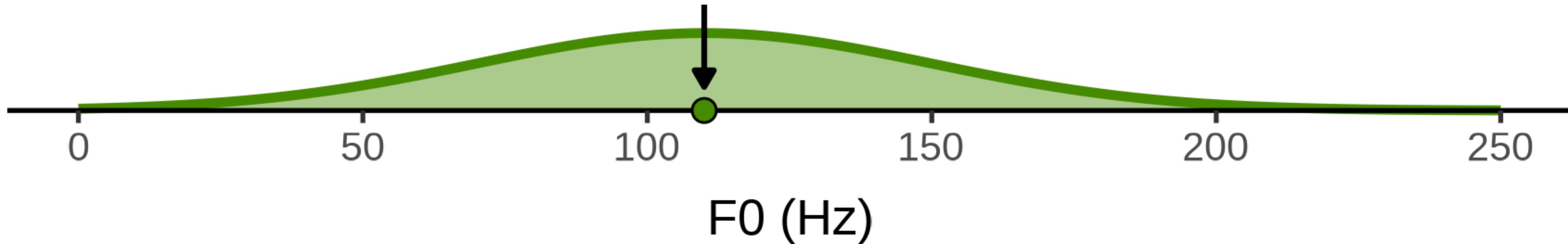
### The likelihood function for $\mu$



## Likelihood functions

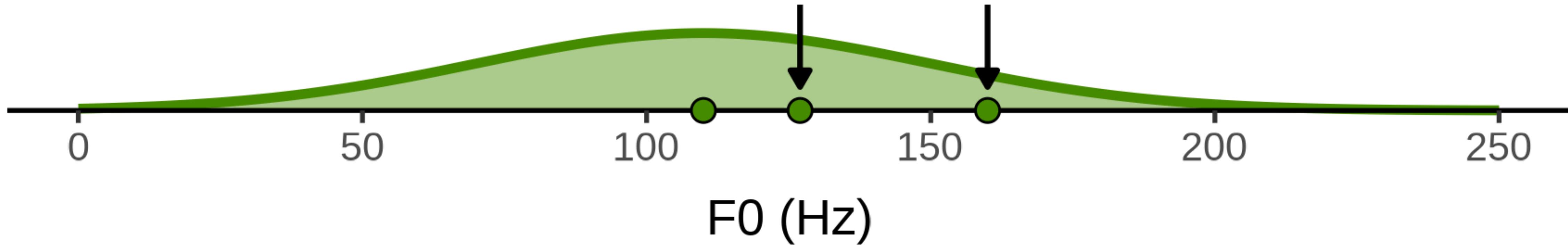
Knowing the data,  
what is the **likelihood** of  
the  $\mu$  parameter?

This one data point is **compatible** with  
a **whole range** of values for  $\mu$



## Likelihood functions

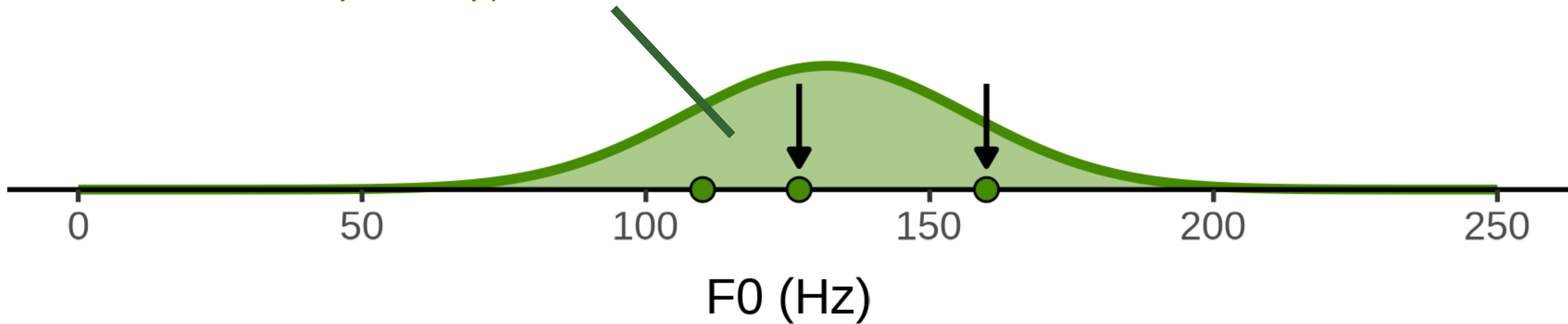
Knowing the data,  
what is the **likelihood** of  
the  $\mu$  parameter?



## Likelihood functions

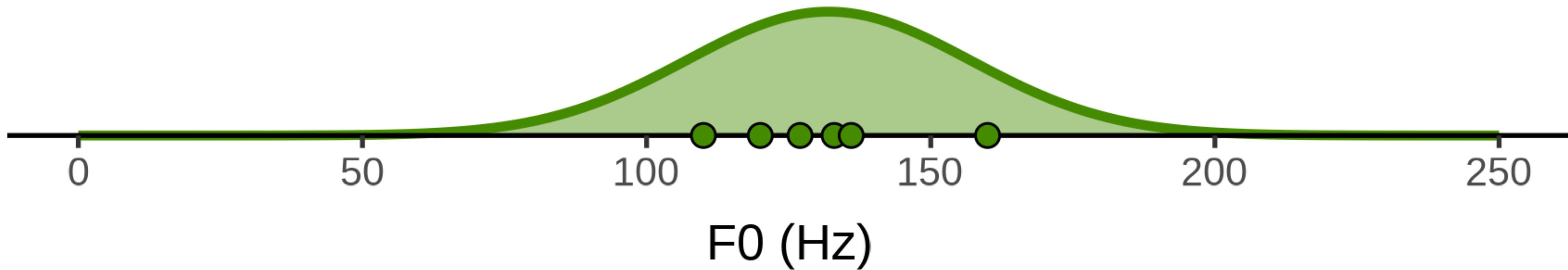
Knowing the data,  
what is the **likelihood** of  
the  $\mu$  parameter?

more data make the likelihood function  
**narrower** (you are accumulating evidence  
for a specific  $\mu$ )



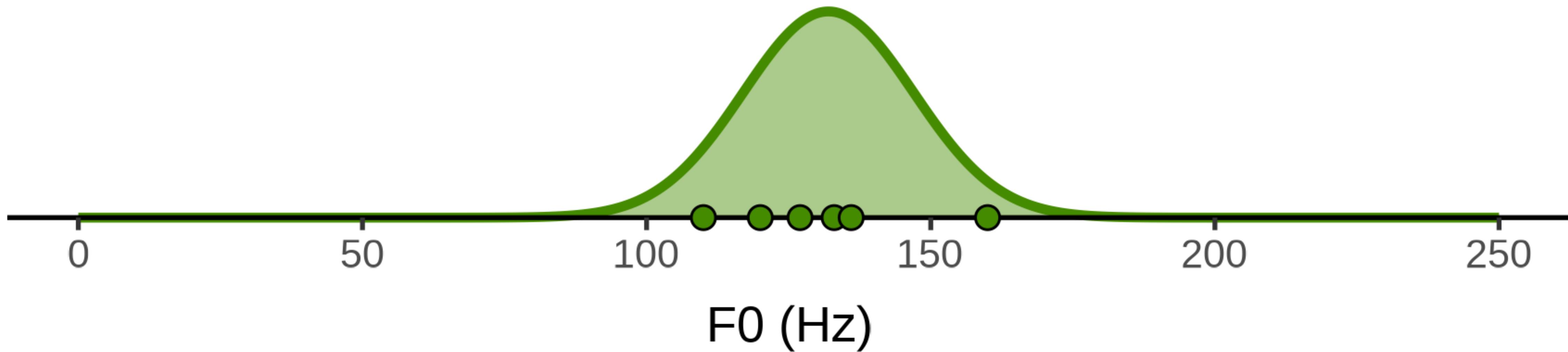
## Likelihood functions

Knowing the data,  
what is the **likelihood** of  
the  $\mu$  parameter?



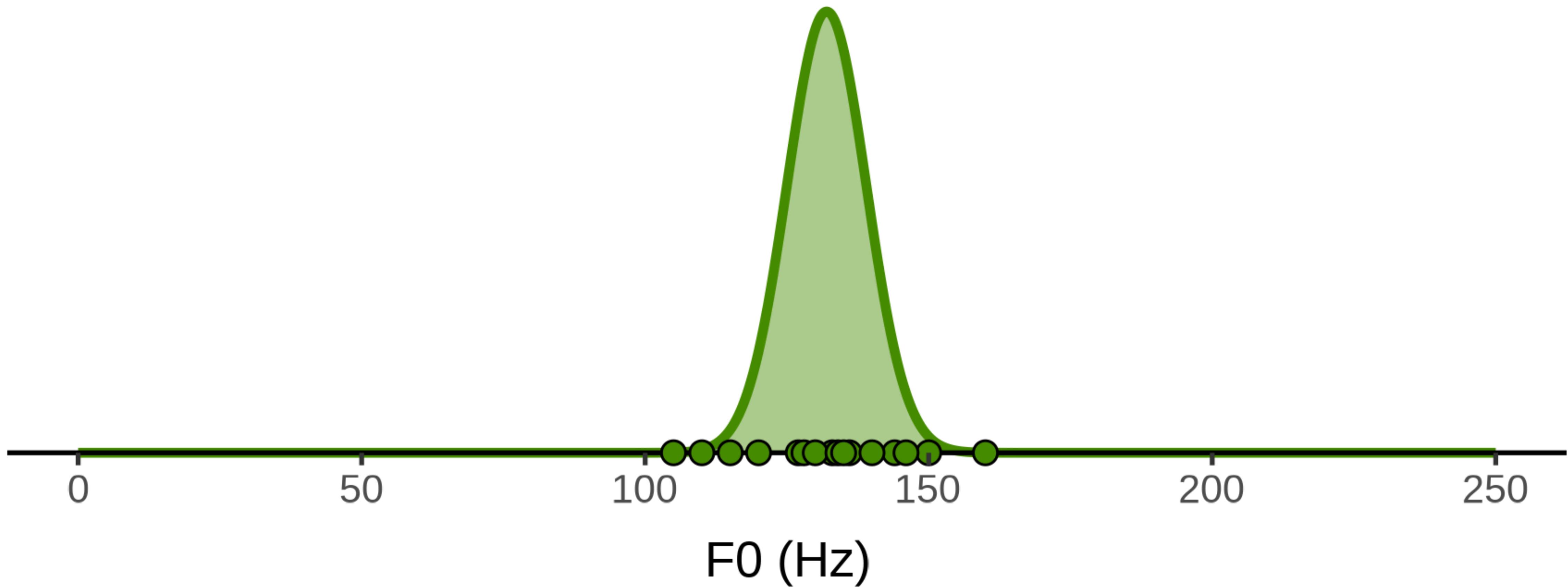
## Likelihood functions

Knowing the data,  
what is the **likelihood** of  
the  $\mu$  parameter?



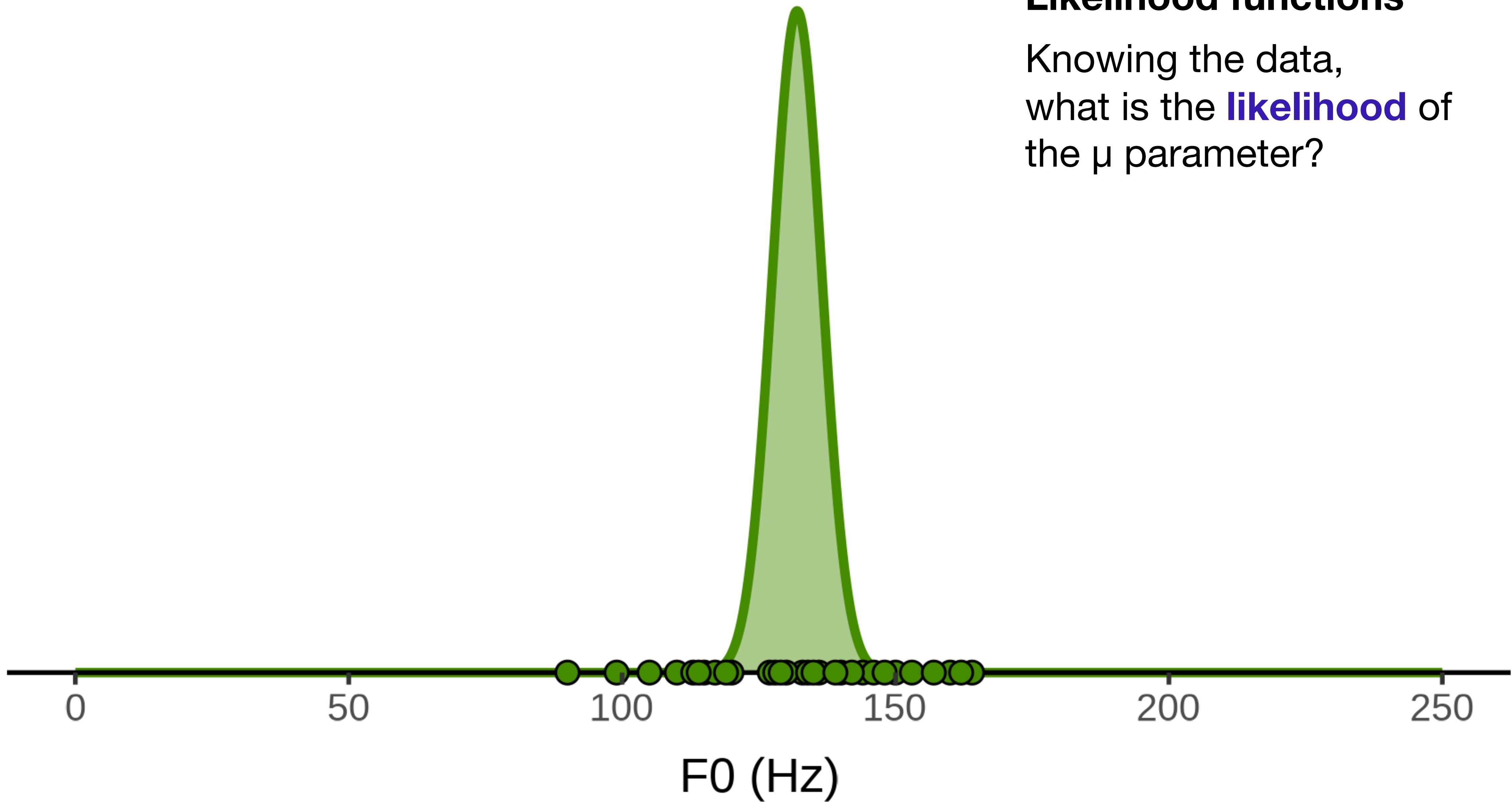
## Likelihood functions

Knowing the data,  
what is the **likelihood** of  
the  $\mu$  parameter?



## Likelihood functions

Knowing the data,  
what is the **likelihood** of  
the  $\mu$  parameter?



# This is Bayes Theorem

$$P(\mathbf{H} | \mathbf{E}) = \frac{P(\mathbf{H}) P(\mathbf{E} | \mathbf{H})}{P(\mathbf{E})}$$

**prior:** initial degree of belief in hypothesis

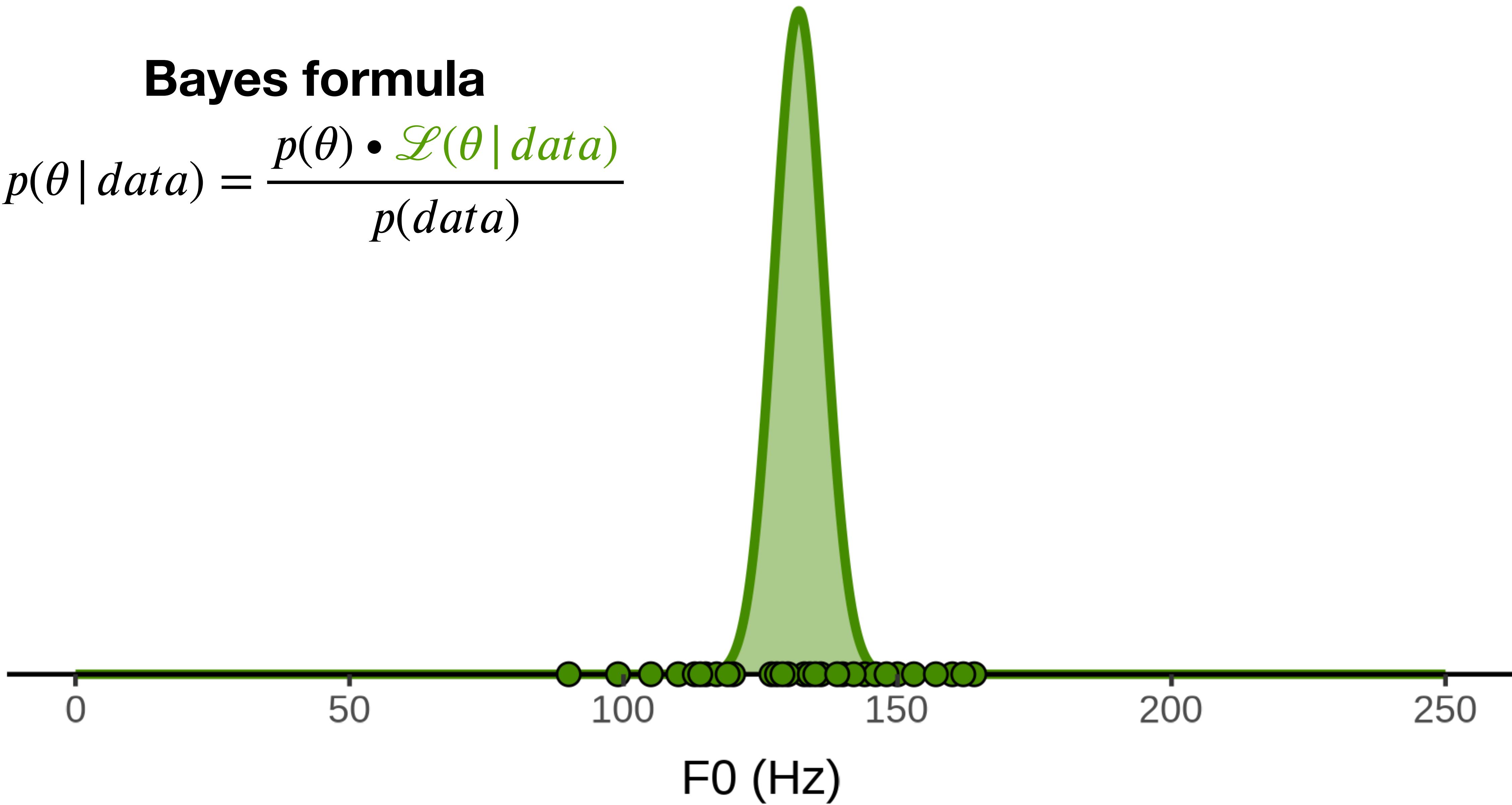
**likelihood:** the probability of the evidence given the hypothesis

**posterior:** degree of belief in hypothesis, after seeing evidence

The diagram illustrates the components of Bayes' Theorem. It shows the formula  $P(\mathbf{H} | \mathbf{E}) = \frac{P(\mathbf{H}) P(\mathbf{E} | \mathbf{H})}{P(\mathbf{E})}$ . Three arrows point from text definitions to specific terms in the formula: one arrow points from the 'prior' definition to  $P(\mathbf{H})$ ; another arrow points from the 'likelihood' definition to  $P(\mathbf{E} | \mathbf{H})$ ; and a third arrow points from the 'posterior' definition to  $P(\mathbf{H} | \mathbf{E})$ .

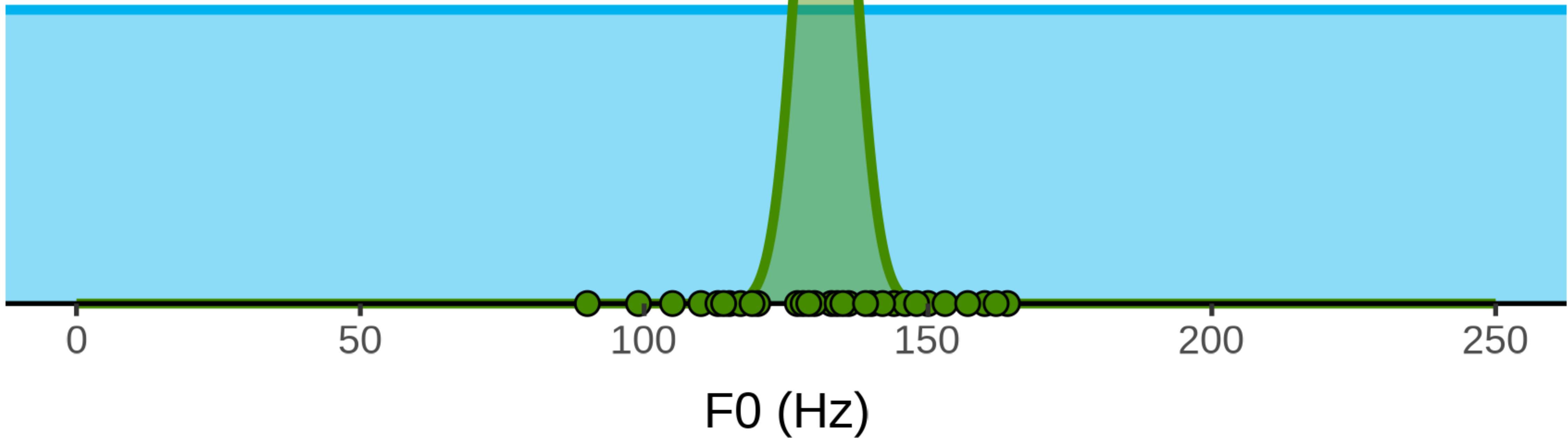
## Bayes formula

$$p(\theta | data) = \frac{p(\theta) \cdot \mathcal{L}(\theta | data)}{p(data)}$$

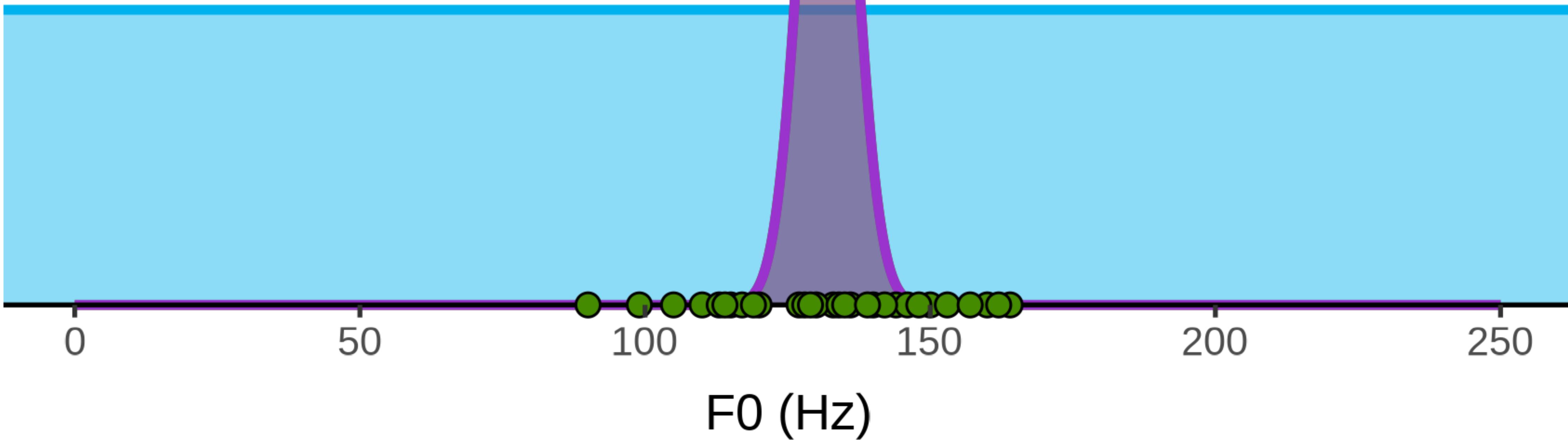


$$p(\theta | data) = \frac{p(\theta) \cdot \mathcal{L}(\theta | data)}{p(data)}$$

$\mu \sim Uniform(-\infty, +\infty)$   
all values are equally probable



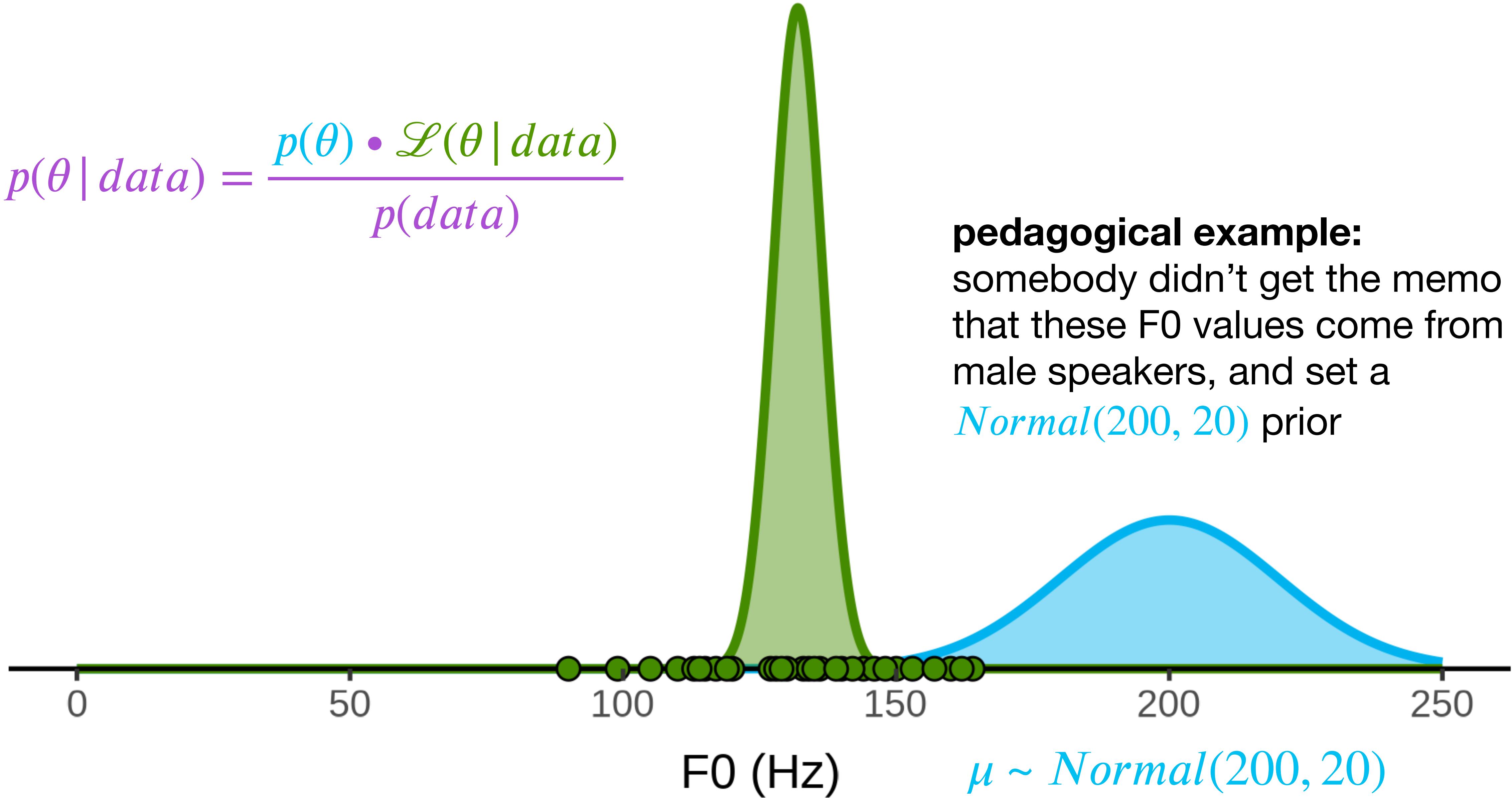
$$p(\theta | data) = \frac{p(\theta) \cdot \mathcal{L}(\theta | data)}{p(data)}$$



**because the flat prior doesn't do anything, the posterior distribution = likelihood distribution**

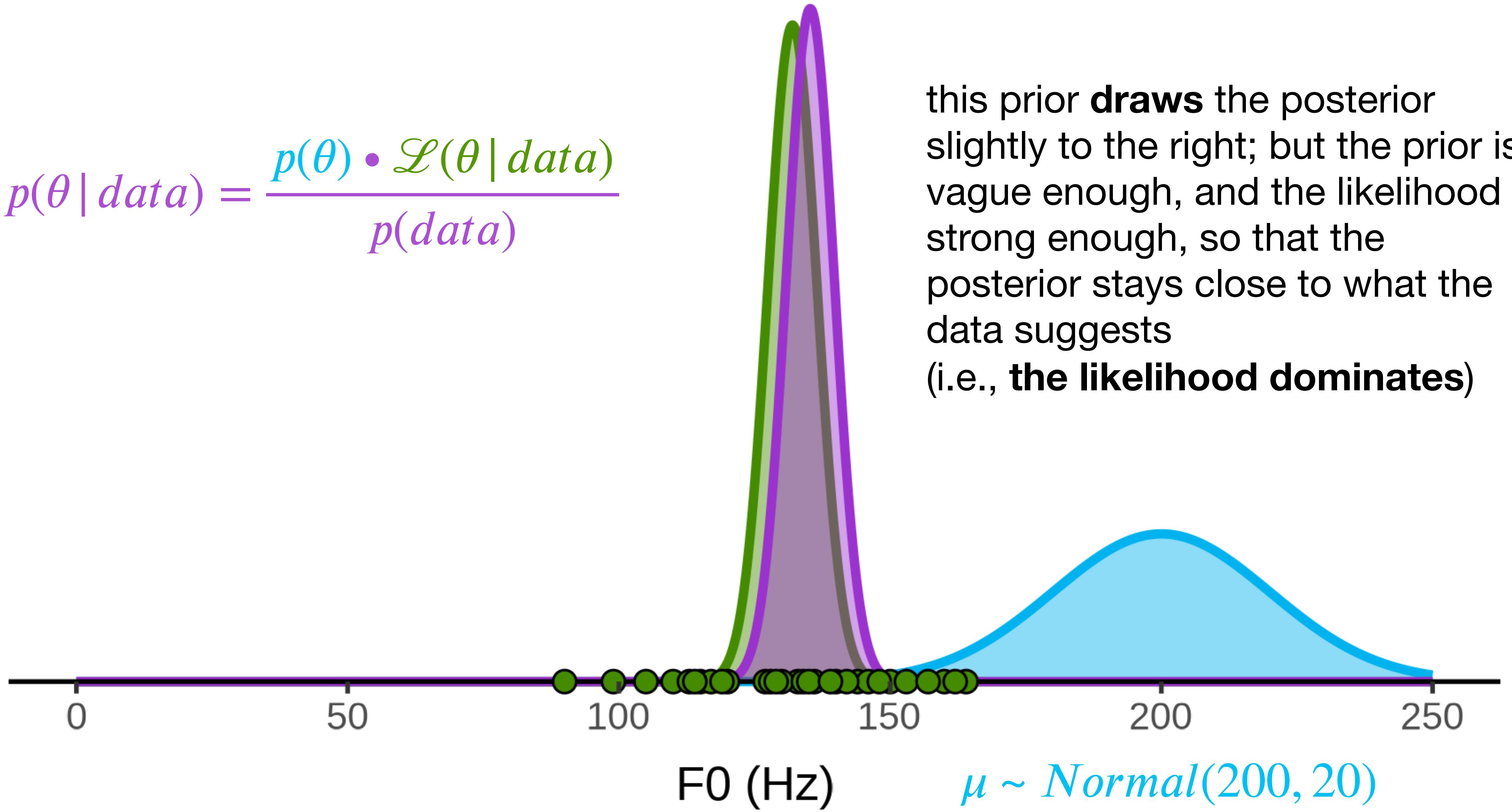
(i.e., everything you base your conclusions on is entirely driven by the evidence)

$$p(\theta | data) = \frac{p(\theta) \cdot \mathcal{L}(\theta | data)}{p(data)}$$

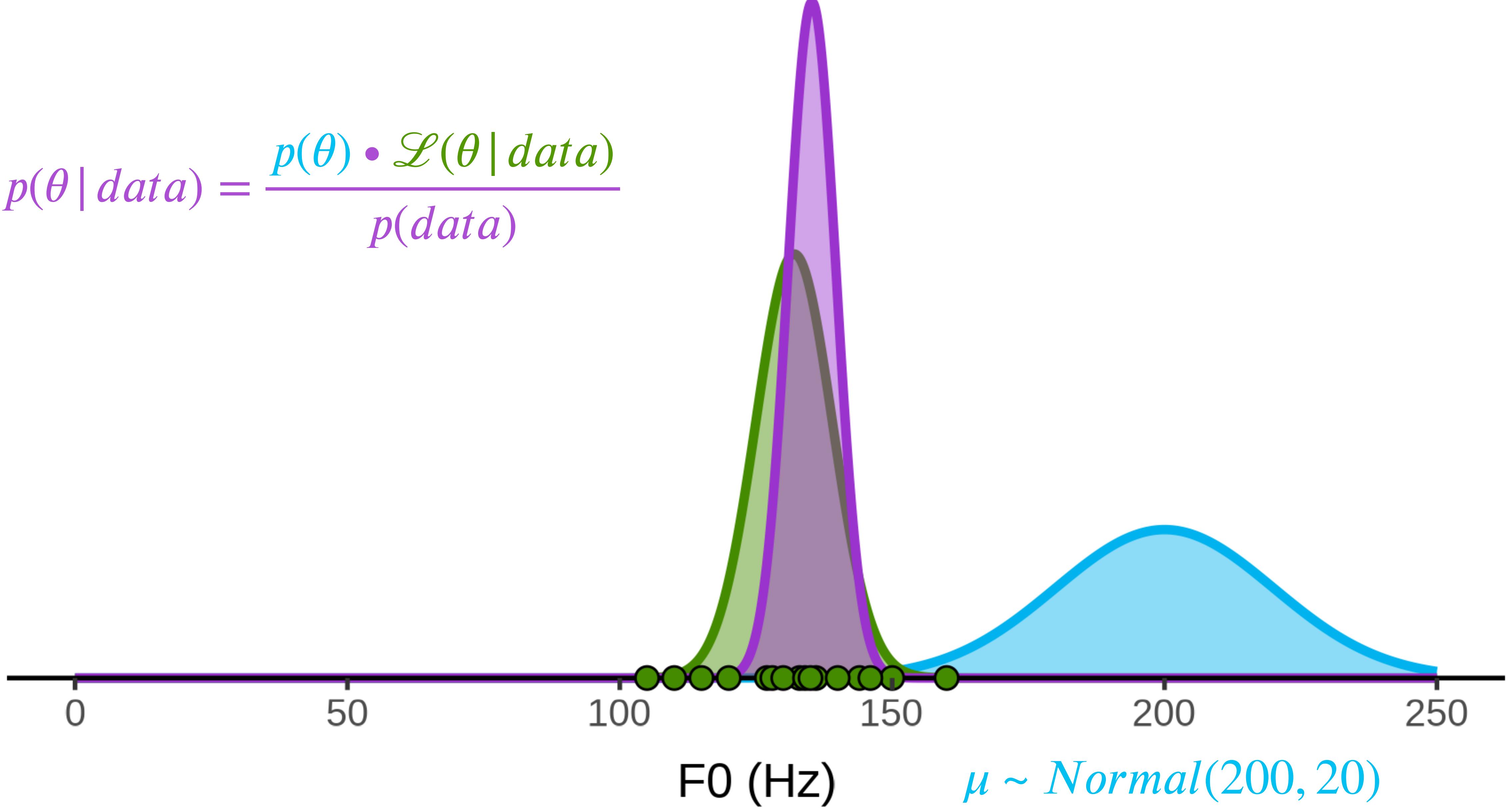


**pedagogical example:**  
somebody didn't get the memo  
that these F0 values come from  
male speakers, and set a  
*Normal(200, 20)* prior

$$p(\theta | \text{data}) = \frac{p(\theta) \cdot \mathcal{L}(\theta | \text{data})}{p(\text{data})}$$

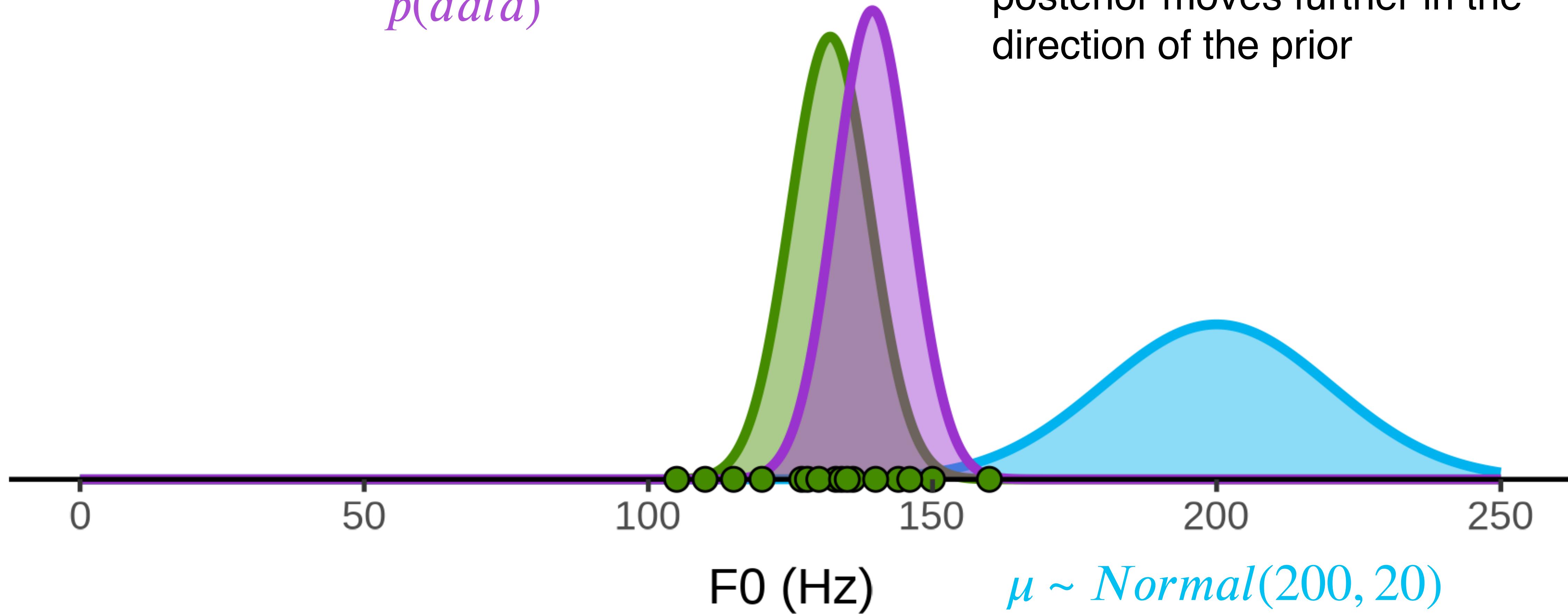


$$p(\theta | data) = \frac{p(\theta) \cdot \mathcal{L}(\theta | data)}{p(data)}$$

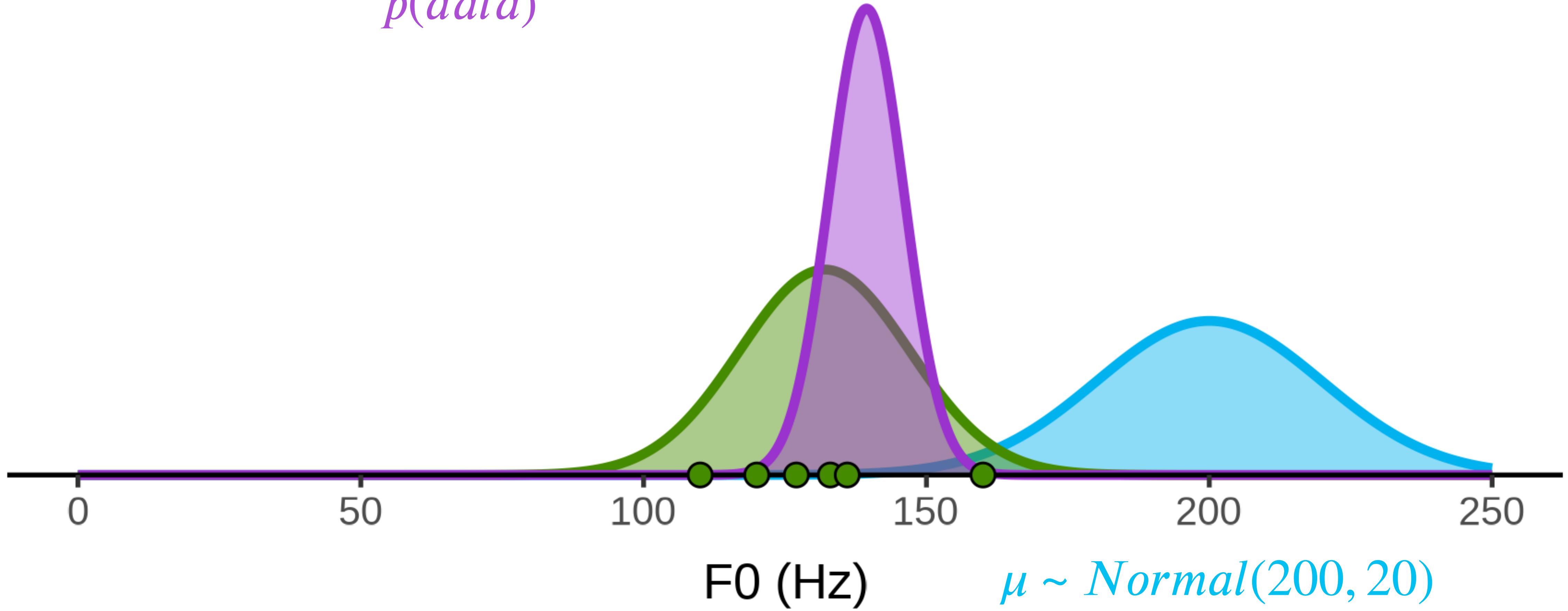


$$p(\theta | data) = \frac{p(\theta) \cdot \mathcal{L}(\theta | data)}{p(data)}$$

with **less data**, the likelihood becomes weaker, and the posterior moves further in the direction of the prior

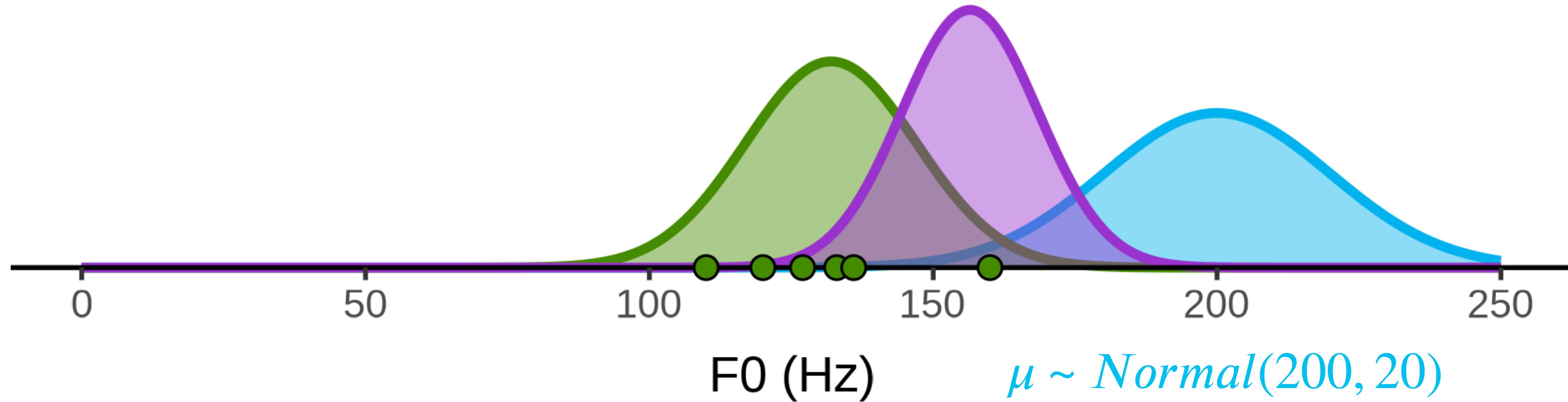


$$p(\theta | data) = \frac{p(\theta) \cdot \mathcal{L}(\theta | data)}{p(data)}$$



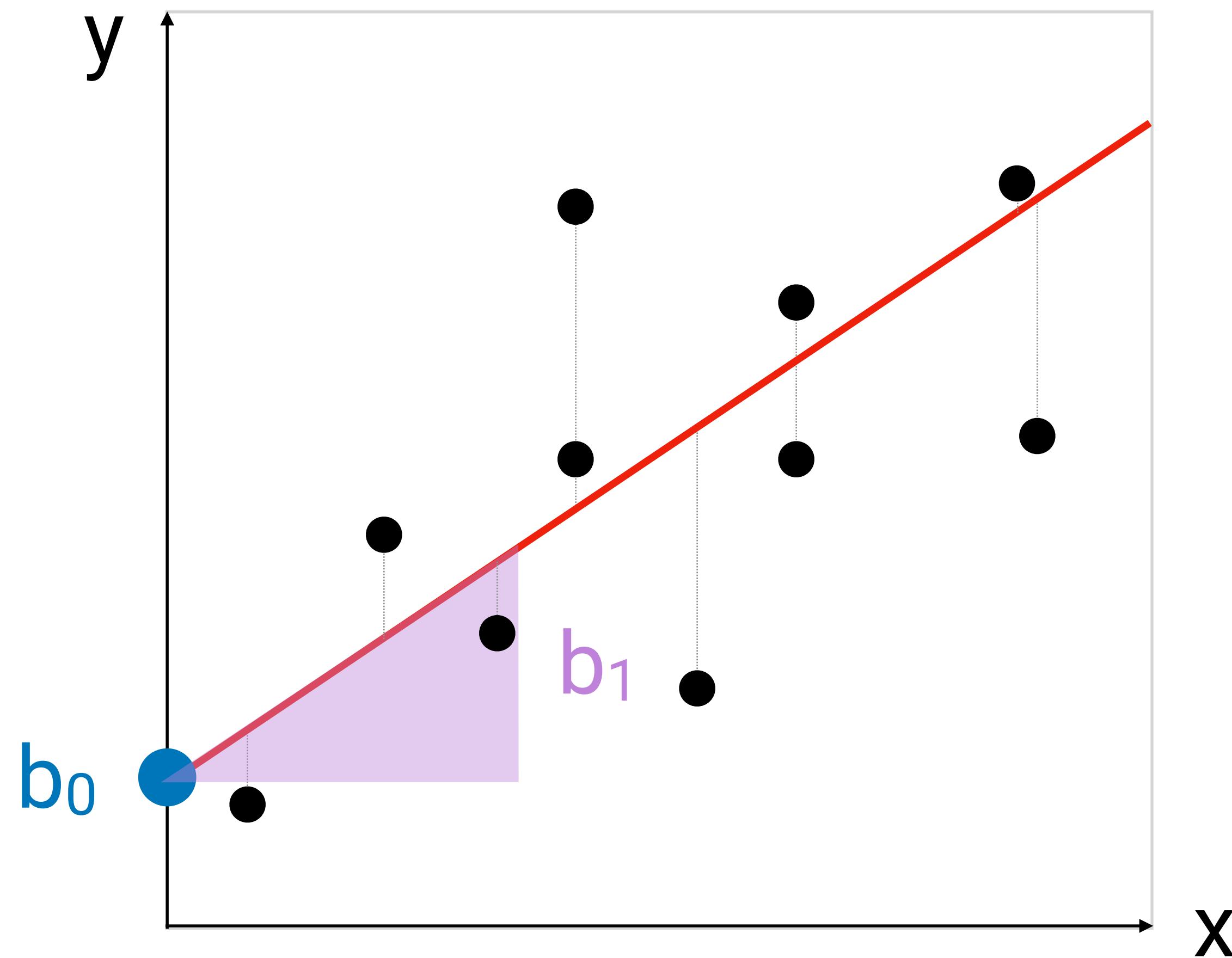
$$p(\theta | data) = \frac{p(\theta) \cdot \mathcal{L}(\theta | data)}{p(data)}$$

with even weaker likelihood,  
the prior exerts an even  
stronger pull on the posterior



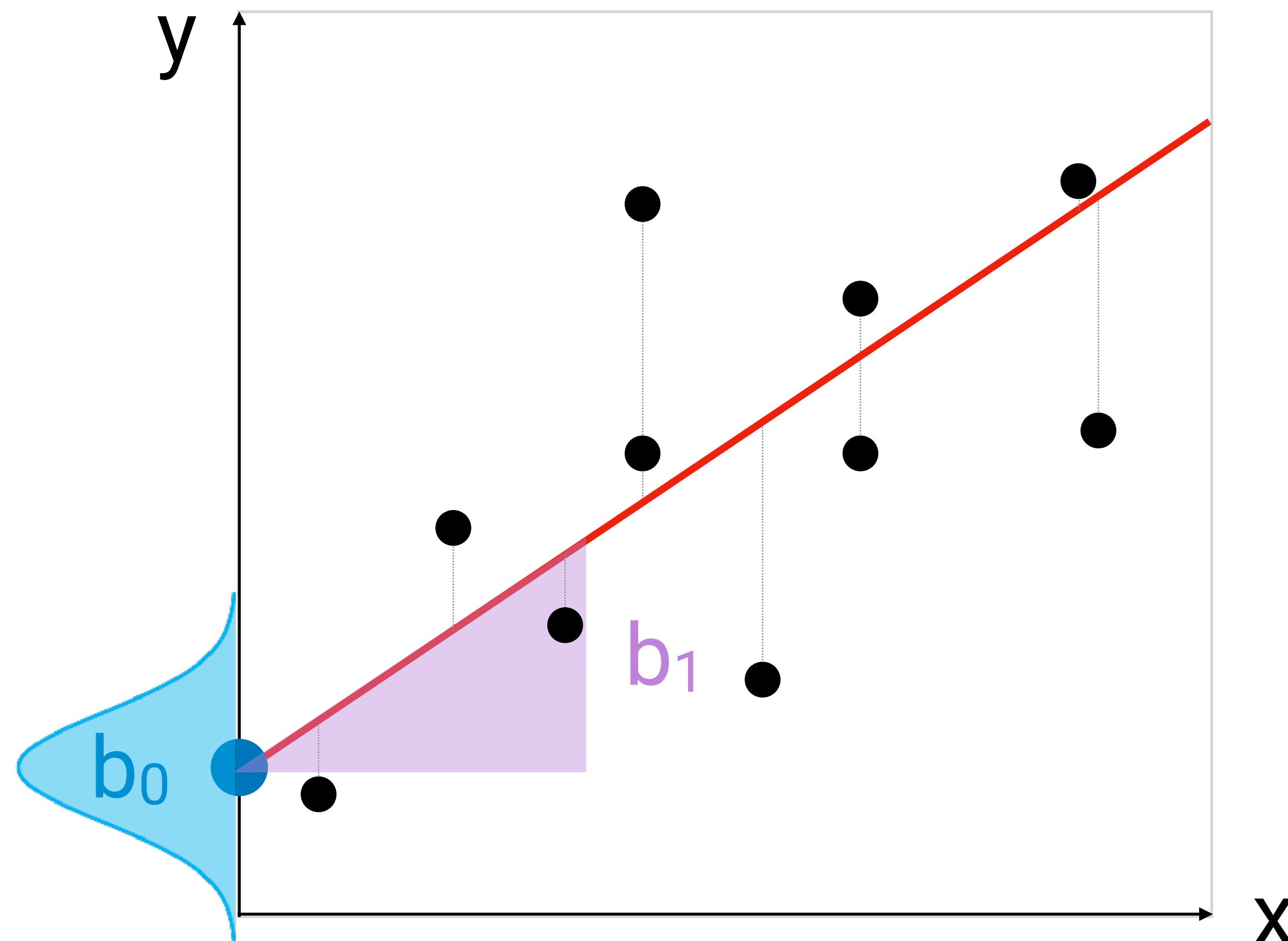
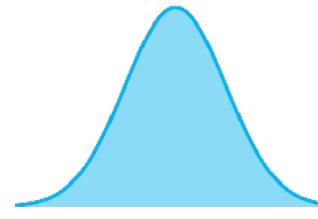
# Priors in linear models

$$y = b_0 + b_1 * x + \sigma$$



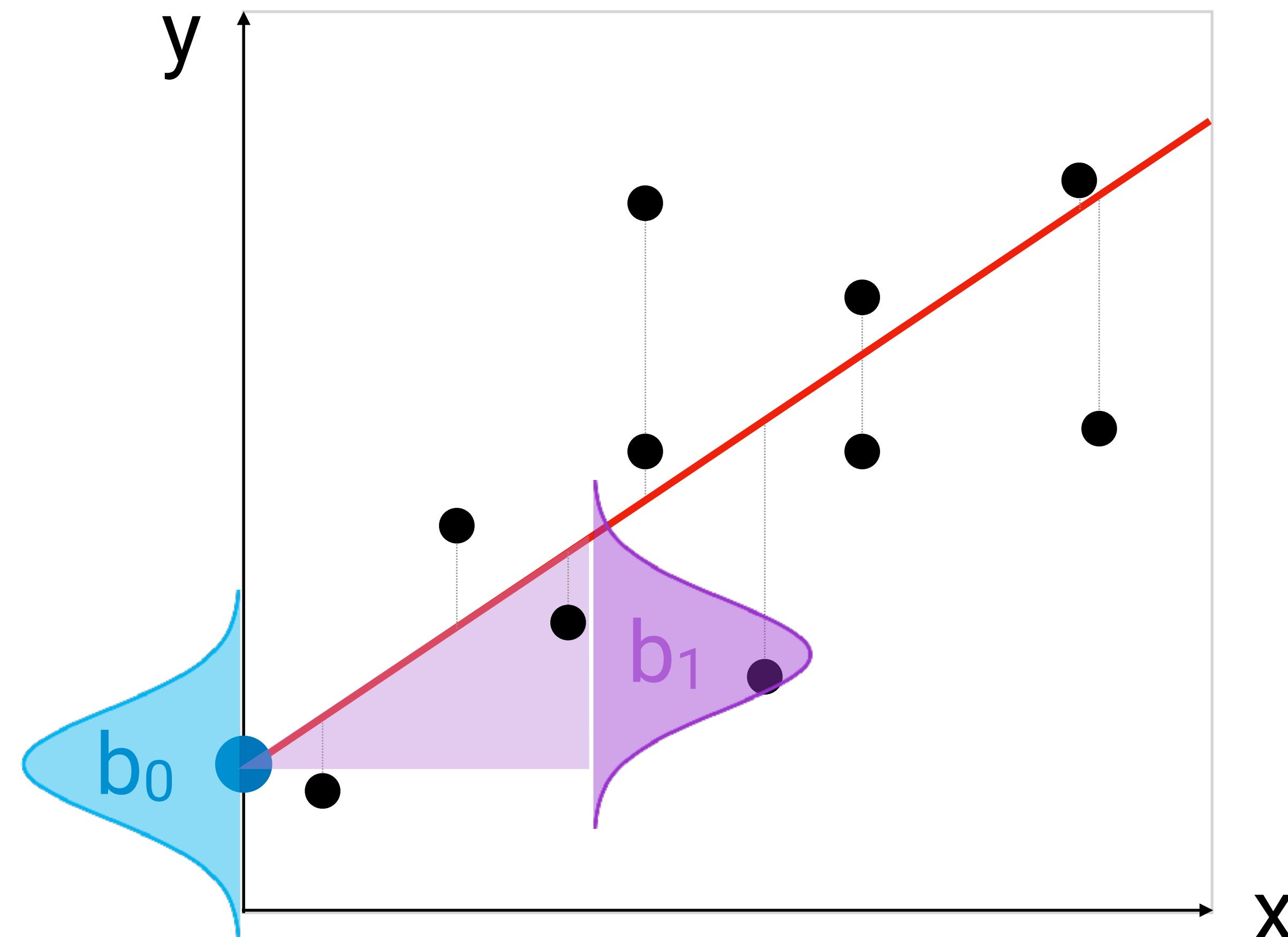
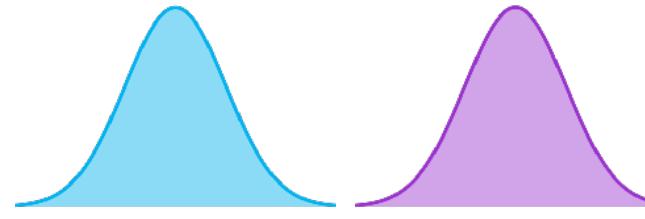
# Priors in linear models

$$y = b_0 + b_1 * x + \sigma$$



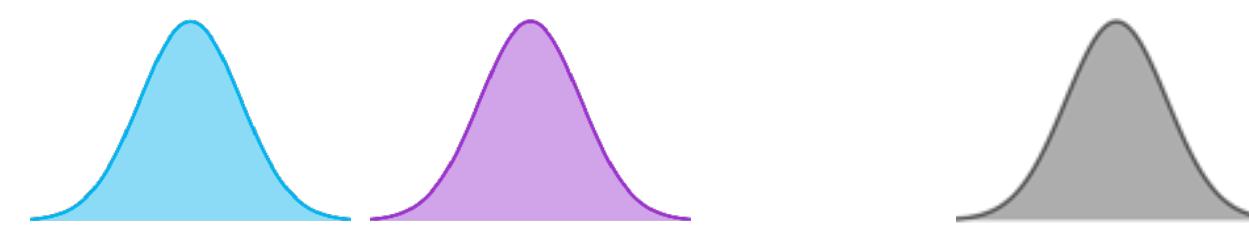
# Priors in linear models

$$y = b_0 + b_1 * x + \sigma$$

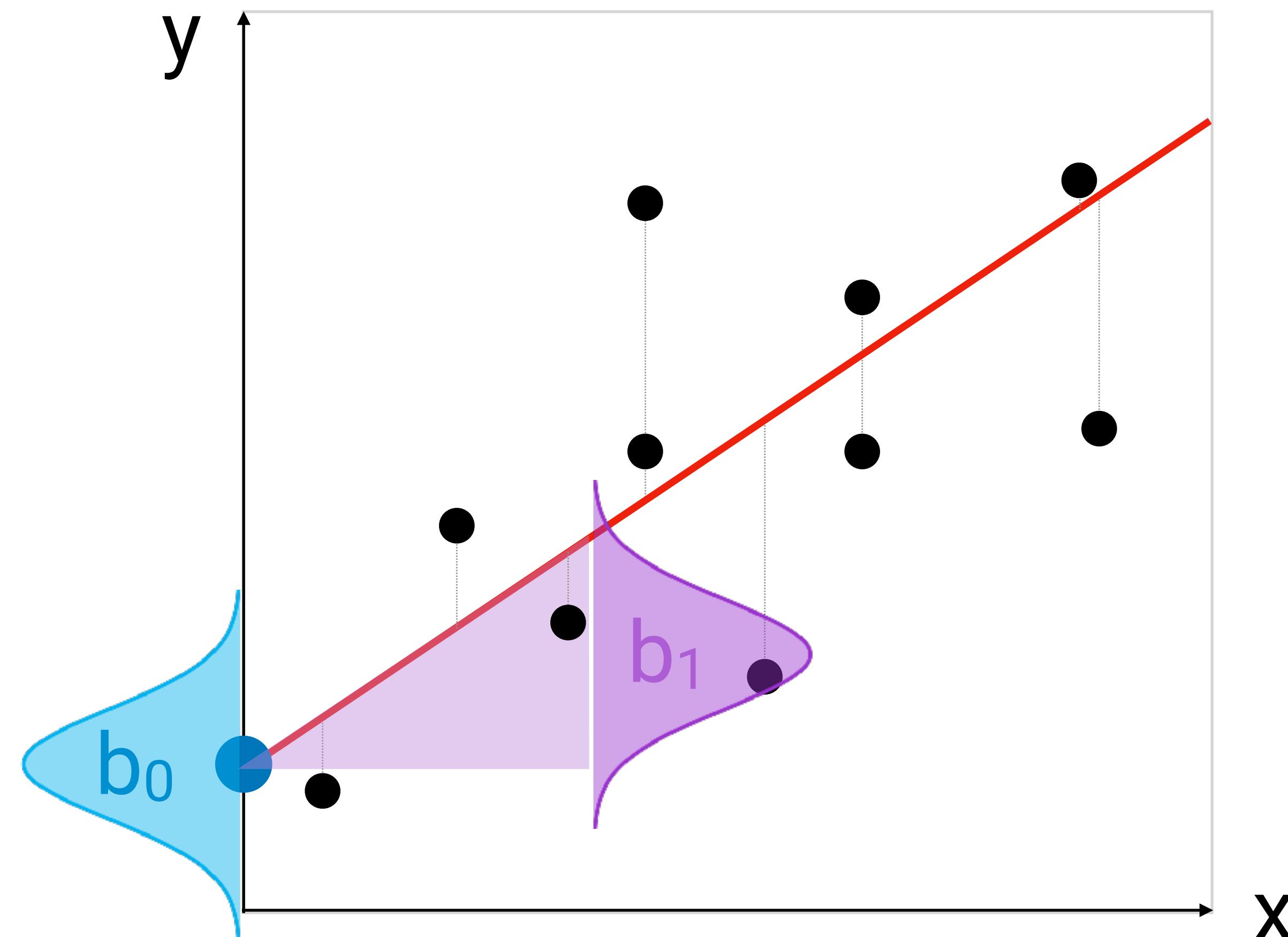


# Priors in linear models

$$y = b_0 + b_1 * x + \sigma$$

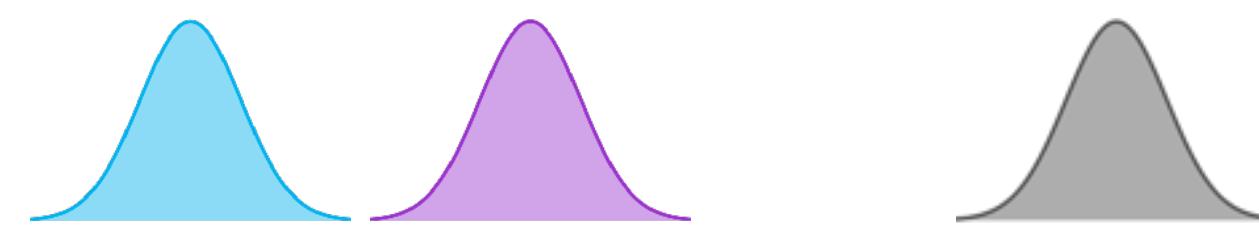


every parameter is modelled  
as a probability distribution  
and receives a prior



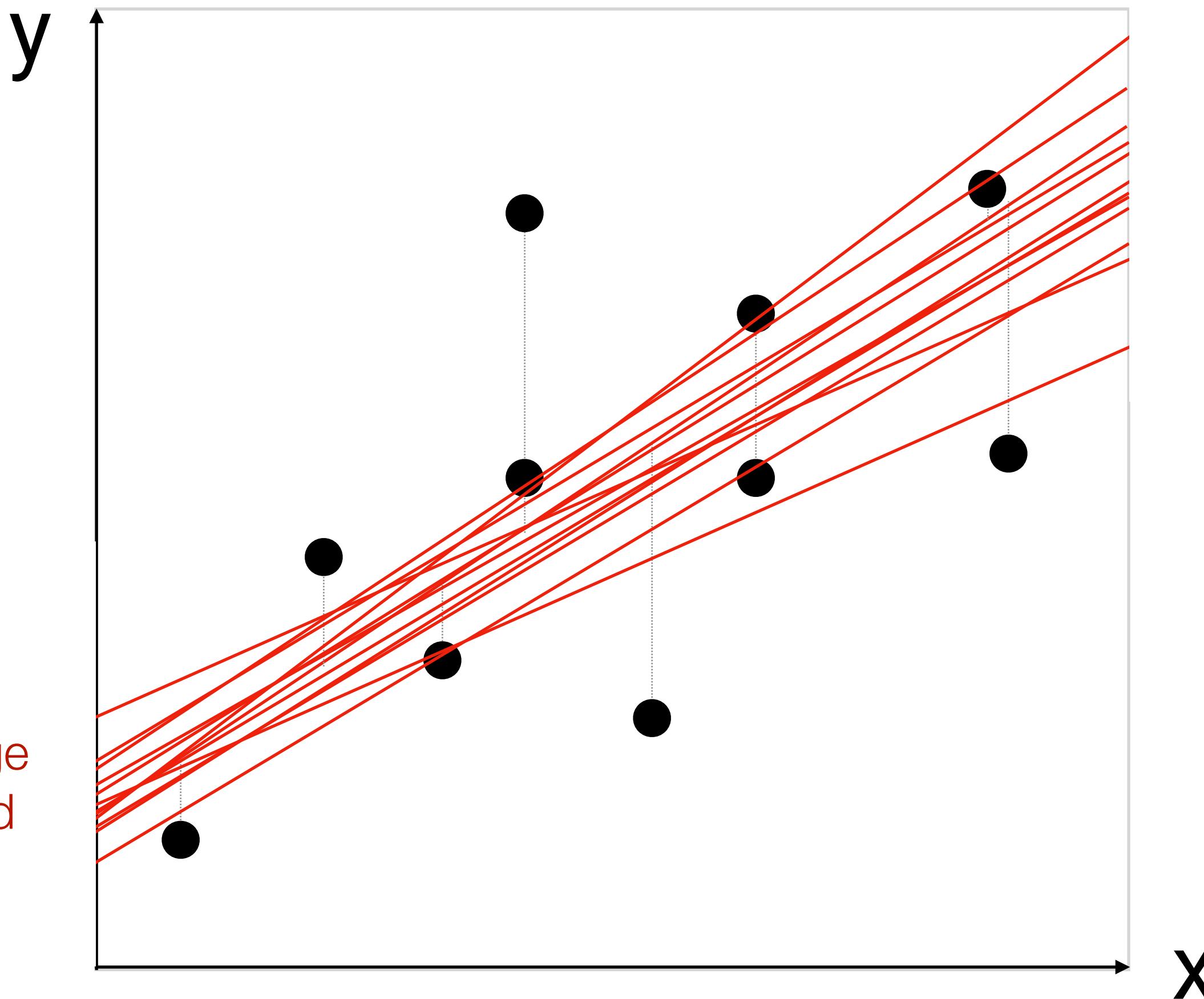
# Priors in linear models

$$y = b_0 + b_1 * x + \sigma$$



every parameter is modelled  
as a probability distribution  
and receives a prior

the model predicts a range  
of plausible intercepts and  
slopes, i.e. a range of  
plausible linear models





# NHST vs Bayes



# Null hypothesis significance testing

- ✖ often does **not** allow **appropriate** use,



Warning message:

```
In checkConv(attr(opt, "derivs"), opt$par, ctrl = control$checkConv, :  
  Model failed to converge with max|grad| = 0.0139723 (tol = 0.002, component 1)
```

```
boundary (singular) fit: see help('isSingular')
```



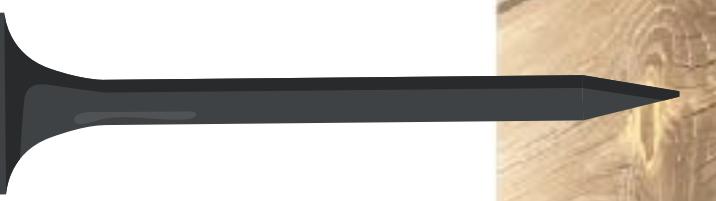
# Null hypothesis significance testing

- ✖ often does **not** allow **appropriate** use,
- ✖ is **not intuitive**,

**89%**

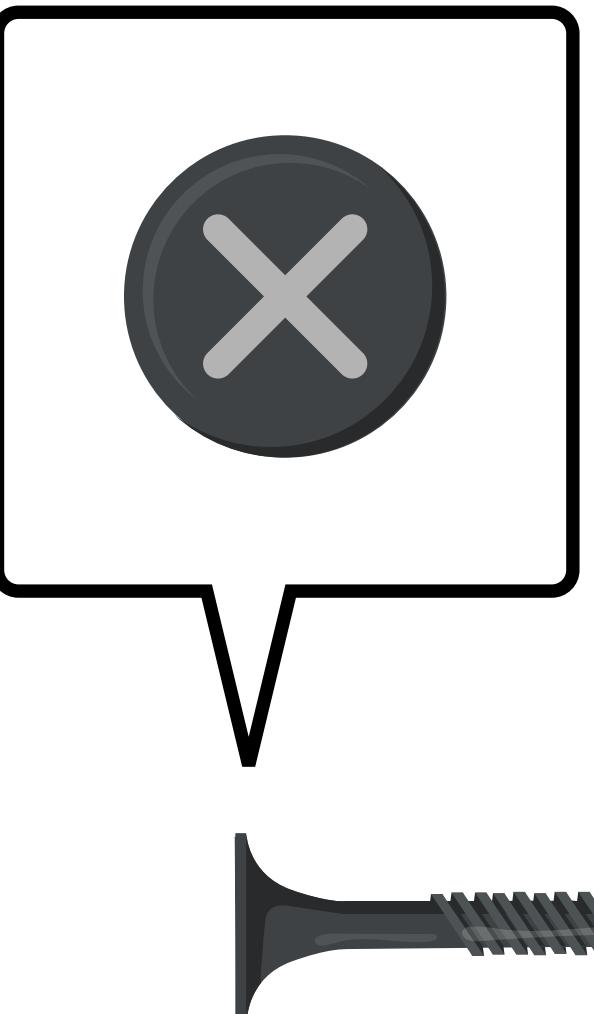
of books that covered statistical significance defined or explained it incorrectly

Cassidy et al. (2019)



# Null hypothesis significance testing

- ✖ often does **not** allow **appropriate** use,
- ✖ is **not intuitive**,
- ✖ and **cannot** provide an **answer** to the **questions** we are interested in.





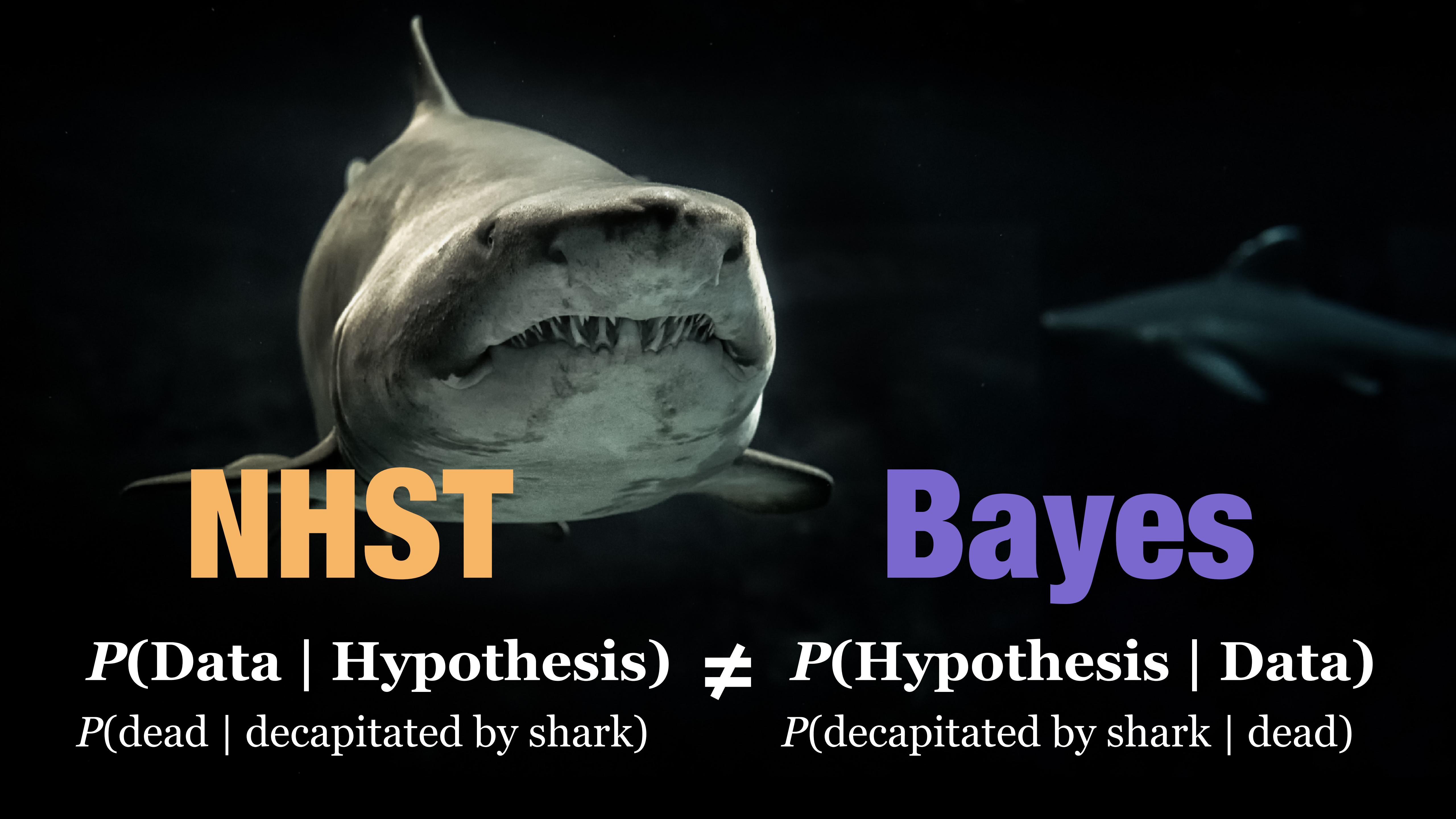
# Bayesian Inference

- ➊ robust inference
- ➋ intuitive
- ➌ flexible

Bürkner, P. C. (2017). brms: An R package for Bayesian multilevel models using Stan. *Journal of statistical software*, 80(1), 1-28.

brms





# NHST

$P(\text{Data} \mid \text{Hypothesis}) \neq P(\text{Hypothesis} \mid \text{Data})$

$P(\text{dead} \mid \text{decapitated by shark})$

# Bayes

$P(\text{decapitated by shark} \mid \text{dead})$

# NHST

$\neq$

# Bayesian

$P(\text{Data} \mid \text{Theory})$

no prior knowledge

quantifies long-run probability  
of finding a false positive

hard cut-off decisions

yields an estimate of one true  
parameter value

$P(\text{Theory} \mid \text{Data})$

incorporates prior knowledge

quantifies uncertainty around  
possible parameter values

gradual assessment of evidence

yields a distribution of plausible  
parameter values

# NHST

$\neq$

# Bayesian

$P(\text{Data} \mid \text{Theory})$

no prior knowledge

quantifies long-run probability  
of finding a false positive

hard cut-off decisions

yields an estimate of one true  
parameter value

$P(\text{Theory} \mid \text{Data})$

incorporates prior knowledge

quantifies uncertainty around  
possible parameter values

gradual assessment of evidence

yields a distribution of plausible  
parameter values

# Bayesian

- 👍 very **flexible** in terms of model architecture
- 👍 not limited by optimization constraints (no “**convergence failures**”)
- 👍 not limited to categorical decision procedure
- 👎 computationally expensive
- 👎 one more layer of researcher degrees of freedom
- 👎 **more “thinking”** required

# flexibility

one and the same framework for everything you need

## types of **error distributions**

gaussian, binomial,  
ordinal, multinomial,  
etc.

## levels of **covariance**

simple regression, multiple  
regression, mixed-effect  
regression, etc.

## types of **fitting procedures**

univariate, multivariate,  
mixture, etc.



# convergence issues?

## lmer()

Warning message:

```
In checkConv(attr(opt, "derivs"), opt$par, ctrl = control$checkConv, :  
  Model failed to converge with max|grad| = 0.0139723 (tol = 0.002, component 1)
```

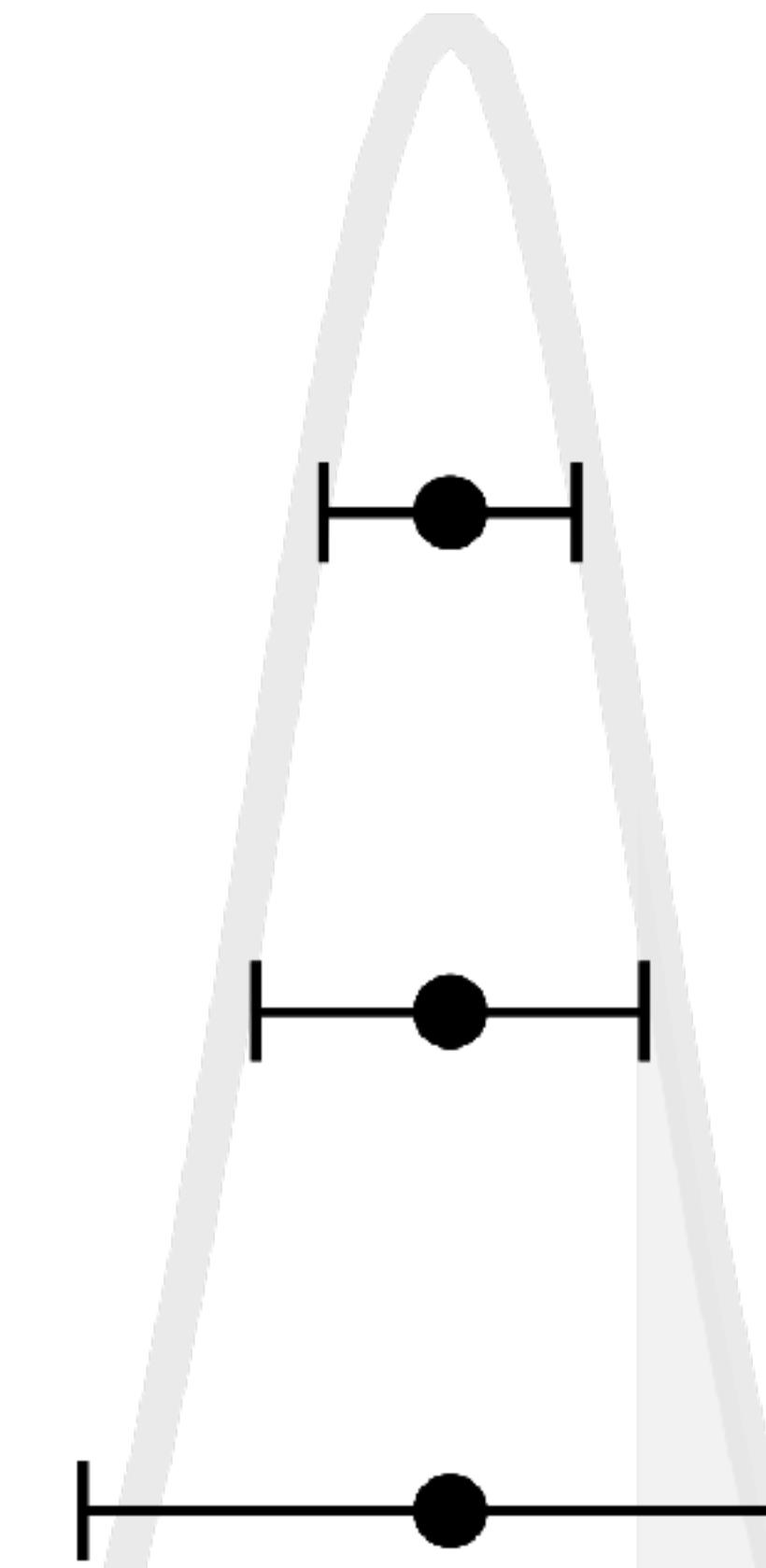
# DO I BRING MY UMBRELLA?



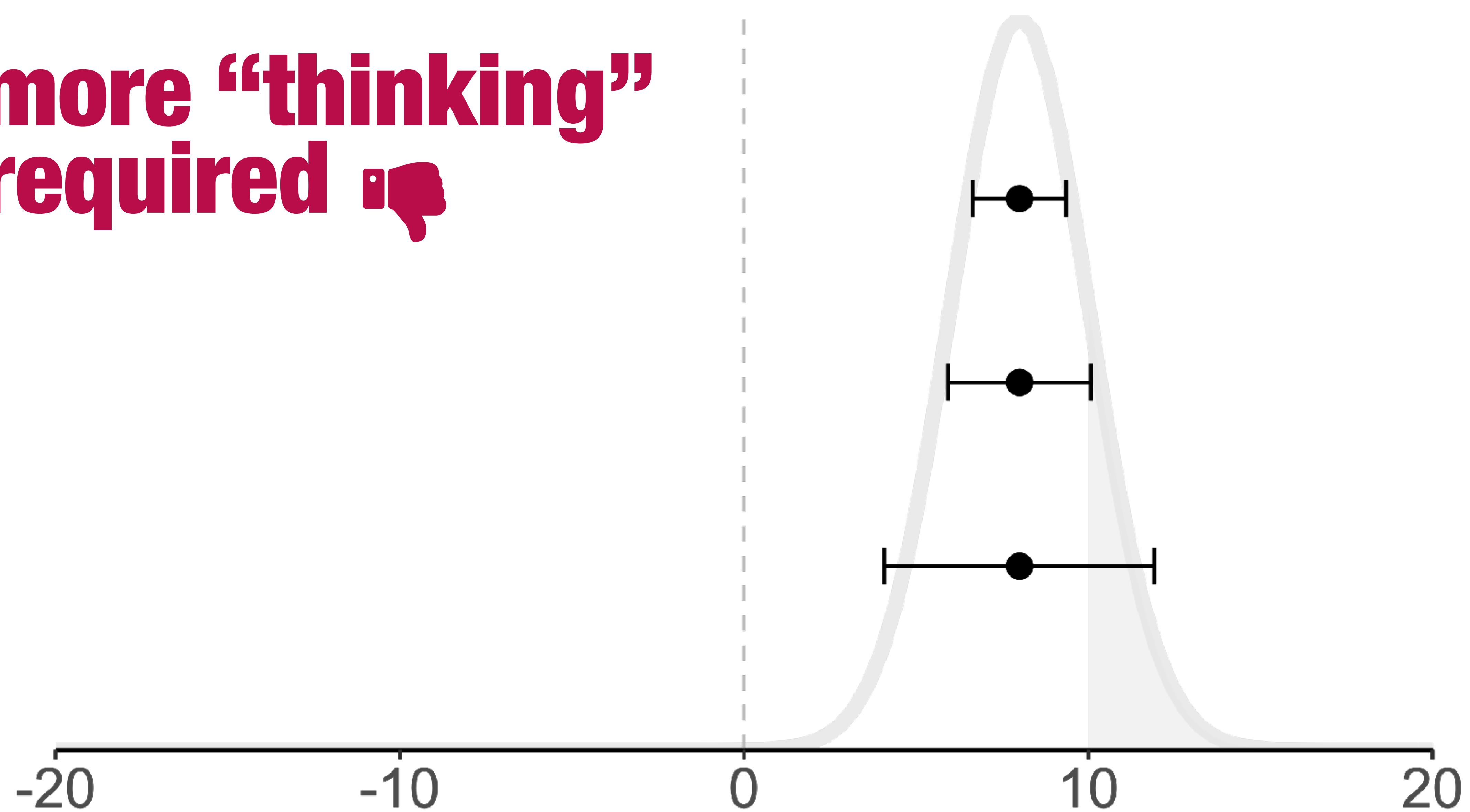
# ARE MY RESULTS SIGNIFICANT?

A dark, moody photograph of two children standing in a field at sunset. The child on the left stands with their back to the viewer, while the child on the right holds a black umbrella over both of them. The sky is filled with heavy, dark clouds, with some lighter areas showing the setting sun.

**no decision  
procedure  
required** 



**more “thinking”  
required** 🙅





# computationally expensive

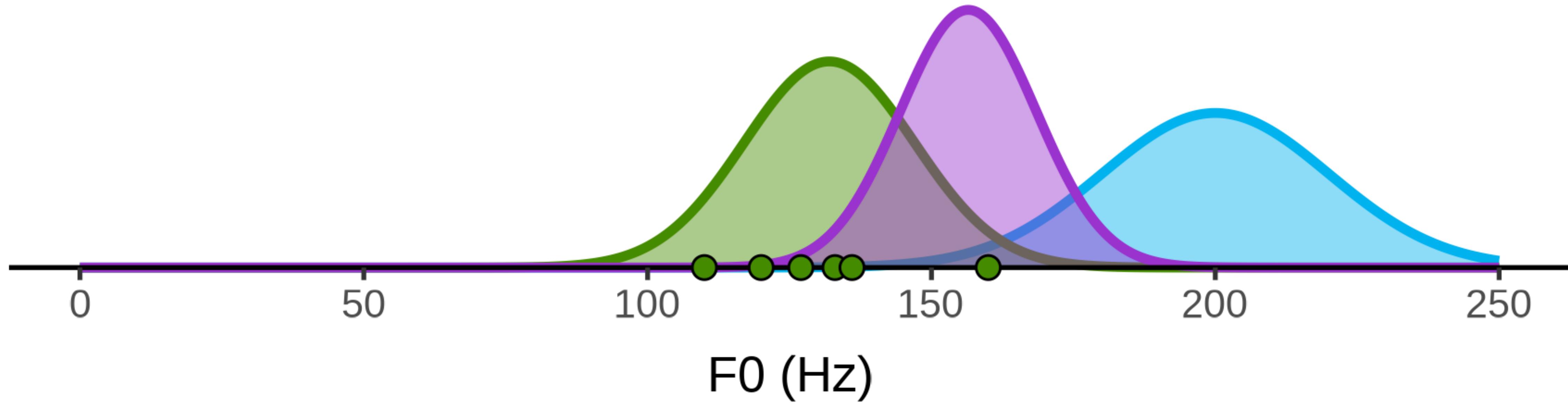
$$Pr(\text{Theory} \mid \text{Data}) = \frac{Pr(\text{Data} \mid \text{Theory}) \times Pr(\text{Theory})}{Pr(\text{Data})}$$

$$\Pr(\text{Data}) = \int \Pr(\text{Data}, \text{Theory}) d\text{Theory}$$

can be **intractable** to solve, but...

can be **approximated** with clever algorithms

**more degrees of  
freedom? **



# Bayesian

- 👍 very **flexible** in terms of model architecture
- 👍 not limited by optimization constraints (no “**convergence failures**”)
- 👍 not limited to categorical decision procedure
- 👎 computationally expensive
- 👎 one more layer of researcher degrees of freedom
- 👎 **more “thinking”** required



**Just one more thing...**

# The technical implementation is easy

```
xlm <- lm(F0 ~ 1, data = pitch)  
summary(xlm)
```

```
xblm <- brm(F0 ~ 1, data = pitch)  
summary(xblm)
```

## Regression Coefficients:

	Estimate	Est.Error	l-95%	CI	u-95%	CI	Rhat	Bulk_ESS	Tail_ESS
Intercept	124.12	4.15	115.73	132.21	1.00	2713	2508		

## Further Distributional Parameters:

	Estimate	Est.Error	l-95%	CI	u-95%	CI	Rhat	Bulk_ESS	Tail_ESS
sigma	23.32	2.99	18.37	29.96	1.00	2722	2271		