

Scientific Machine Learning

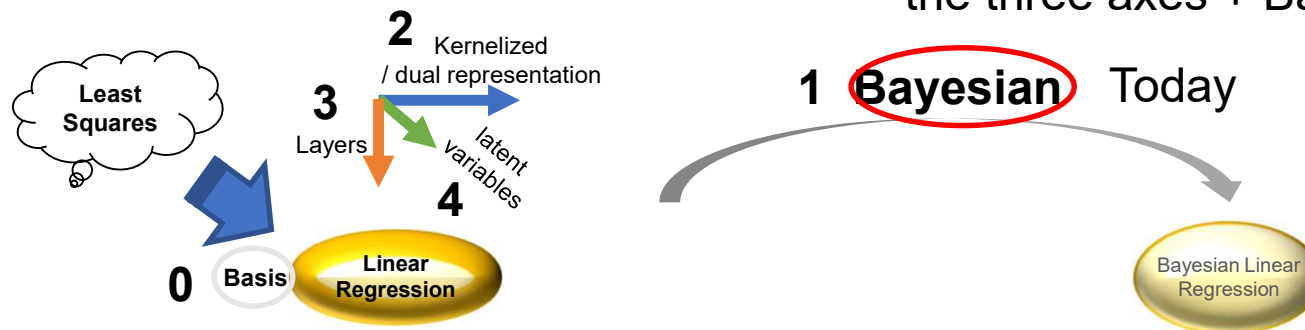
Lecture 5: Bayesian Statistics (1/2)

Dr. Daigo Maruyama

Prof. Dr. Ali Elham

Key Components

the three axes + Bayesian axis = four axes



Lecture content

- Bayesian approach
- Summary (update of the current process)

There are no references about the context of this lecture.



Process in General (Review)

1. Define a probabilistic model

$$p(y|x, \mathbf{w})$$

← Model definition

2. MLE

$$p(\mathbf{y}|\mathbf{X}, \mathbf{w}) \xrightarrow{\text{MLE}} \hat{\mathbf{w}}$$

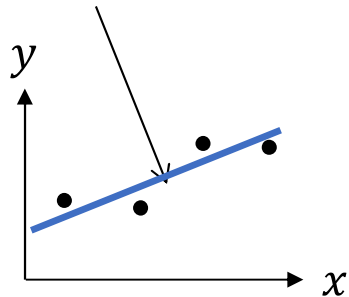
← Learning process

3. New prediction

$$p(y_{\text{new}}|x_{\text{new}}, \hat{\mathbf{w}})$$

← Prediction process

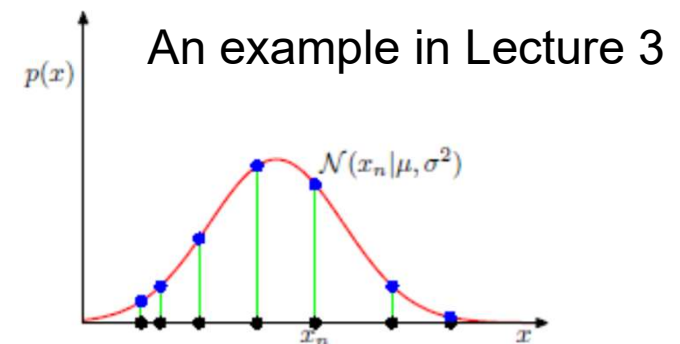
$$y = w_1 x + w_0 = \mathbf{w}^T \boldsymbol{\phi}(x)$$



Data: \mathcal{D}

$$\mathbf{X} = (\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(N)})^T$$

$$\mathbf{y} = (y^{(1)}, y^{(2)}, \dots, y^{(N)})^T$$



Bayes' theorem (REVIEW of Lecture 2)

sum rule $P(Y) = \sum_X P(X, Y)$

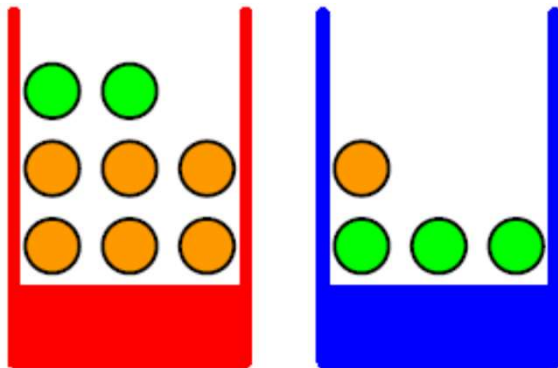
product rule $P(X, Y) = P(X|Y)P(Y)$ $P(X, Y) = P(Y|X)P(X)$

Bayes' theorem

$$P(Y|X) = \frac{P(X|Y)P(Y)}{P(X)}$$

Interpretation is important.

Let us consider
time flow / causality



Bayesian Approach

Relationship between parameters \mathbf{w} and data \mathcal{D}

Bayes' theorem $P(\mathbf{w}|\mathcal{D}) = \frac{P(\mathcal{D}|\mathbf{w})P(\mathbf{w})}{P(\mathcal{D})}$

Likelihood function $P(\mathcal{D}|\mathbf{w})$

So far:

1. Define a probabilistic model
2. Then, MLE (Maximum Likelihood Estimation)

extended

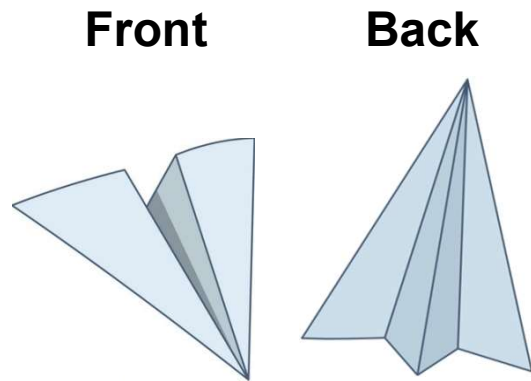
In this process, the goal was to determine the likelihood function, then to maximize it.

Bayesian Approach

Example:

You threw a paper plane three times.

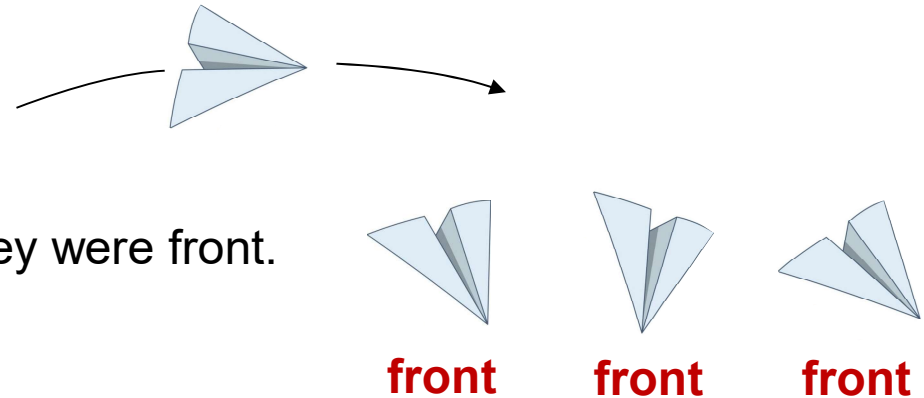
Observed Data: all the three times they were front.



probability
 μ

probability
 $1 - \mu$

$$0 \leq \mu \leq 1$$



What is the probability μ ?

You can imagine **a coin** as well.

→ Then, $\mu = 0.5$?

But why?

Bayesian Approach

Let's do the MLE.

1. Define a probabilistic model
2. Then, MLE

For the multiple trials for the binary result, the probability distribution is



Binomial distribution

$$p(m|\mu) = \binom{3}{m} \mu^m (1 - \mu)^{3-m}$$

$$0 \leq \mu \leq 1$$

Probabilistic model

strictly, $p(m|N = 3, \mu)$ ($N = 3$ is fixed and omitted.)

Now we could obtain data \mathcal{D} :

$$\mathcal{D}: m = 3$$

$$p(m = 3|\mu) = \binom{3}{3} \mu^3 (1 - \mu)^0 = \mu^3$$

Likelihood function



MLE:

maximize $p(m = 3|\mu)$

and find $\hat{\mu}$ where it maximize p

$$\hat{\mu} = \max_{\mu} p(m = 3|\mu)$$

Bayesian Approach

$$\begin{aligned}\hat{\mu} &= \max_{\mu} p(m = 3 | \mu) \\ &= \max_{\mu} \mu^3\end{aligned}\quad 0 \leq \mu \leq 1$$

➡ $\hat{\mu} = 1$ The probability where the front is observed is 1.

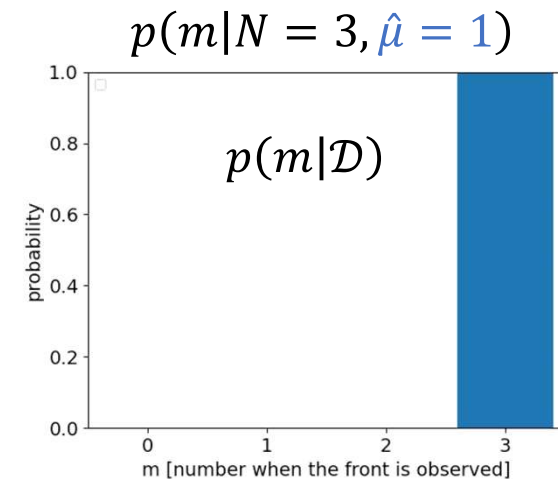
Probabilistic model

$$p(m | \mu) = \binom{3}{m} \mu^m (1 - \mu)^{3-m}$$

What is the probability of m ?

GOAL: Let's use this for prediction (for arbitrary m).

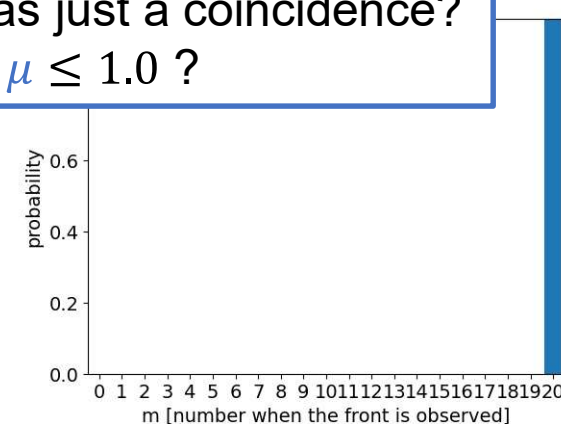
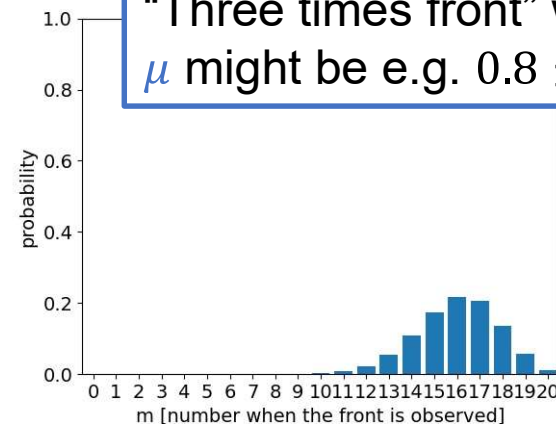
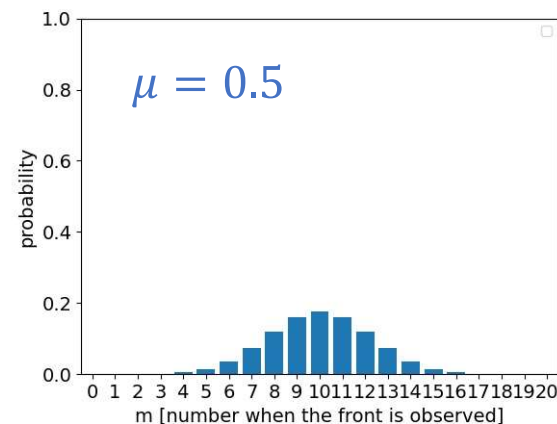
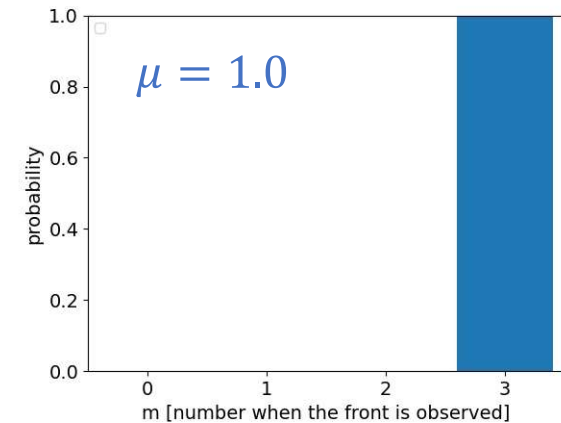
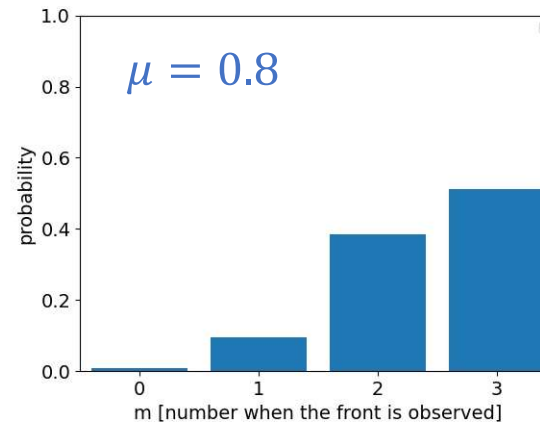
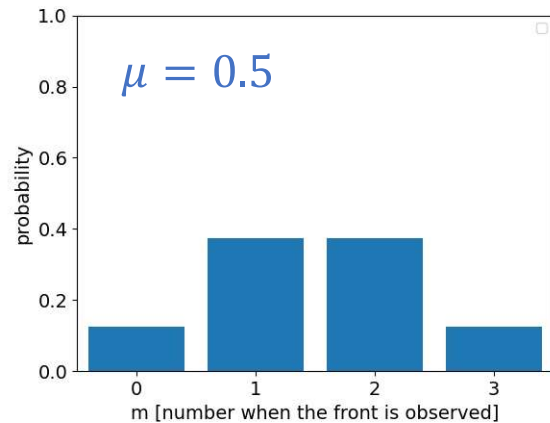
$$\begin{aligned}p(m | \hat{\mu} = 1) &= \binom{3}{m} 1^m (1 - 1)^{3-m} \\ &= \binom{3}{m}\end{aligned}$$



Bayesian Approach

$$0 \leq \mu \leq 1$$

$$N = 3$$



$\mu = 0.8$ seems to be more realistic?
“Three times front” was just a coincidence?
 μ might be e.g. $0.8 \leq \mu \leq 1.0$?

Bayesian Approach

What was the problem?

- Maybe too few trials

➡ If we try $N = 20$, the front might observed less than that at least not 20.

- Difficult to uniquely determine the parameter μ

➡ μ might be between 0.8 and 1.0?

- Contrary to your belief? or your past experiences?

➡ The true μ might be around 0.5 in my belief / past experiences.

Bayesian Approach

STOP before doing MLE!

We can manage to use all the information of the likelihood function.

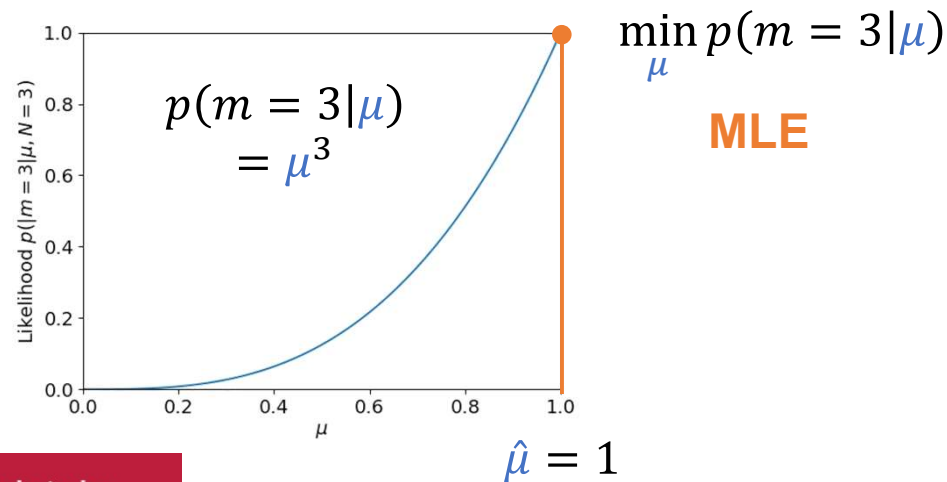
Data \mathcal{D} : $m = 3$

$$p(m = 3|\mu) = \binom{3}{3} \mu^3 (1 - \mu)^0 = \mu^3$$

$$N = 3$$

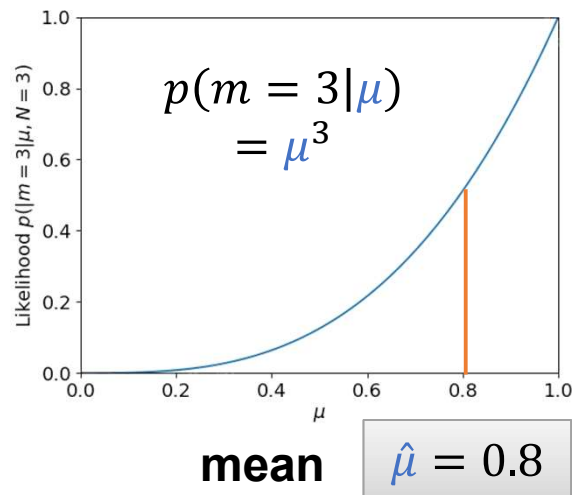
$$0 \leq \mu \leq 1$$

Likelihood function



Bayesian Approach

Likelihood function



How about taking **the expectation**?

We just have started to be conscious of **the distribution**!

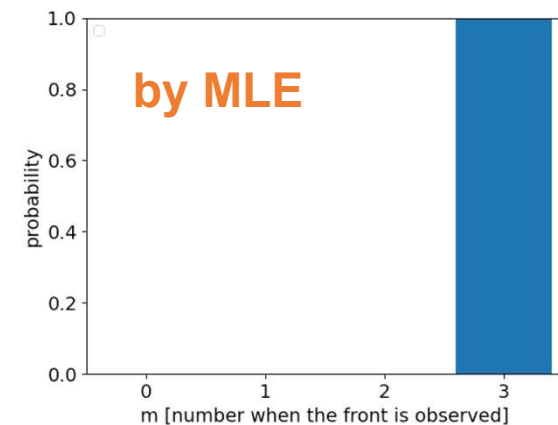
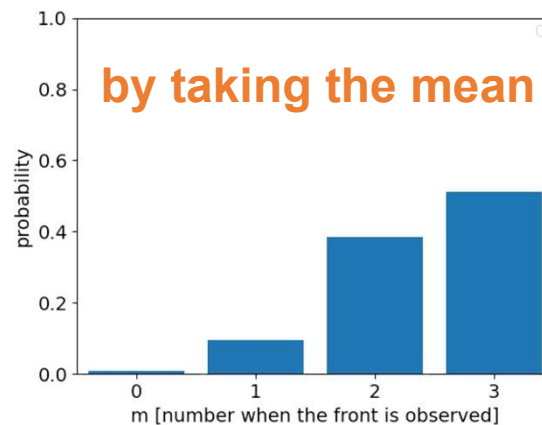


Bayesian perspectives

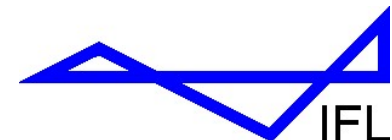
$\mu = 0.8$

$\mu = 1.0$

MLE is not always
the best option?

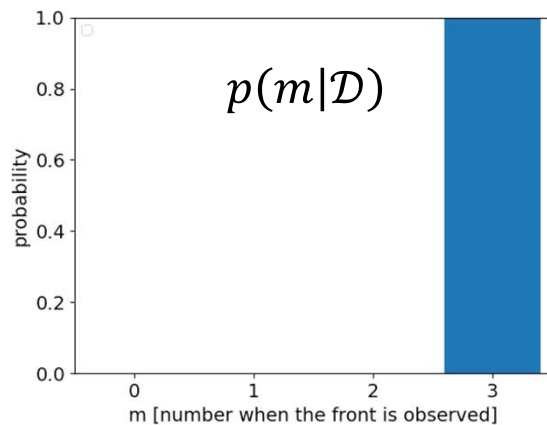
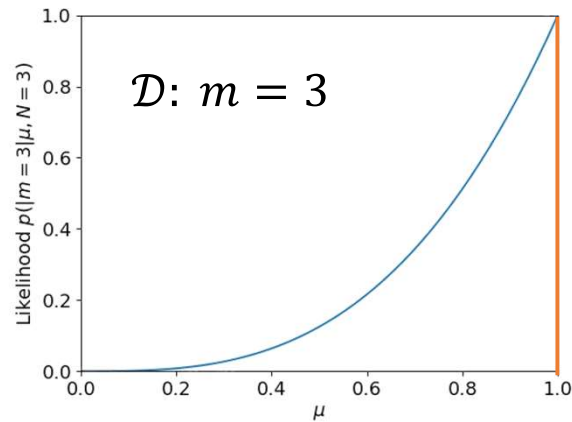


Technische
Universität
Braunschweig

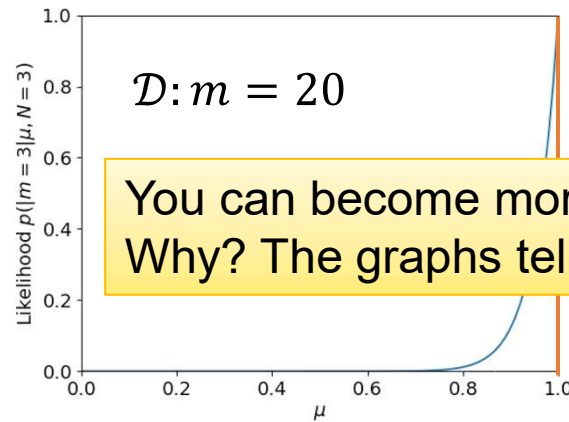


Bayesian Approach

$N = 3$



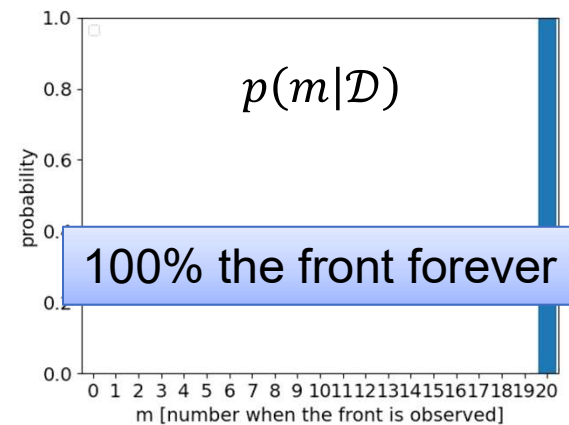
$N = 20$



all $\hat{\mu} = 1$

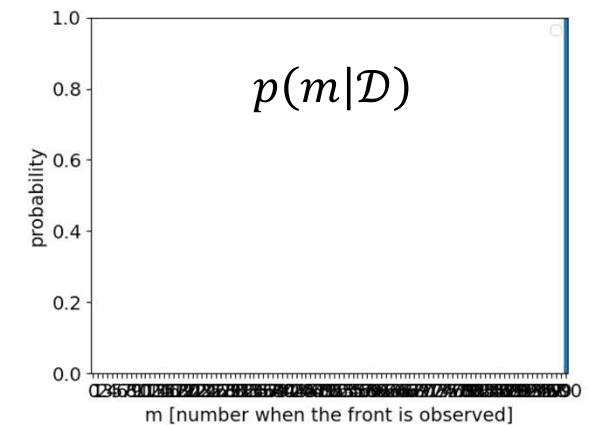
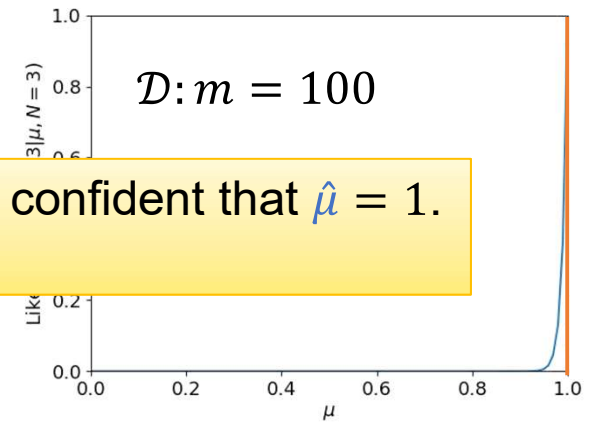


by MLE



100% the front forever

$N = 100$

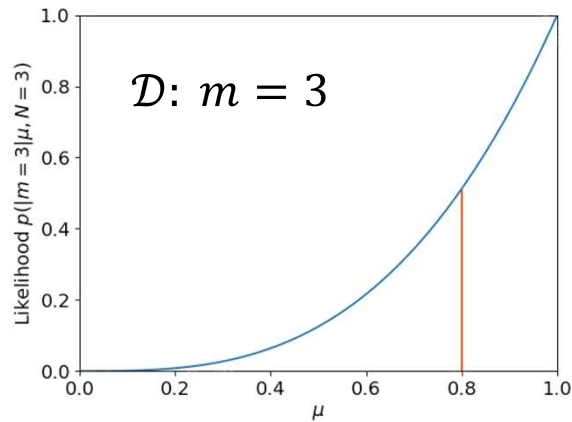


You can become more confident that $\hat{\mu} = 1$. Why? The graphs tell.

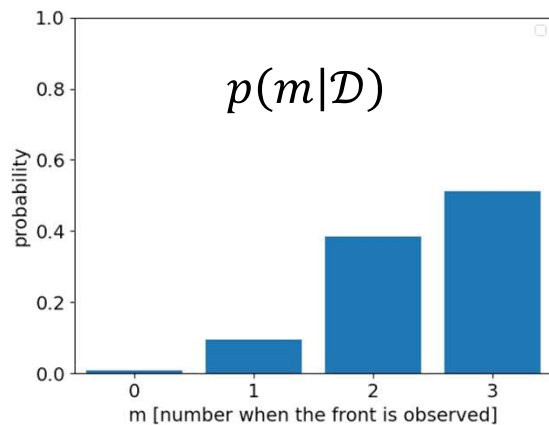
Bayesian Approach

by taking the mean

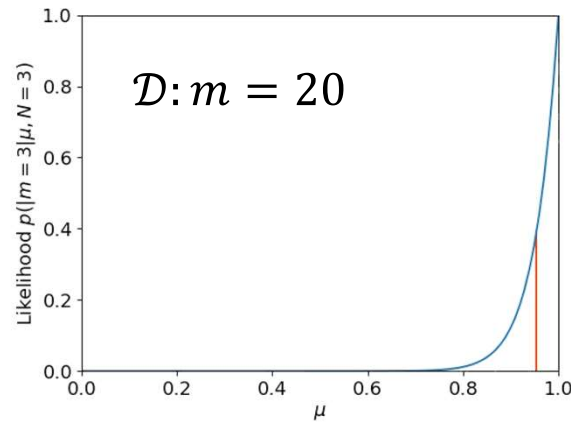
$N = 3$



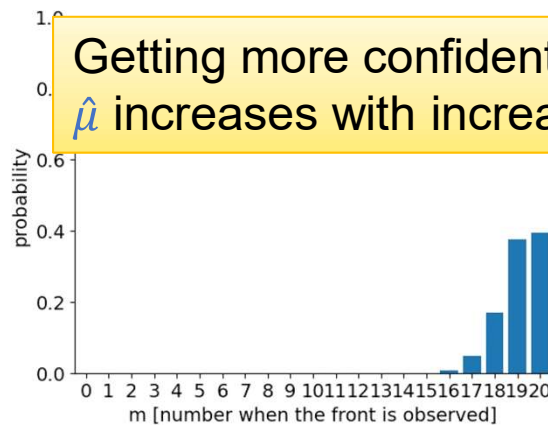
$$\hat{\mu} = 0.8$$



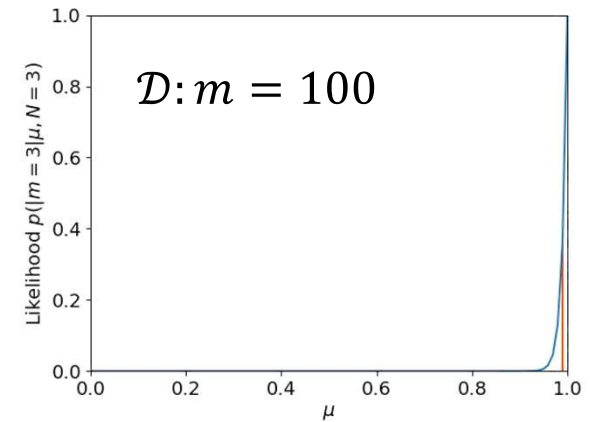
$N = 20$



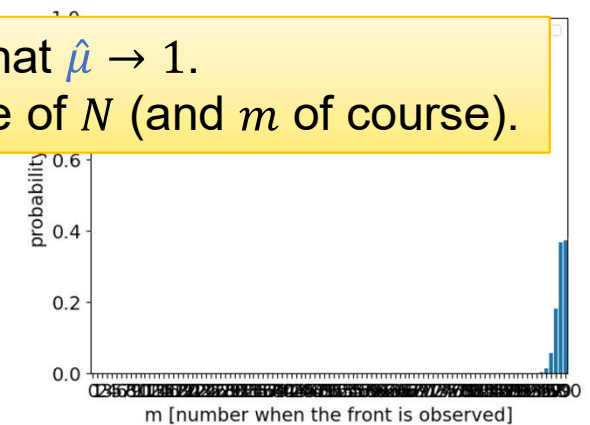
$$\hat{\mu} = 0.954$$



$N = 100$



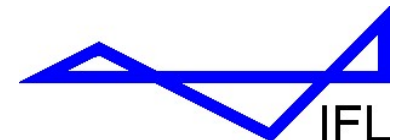
$$\hat{\mu} = 0.990$$



Getting more confident that $\hat{\mu} \rightarrow 1$.
 $\hat{\mu}$ increases with increase of N (and m of course).



Technische
Universität
Braunschweig



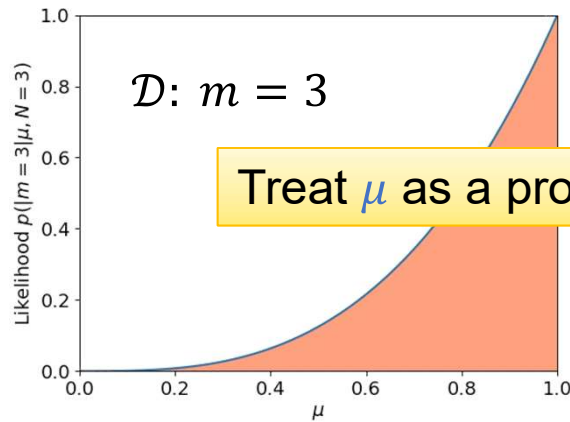
Bayesian Approach

by using all the information

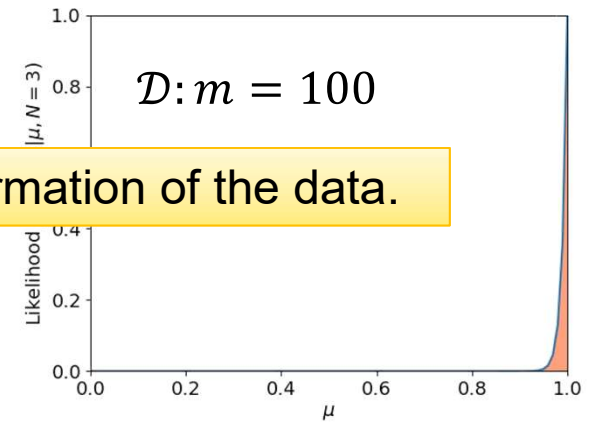
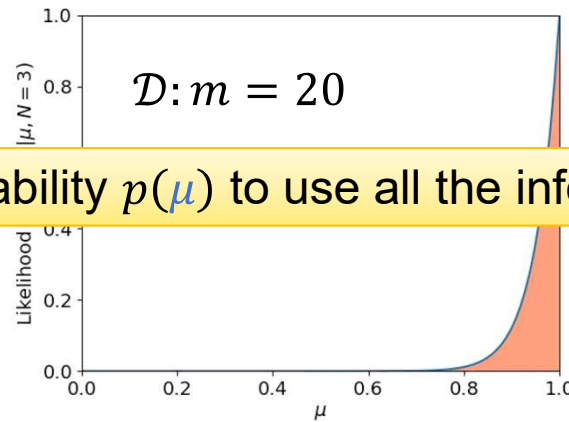
$N = 3$

$N = 20$

$N = 100$



Treat μ as a probability $p(\mu)$ to use all the information of the data.



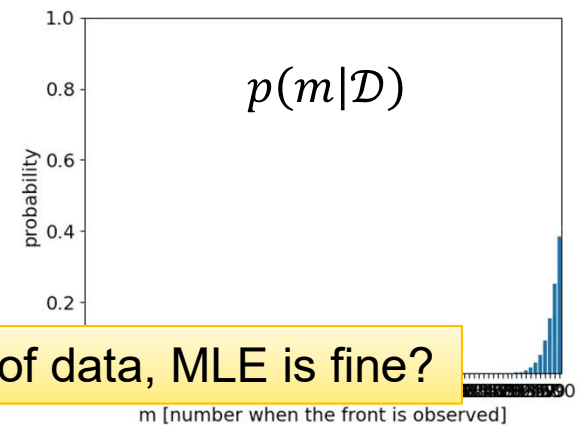
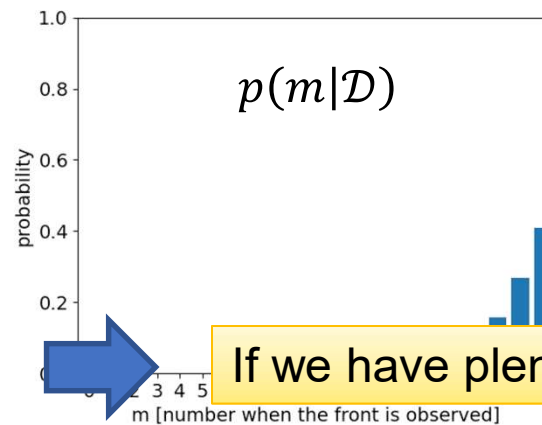
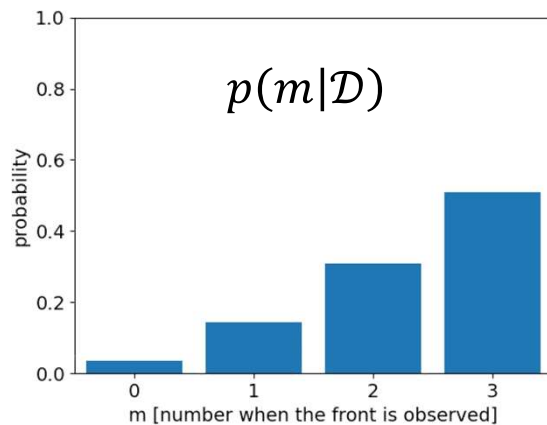
$p(\mu|m=3)$



$p(\mu|m=20)$



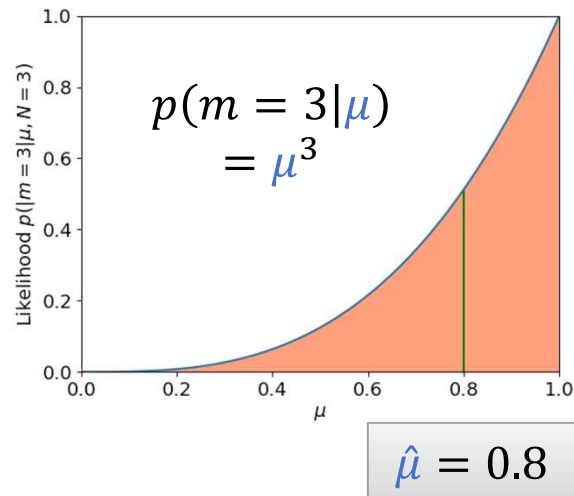
$p(\mu|m=100)$



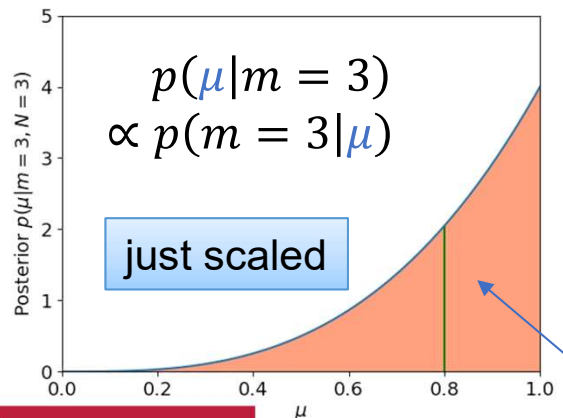
If we have plenty of data, MLE is fine?

Bayesian Approach

Likelihood function



Posterior function



We have to treat $p(\mu|m=3)$, not $p(m=3|\mu)$.

Bayes' theorem

$$p(\mu|m) = \frac{p(m|\mu)p(\mu)}{p(m)} = \frac{\text{likelihood} \times \text{prior}}{p(m)}$$

$$= p(m|\mu) \times \frac{p(\mu)}{p(m)}$$

no info in advance
no relation with μ

constant (scaling)

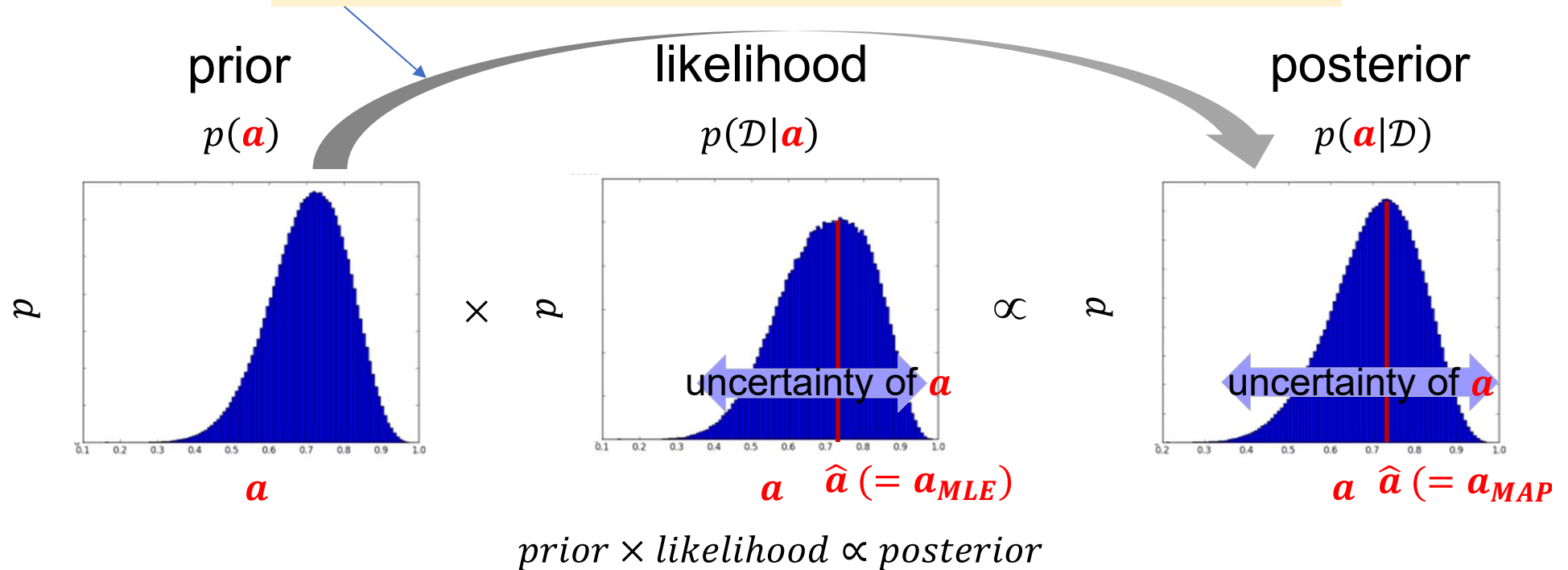
$\text{posterior} \propto \text{likelihood}$

more generalized

$\text{posterior} \propto \text{likelihood} \times \text{prior}$

Bayesian Approach

The prior (of a) is **updated** to the posterior when data given.

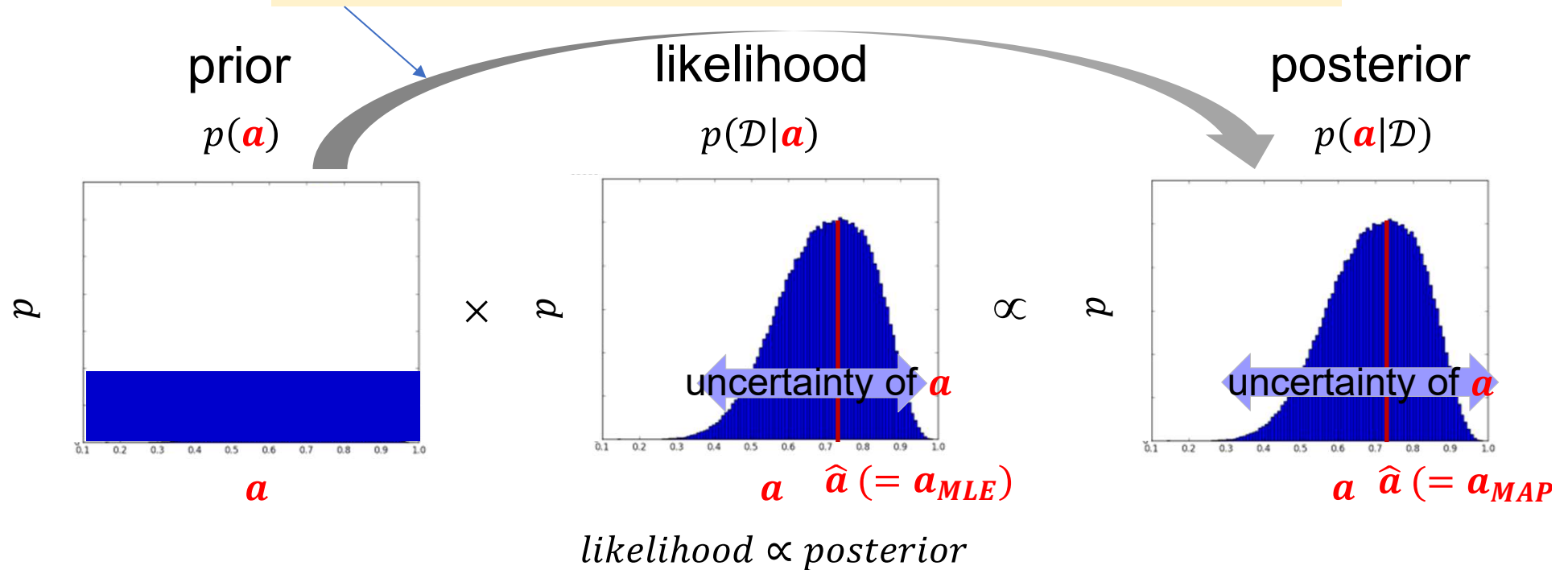


Bayes' theorem: update of probability

- Everything can be treated as a probability.
- The probability is updated by information.

Bayesian Approach

The prior (of a) is **updated** to the posterior when data given.



Bayes' theorem: update of probability

- Everything can be treated as a probability.
- The probability is updated by information.

Bayesian Approach

Another example:

You threw a coin three times.

Observed Data: all the three times they were heads.



heads heads heads

Heads

Tails



probability

μ



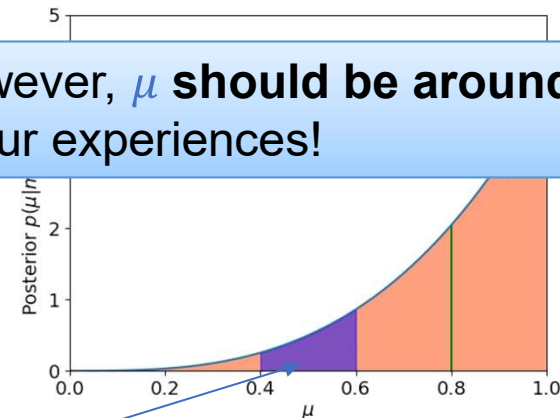
probability

$1 - \mu$



We know several concepts to evaluate μ .

However, μ should be around 0.5 in our experiences!



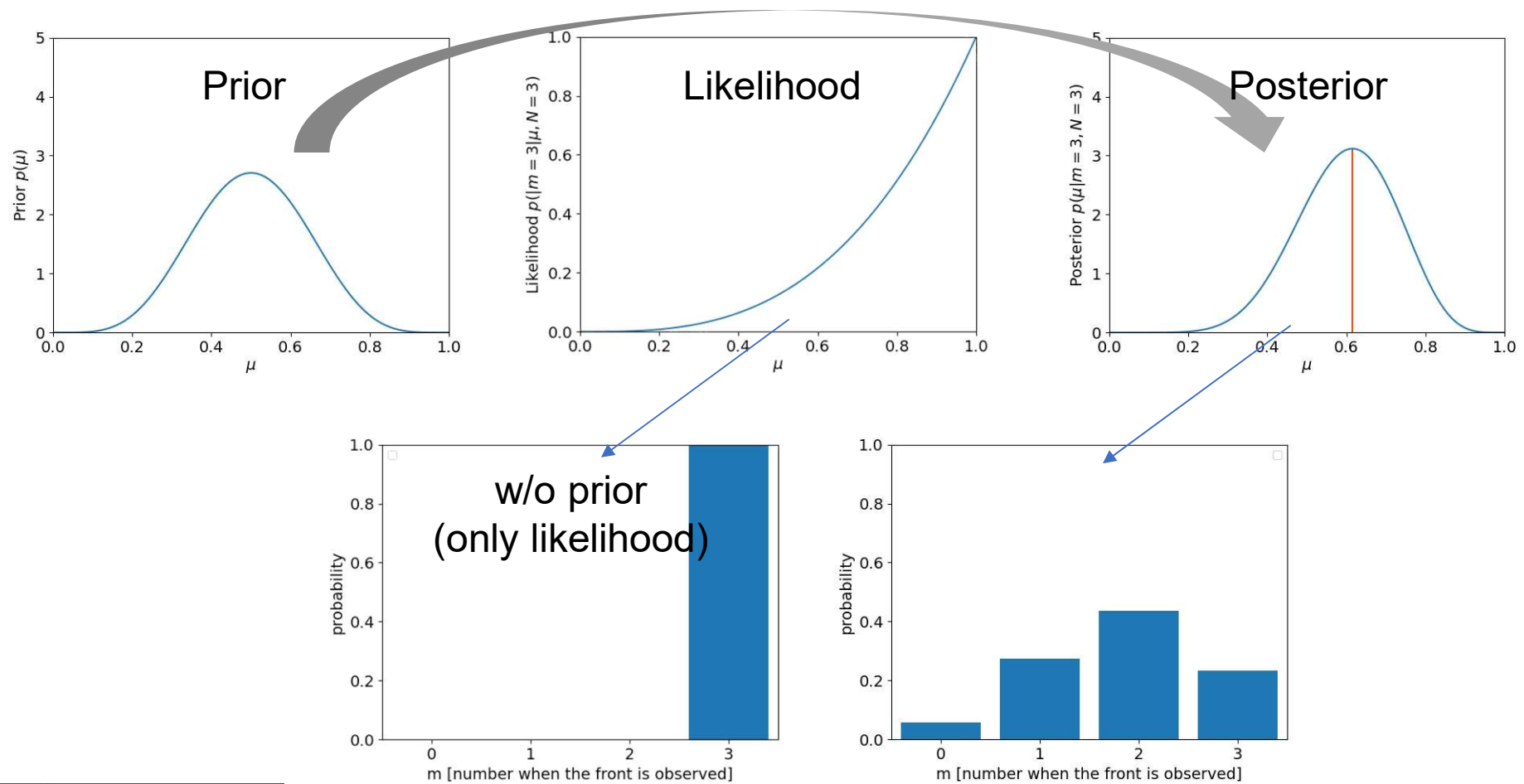
$$0 \leq \mu \leq 1$$

But the probability when $0.4 \leq \mu \leq 0.6$ is only 10.4%.

IFL

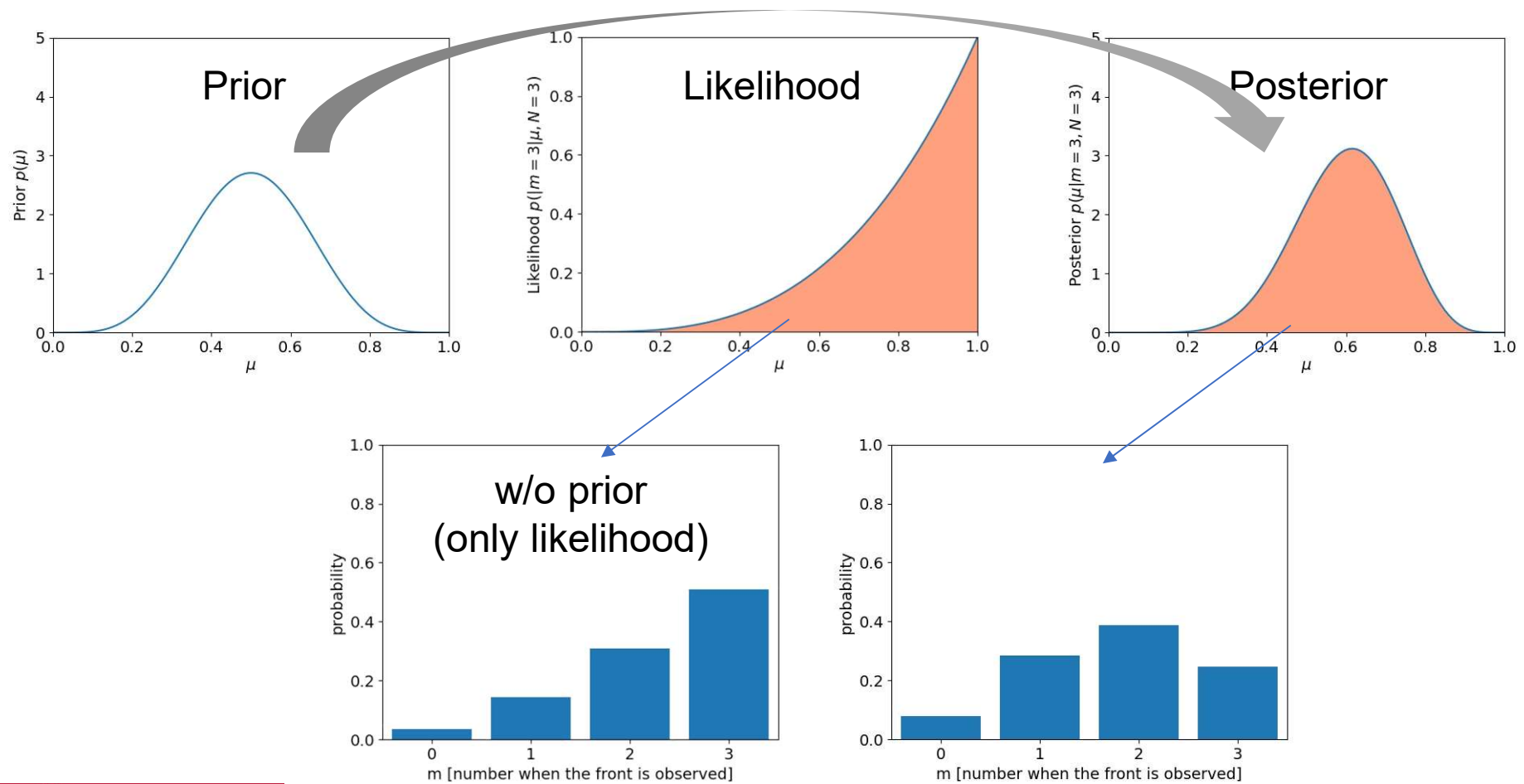
Bayesian Approach

If we have strong belief that μ is around 0.5.



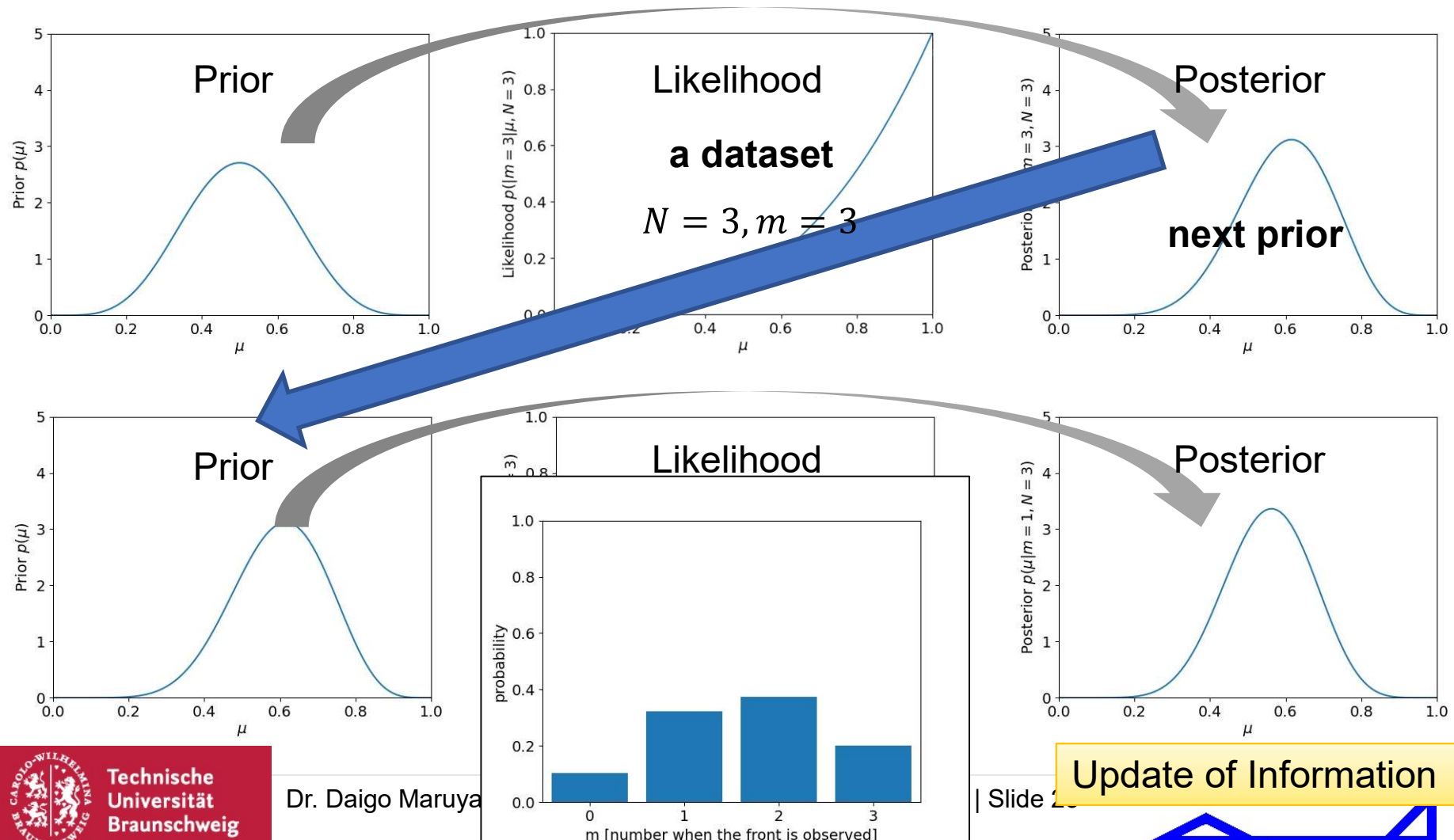
Bayesian Approach

If we have strong belief that μ is around 0.5.



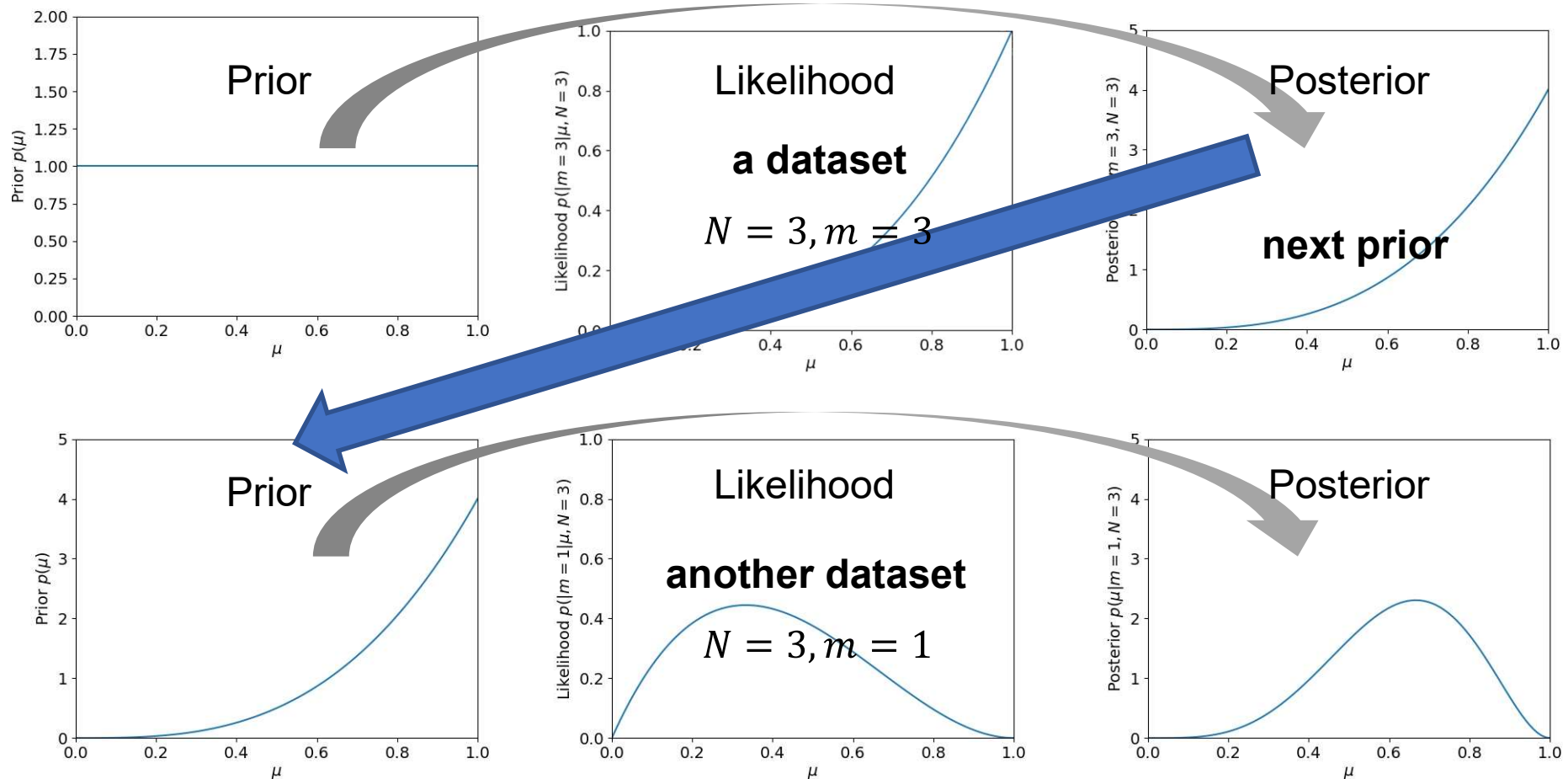
Bayesian Approach

If we have strong belief that μ is around 0.5.



Bayesian Approach

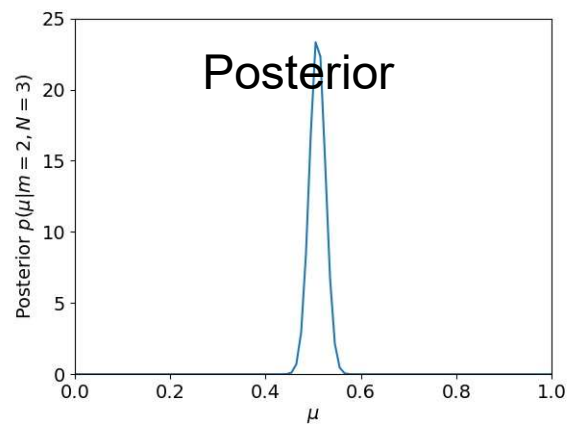
Even if we start from “no information”.



Update of Information

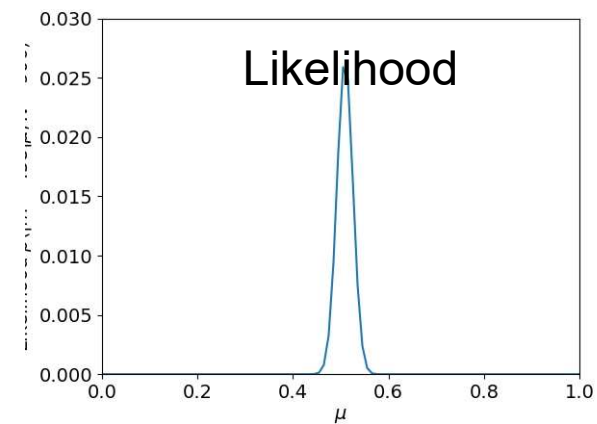
Bayesian Approach

... after 300 datasets
(300 times 3 (=900) trials) ...



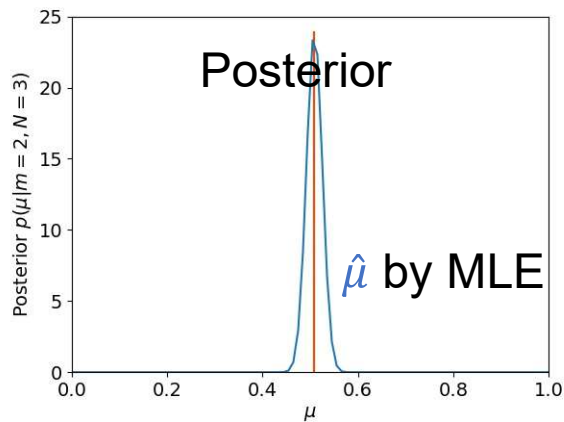
due to the law of large numbers

$N = 900, m = 458$

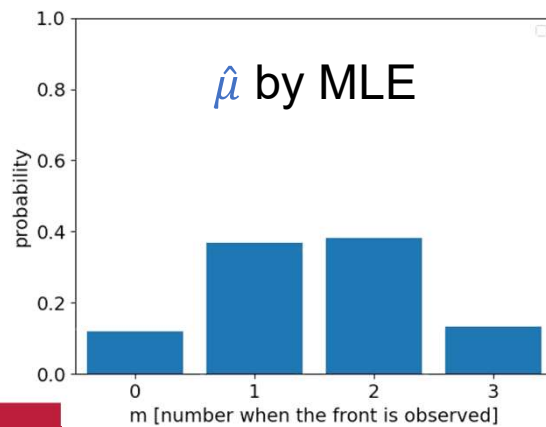
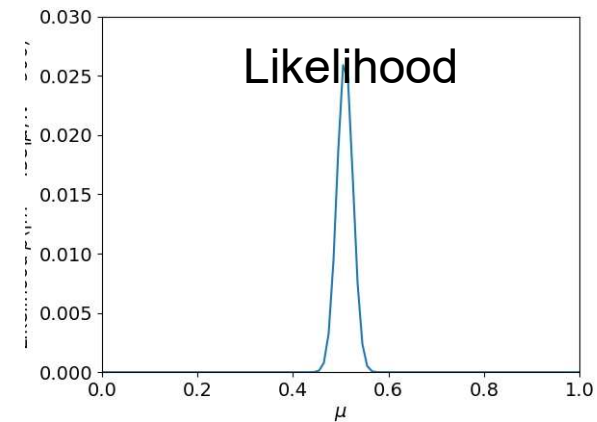


Bayesian Approach

... after 300 datasets
(300 times 3 (=900) trials) ...



$N = 900, m = 458$



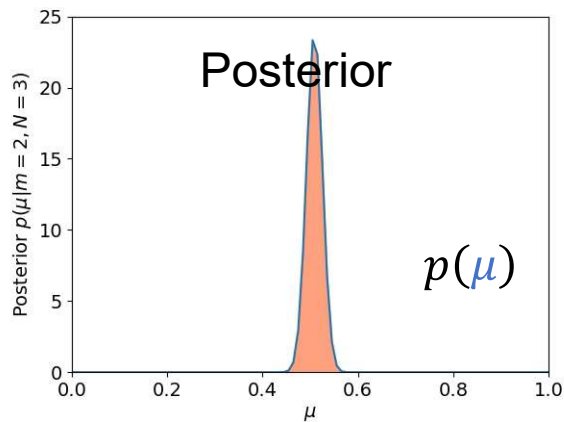
It does not influence to the result if we use,

- $\hat{\mu}$
- $p(\mu)$

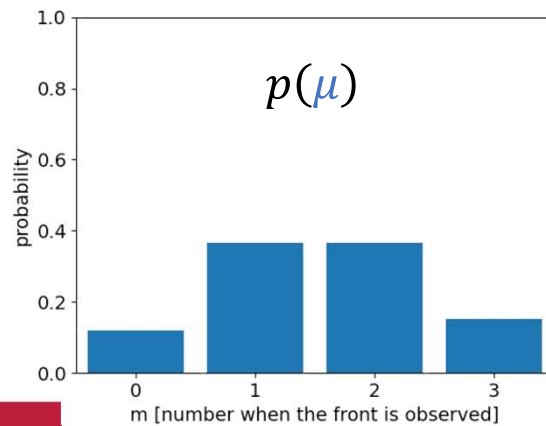
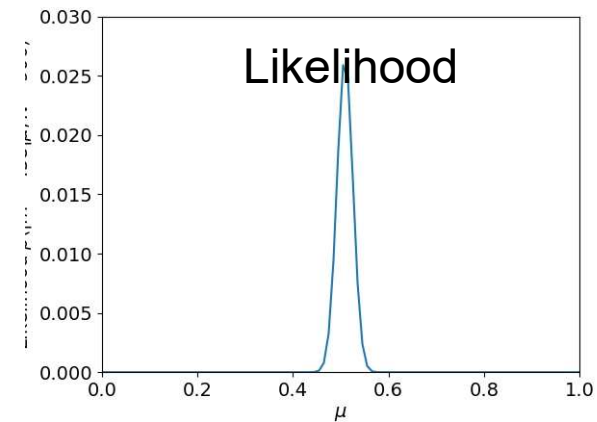
- Uncertainty
- Big data

Bayesian Approach

... after 300 datasets
(300 times 3 (=900) trials) ...



$N = 900, m = 458$



It does not influence to the result if we use,

- $\hat{\mu}$
- $p(\mu)$

- Uncertainty
- Big data

Bayesian Approach

Data treatment Parameter treatment	Batch
Point estimate	Frequentist approach



Bayesian Approach

<div>Data treatment</div> <div>Parameter treatment</div>	Batch	Sequential
Point estimate	Frequentist approach	Bayesian approach
Probability Distribution	Bayesian approach	Bayesian approach



Bayesian Approach

1. Define a probabilistic model
2. Then, **MLE**

Probabilistic model

$$p(m|\mu) = \binom{3}{m} \mu^m (1 - \mu)^{3-m}$$



A probability of m
when data \mathcal{D} is given

$$p(m|\mathcal{D}) = p(m|\hat{\mu}) \quad \text{goal}$$



Likelihood function

$$\mathcal{D}: m_{Data} = 3$$

$$p(\mathcal{D}|\mu) = \mu^3$$



• **Deterministic** $\hat{\mu}$

by MLE

Bayesian Approach

1. Define a probabilistic model
2. Then, **compute the posterior** (several choices)

Probabilistic model

$$p(m|\mu) = \binom{3}{m} \mu^m (1 - \mu)^{3-m}$$



$$p(m|\mathcal{D}) = \int p(m|\mu) p(\mu|\mathcal{D}) d\mu$$



Why this formulation?

Likelihood function

$$\mathcal{D}: m_{Data} = 3$$

$$p(\mathcal{D}|\mu) = \mu^3$$



Prior distribution

$$p(\mu): \text{arbitrary}$$



Posterior distribution $p(\mu|\mathcal{D})$

- **Deterministic** $\hat{\mu}$
Point estimate
- **Stochastic** $p(\mu)$
Probability distribution

by MLE
by mean

Probability Theory (Rules of Probability) – Review Lecture 2

sum rule $p(y) = \int p(x, y) dx$

product rule $p(x, y) = p(x|y)p(y)$

Bayes' theorem $p(y|x) = \frac{p(x|y)p(y)}{p(x)}$

marginal distribution

$$\begin{aligned} p(y) &= \int p(x, y) dx \\ &= \int p(y|x)p(x) dx \end{aligned}$$

We need to get familiar with this transformation process.

Bayesian Approach

$$p(y) = \int p(y|x)p(x)dx$$

by the sum and product rules

$y \rightarrow m$
 $x \rightarrow \mu$

$$p(m) = \int \underbrace{p(m|\mu)}_{\text{Probabilistic model}} \underbrace{p(\mu)}_{\text{A probability of } \mu} d\mu$$

Probabilistic
model

A probability of μ
prior
posterior

$p(\mu) \rightarrow p(\mu|\mathcal{D})$

$$\underbrace{p(m|\mathcal{D})}_{\text{A probability of } m \text{ when data } \mathcal{D} \text{ is given}} = \int \underbrace{p(m|\mu)}_{\text{Probabilistic model}} \underbrace{p(\mu|\mathcal{D})}_{\text{A probability of } \mu \text{ when data } \mathcal{D} \text{ is given}} d\mu$$

Posterior

A probability of μ
when data \mathcal{D} is given

A probability of m
when data \mathcal{D} is given



Bayesian Approach

1. Define a probabilistic model
2. Then, **compute the posterior** (several choices)

parametric probability distributions (see Lecture 3)

↓
Your model for the output,
which is parametrized by parameters

← 1. Define a **probabilistic model**

↓

$$p(\text{output}) = \int (\text{Probabilistic Model}) \times (\text{Posterior}) d(\text{parameter})$$

↑
Information of the parameters - a probability distribution
(obtained by data)

↑
Predictive distribution (goal)

↑
2. Then, **compute the posterior**


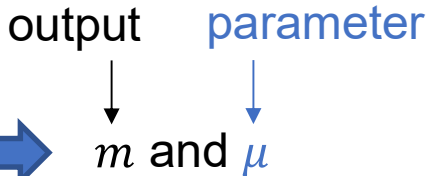
Bayesian Approach

Another Meaning / Aspect / Interpretation

$$\begin{aligned} p(m|\mathcal{D}) &= \int p(m|\mu)p(\mu|\mathcal{D})d\mu \\ &= \int p(m, \mu|\mathcal{D})d\mu \end{aligned}$$

REVIEW of Lecture 2:
Meaning of

- Joint distribution
- Marginal distribution

1. Clarify all the stochastic variables   m and μ
2. Try to find **the joint probability** $p(m, \mu)$, which include all the information.
3. Our objective is only **the marginal distribution** $p(m)$ $p(m) = \int p(m, \mu)d\mu$

We are not interested in the parameter μ anymore.

Bayesian Approach

$$p(m|\mathcal{D}) = \int p(m|\mu)p(\mu|\mathcal{D})d\mu$$

We know the **probabilistic model** $p(m|\mu)$ (because “we” defined it by ourselves).
But how can we obtain the **posterior** $p(\mu|\mathcal{D})$?

We know the **likelihood function** $p(\mathcal{D}|\mu)$.



We have been computing this to maximize it.

Posterior

$$p(\mu|\mathcal{D}) = \frac{p(\mathcal{D}|\mu)p(\mu)}{p(\mathcal{D})}$$

Likelihood

$$p(\mathcal{D}|\mu) = \mu^3$$

Prior

$p(\mu)$: arbitrary.

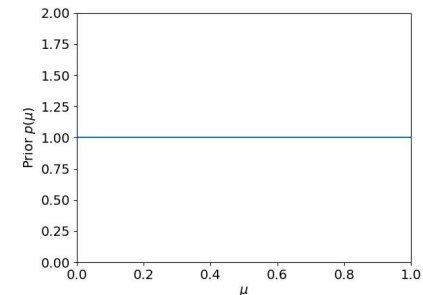
$$\int p(\mu)d\mu = 1$$

Bayesian Approach

$$\begin{aligned} p(\mu|\mathcal{D}) &= \frac{p(\mathcal{D}|\mu)p(\mu)}{p(\mathcal{D})} \\ &= \frac{\mu^3 \times 1}{p(\mathcal{D})} \\ &= 4\mu^3 \end{aligned} \quad \hookrightarrow \quad \int p(\mu|\mathcal{D})d\mu = 1$$

$p(\mu)$ can be arbitrary. e.g. $p(\mu) = 1$

$$\int_0^1 p(\mu)d\mu = 1$$
$$0 \leq \mu \leq 1$$



The posterior $p(\mu|\mathcal{D})$ can be obtained without computing $p(\mathcal{D})$.

Actually $p(\mathcal{D})$ can be computed.

$$\begin{aligned} p(\mathcal{D}) &= \int p(\mathcal{D}|\mu)p(\mu)d\mu \\ &= \int (\mu^3 \times 1)d\mu = \frac{1}{4} \end{aligned}$$

Bayes' theorem
(another expression)

$$p(\mu|\mathcal{D}) = \frac{p(\mathcal{D}|\mu)p(\mu)}{\int p(\mathcal{D}|\mu)p(\mu)d\mu}$$

Bayesian Approach

$$p(m|\mathcal{D}) = \int p(m|\mu)p(\mu|\mathcal{D})d\mu$$

$$= \int \binom{3}{m} \mu$$

The output distributions are non closed-form in real applications.
Exceptions: Gaussian Processes (Lectures 7,8)

$$= 4 \binom{3}{m} \int \underbrace{\mu^{3+m}(1-\mu)^{3-m}}_{\text{any closed-form expressions?}} d\mu$$

= ...

any closed-form
expressions?

$$p(m|\mathcal{D}) = \iiint \cdots d\mu_1 d\mu_2 d\mu_3 \cdots$$

In practical applications, the number of the parameters is large and $p(m|\mathcal{D})$ requires multiple integral.



In practice, numerical approaches are used.
(e.g. MCMC presented in Lecture 2 and Lecture 4
– explained in Lecture 12)

The results in the slide 16 (and the above equation) was obtained by using the MCMC.

Bayesian Approach

Let's try to apply this concept to **the curve fitting problem**.

Probabilistic model

$$p(t|x, \mathbf{w}, \sigma) = \mathcal{N}(t|y(x, \mathbf{w}), \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left\{ -\frac{\{t - y(x, \mathbf{w})\}^2}{2\sigma^2} \right\}$$

Likelihood function

$$p(\mathbf{t}|\mathbf{x}, \mathbf{w}, \sigma) = \prod_{i=1}^N \mathcal{N}(t_i|y(x_i, \mathbf{w}), \sigma^2) = \prod_{i=1}^N \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left[-\frac{\{t_i - y(x_i, \mathbf{w})\}^2}{2\sigma^2} \right]$$

➡ **Posterior distribution** $p(\mathbf{w}, \sigma|\mathbf{x}, \mathbf{t})$

Predictive distribution

$$p(t|x, \mathbf{x}, \mathbf{t}) = \iint p(t|x, \mathbf{w}, \sigma) p(\mathbf{w}, \sigma|\mathbf{x}, \mathbf{t}) d\mathbf{w} d\sigma = ?$$

probabilistic model \times posterior

Summary

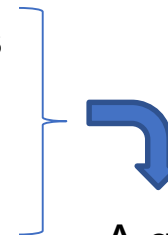
Process

1. Define a probabilistic model
2. Then, **compute the posterior** (several choices)
 - Point estimate
 - Probability distribution ← computations hard

Next lecture

- Connections with the regularization techniques
- Uncertainty
- Large numbers of dataset
- etc.

Using the curve fitting problem



A generalized perspective
by the Bayes approach