

~ "AI goes MAD"

IFT, June 2022

BRYAN ZALDIVAR (IFIC)

---

~ "AI goes MAD"

IFT, June 2022

BRYAN ZALDIVAR (IFIC)

---

A Review on Modern Bayesian Inference

~ "AI goes MAD"

IFT, June 2022

BRYAN ZALDIVAR (IFIC)

---

"AI goes MAD"

IFT, June 2022

BRYAN ZALDIVAR (IFIC)

---

Some open research lines in B.I.



## INTRODUCTION

---

## INTRODUCTION

---

- What is Bayesian Inference?

## INTRODUCTION

---

- What is Bayesian Inference?

A way to update the "degree of believe" of your hypotheses after seeing some data

## INTRODUCTION

---

- What is Bayesian Inference?

A way to update the "degree of believe" of your hypotheses after seeing some data

Hyp.  $A$  : e.g.  $A$  = Higgs mass

$p(A)$   $\Rightarrow$  prior dist.

$p(\text{Data} | A)$   $\Rightarrow$  Likelihood of data

## INTRODUCTION

- What is Bayesian Inference?

A way to update the "degree of believe" of your hypotheses after seeing some data

$$\underbrace{p(A|Data)}_{\text{Posterior dist.}} = \frac{p(Data|A) p(A)}{p(Data)}$$

Hyp. A : e.g. A = Higgs mass

$p(A) \Rightarrow$  prior dist.

$p(Data|A) \Rightarrow$  Likelihood of data

## INTRODUCTION

- What is Bayesian Inference?

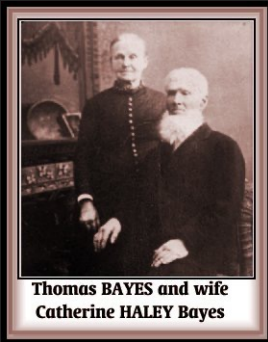
A way to update the "degree of believe" of your hypotheses after seeing some data

$$\underbrace{p(A|Data)}_{\text{Posterior dist.}} = \frac{p(Data|A) p(A)}{p(Data)}$$

Hyp. A : e.g. A = Higgs mass

$p(A)$   $\Rightarrow$  prior dist.

$p(Data|A)$   $\Rightarrow$  Likelihood of data



Thomas BAYES and wife  
Catherine HALEY Bayes

## INTRODUCTION

- What is Bayesian Inference?

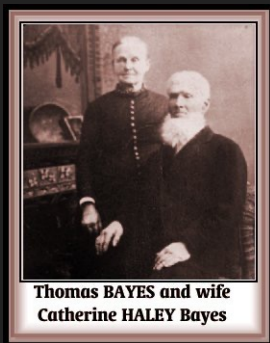
A way to update the "degree of believe" of your hypotheses after seeing some data

$$\underbrace{p(A|Data)}_{\text{Posterior dist.}} = \frac{p(Data|A) p(A)}{p(Data)}$$

Hyp. A : e.g. A = Higgs mass

$p(A) \Rightarrow$  prior dist.

$p(Data|A) \Rightarrow$  Likelihood of data



As opposed to "Frequentist" approach

$$p(A) = \frac{\# \text{ of occurrences of } A}{\text{Total \# of trials}}$$

A : an observable (not a hypothesis)

## INTRODUCTION

---



## INTRODUCTION

---

- What is it used for?

## INTRODUCTION

---

- What is it used for?

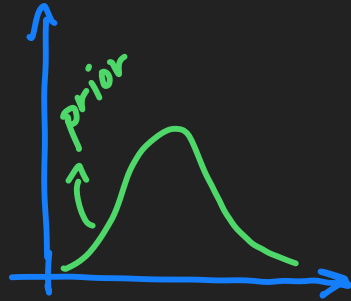
- 1) Making inference on physical parameters

## INTRODUCTION

---

- What is it used for?

- 1) Making inference on physical parameters

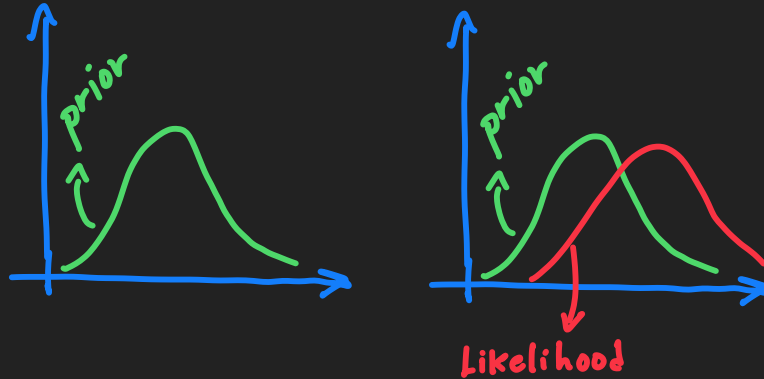


## INTRODUCTION

---

- What is it used for?

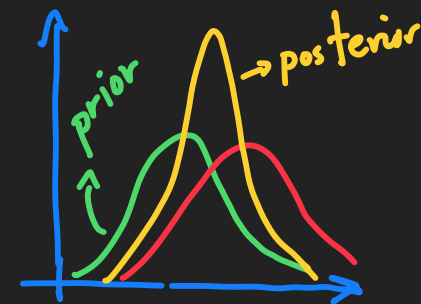
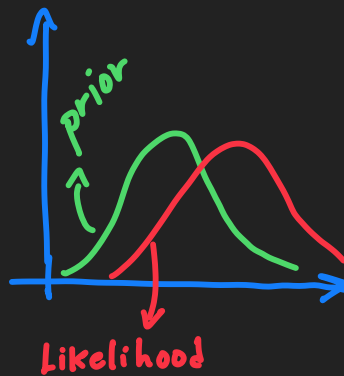
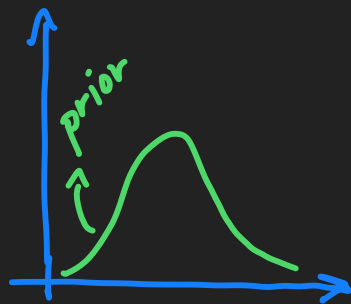
- 1) Making inference on physical parameters



## INTRODUCTION

- What is it used for?

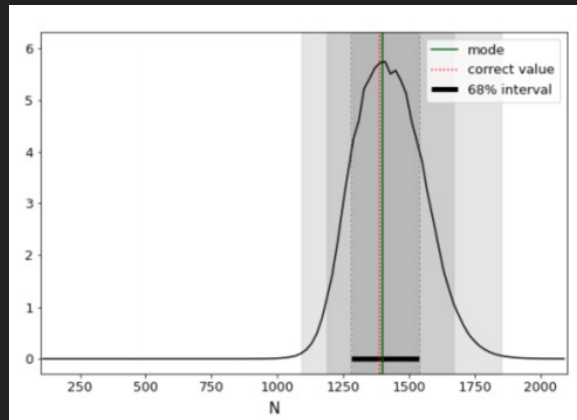
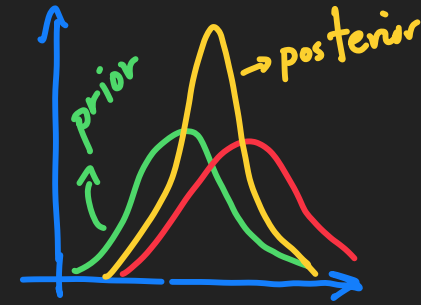
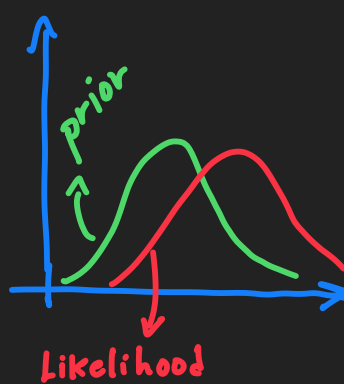
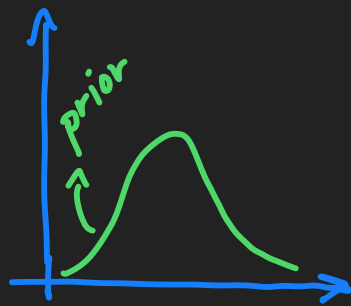
1) Making inference on physical parameters



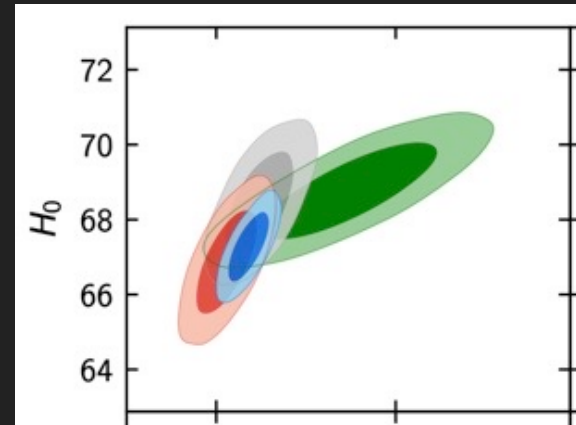
## INTRODUCTION

- What is it used for?

### 1) Making inference on physical parameters



(CREDITS: Dimitriu et al, (to appear))



Planck Collaboration 1807.06209

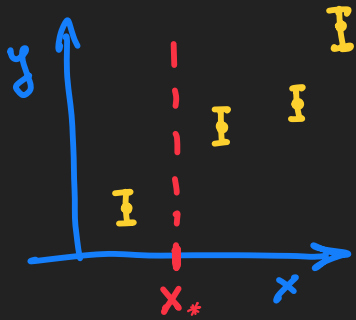
## INTRODUCTION

---

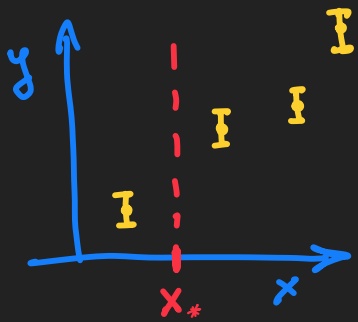
2) Obtaining predictions with uncertainties



### 2) Obtaining predictions with uncertainties



## 2) Obtaining predictions with uncertainties

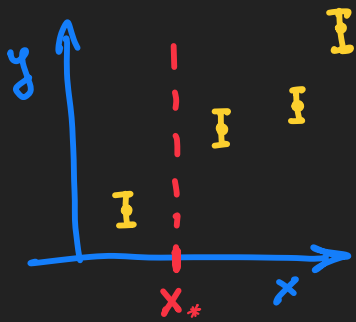


$$P(y_* | x_*, \text{Data}) = \int d\vec{\theta} \left( \begin{array}{c} \text{Likelihood} \\ \text{of } y_* | x_*, \vec{\theta} \end{array} \right) \cdot \left( \begin{array}{c} \text{Posterior of} \\ \vec{\theta} | \text{Data} \end{array} \right)$$

model parameters

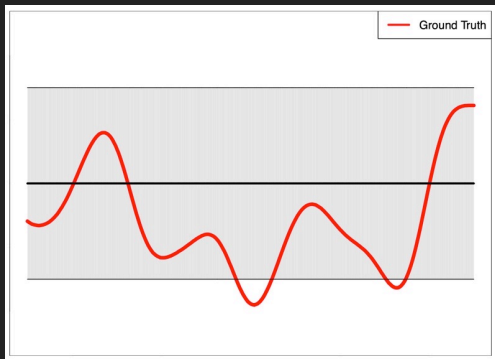
## INTRODUCTION

### 2) Obtaining predictions with uncertainties



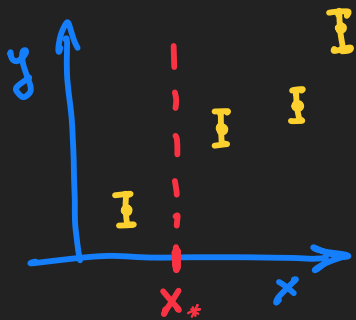
$$P(y_* | x_*, \text{Data}) = \int d\vec{\theta} \left( \begin{array}{c} \text{Likelihood} \\ \text{of } y_* | x_*, \vec{\theta} \end{array} \right) \cdot \left( \begin{array}{c} \text{Posterior of} \\ \vec{\theta} | \text{Data} \end{array} \right)$$

↖ model parameters



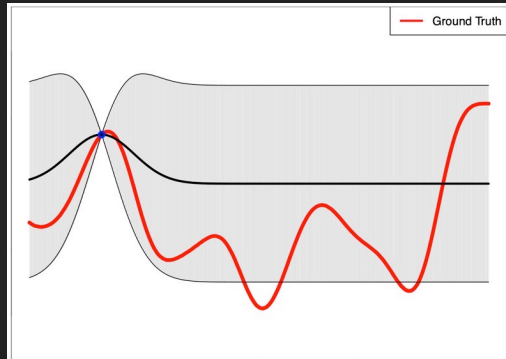
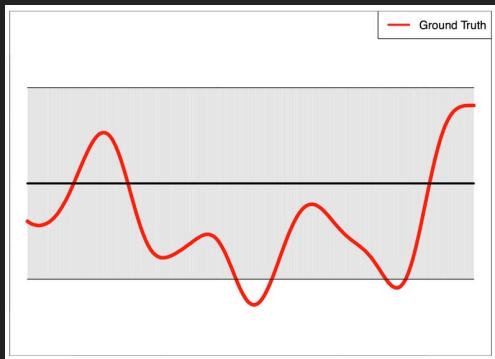
## INTRODUCTION

### 2) Obtaining predictions with uncertainties



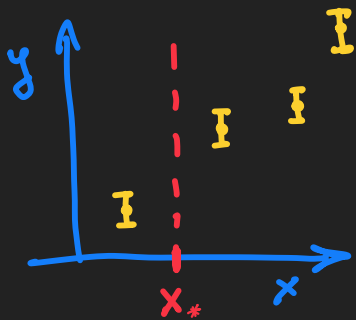
$$P(y_* | x_*, \text{Data}) = \int d\vec{\theta} \left( \text{Likelihood of } y_* | x_*, \vec{\theta} \right) \cdot \left( \text{Posterior of } \vec{\theta} | \text{Data} \right)$$

model parameters



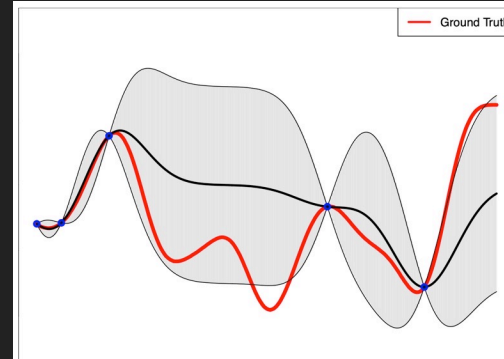
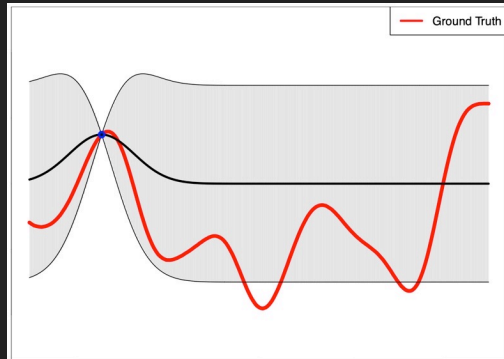
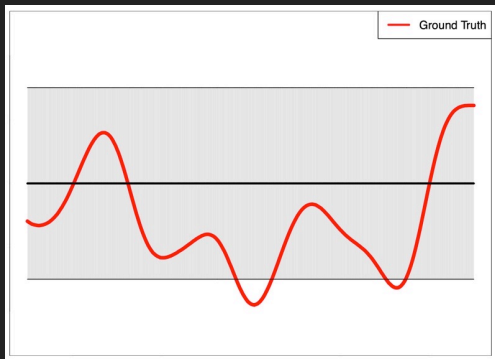
## INTRODUCTION

### 2) Obtaining predictions with uncertainties



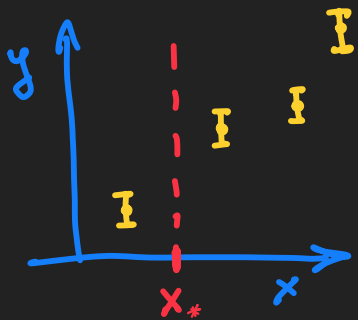
$$P(y_* | x_*, \text{Data}) = \int d\vec{\theta} \left( \text{Likelihood of } y_* | x_*, \vec{\theta} \right) \cdot \left( \text{Posterior of } \vec{\theta} | \text{Data} \right)$$

model parameters



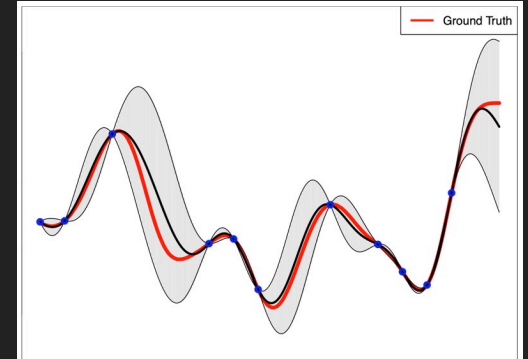
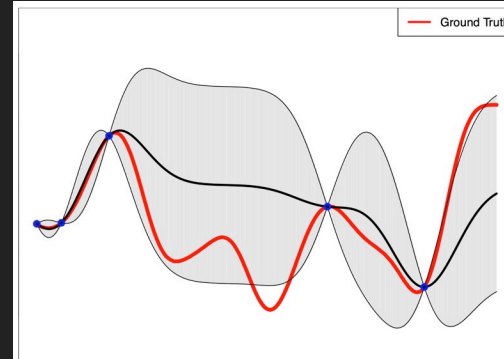
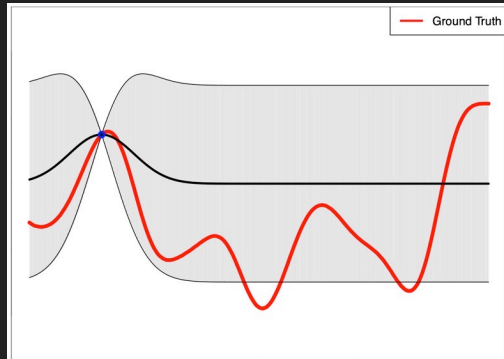
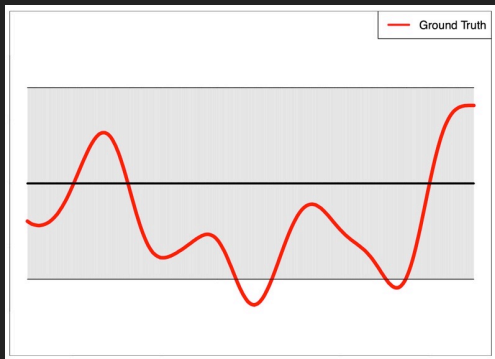
## INTRODUCTION

### 2) Obtaining predictions with uncertainties



$$P(y_* | x_*, \text{Data}) = \int d\vec{\theta} \left( \text{Likelihood of } y_* | x_*, \vec{\theta} \right) \cdot \left( \text{Posterior of } \vec{\theta} | \text{Data} \right)$$

*model parameters* (pointing to  $d\vec{\theta}$ )



## METHODS

---

## METHODS

---

- Bayesian Inference is computationally challenging!



- Bayesian Inference is computationally challenging!

$$p(\vec{\theta} | \text{Data}) = p(\text{Data} | \vec{\theta}) p(\vec{\theta}) / p(\text{Data})$$

- Bayesian Inference is computationally challenging!

$$p(\vec{\theta} | \text{Data}) = p(\text{Data} | \vec{\theta}) p(\vec{\theta}) / p(\text{Data})$$

$$p(\text{Data}) = \int d\vec{\theta} \, p(\text{Data} | \vec{\theta}) p(\vec{\theta})$$

## METHODS

---

- Bayesian Inference is computationally challenging!

$$p(\vec{\theta} | \text{Data}) = p(\text{Data} | \vec{\theta}) p(\vec{\theta}) / p(\text{Data})$$

$$p(\text{Data}) = \int d\vec{\theta} \, p(\text{Data} | \vec{\theta}) p(\vec{\theta}) \quad \left. \vphantom{\int d\vec{\theta} \, p(\text{Data} | \vec{\theta}) p(\vec{\theta})} \right\} \begin{array}{l} \text{typically} \\ \text{intractable} \end{array}$$

## METHODS

---

- Bayesian Inference is computationally challenging!

$$p(\vec{\theta} | \text{Data}) = p(\text{Data} | \vec{\theta}) p(\vec{\theta}) / p(\text{Data}) \rightarrow \text{Needs to be approximated}$$

$$p(\text{Data}) = \int d\vec{\theta} \, p(\text{Data} | \vec{\theta}) p(\vec{\theta}) \quad \left. \vphantom{\int d\vec{\theta} \, p(\text{Data} | \vec{\theta}) p(\vec{\theta})} \right\} \begin{array}{l} \text{typically} \\ \text{intractable} \end{array}$$

## METHODS

---

- Bayesian Inference is computationally challenging!

$$p(\vec{\theta} | \text{Data}) = p(\text{Data} | \vec{\theta}) p(\vec{\theta}) / p(\text{Data}) \rightarrow \text{Needs to be approximated}$$

$$p(\text{Data}) = \int d\vec{\theta} p(\text{Data} | \vec{\theta}) p(\vec{\theta}) \quad \left. \vphantom{\int d\vec{\theta} p(\text{Data} | \vec{\theta}) p(\vec{\theta})} \right\} \begin{array}{l} \text{typically} \\ \text{intractable} \end{array}$$

- Approximate methods

- Bayesian Inference is computationally challenging!

$$p(\vec{\theta} | \text{Data}) = p(\text{Data} | \vec{\theta}) p(\vec{\theta}) / p(\text{Data}) \rightarrow \text{Needs to be approximated}$$

$$p(\text{Data}) = \int d\vec{\theta} p(\text{Data} | \vec{\theta}) p(\vec{\theta}) \quad \left. \vphantom{\int d\vec{\theta} p(\text{Data} | \vec{\theta}) p(\vec{\theta})} \right\} \begin{array}{l} \text{typically} \\ \text{intractable} \end{array}$$

- Approximate methods

### ① Sampling methods

e.g. MCMC,  
ABC, ...

- Bayesian Inference is computationally challenging!

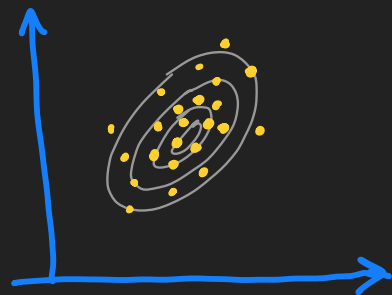
$$p(\vec{\theta} | \text{Data}) = p(\text{Data} | \vec{\theta}) p(\vec{\theta}) / p(\text{Data}) \rightarrow \text{Needs to be approximated}$$

$$p(\text{Data}) = \int d\vec{\theta} p(\text{Data} | \vec{\theta}) p(\vec{\theta}) \quad \left. \vphantom{\int d\vec{\theta} p(\text{Data} | \vec{\theta}) p(\vec{\theta})} \right\} \begin{array}{l} \text{typically} \\ \text{intractable} \end{array}$$

- Approximate methods

### ① Sampling methods

e.g. MCMC,  
ABC, ...



- Bayesian Inference is computationally challenging!

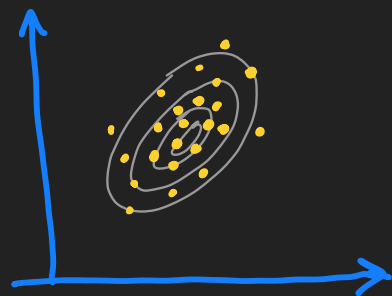
$$p(\vec{\theta} | \text{Data}) = p(\text{Data} | \vec{\theta}) p(\vec{\theta}) / p(\text{Data}) \rightarrow \text{Needs to be approximated}$$

$$p(\text{Data}) = \int d\vec{\theta} p(\text{Data} | \vec{\theta}) p(\vec{\theta}) \quad \left. \vphantom{\int d\vec{\theta} p(\text{Data} | \vec{\theta}) p(\vec{\theta})} \right\} \begin{array}{l} \text{typically} \\ \text{intractable} \end{array}$$

- Approximate methods

### ① Sampling methods

e.g. MCMC,  
ABC, ...



### ② Optimization methods

e.g. Laplace,  
Variational Inference,  
N. ratio-estimation, ...



## METHODS

- Bayesian Inference is computationally challenging!

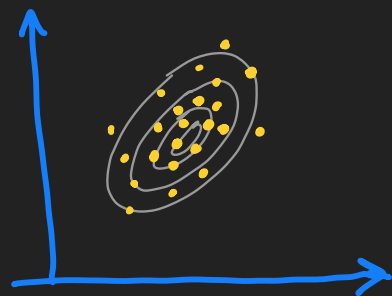
$$p(\vec{\theta} | \text{Data}) = p(\text{Data} | \vec{\theta}) p(\vec{\theta}) / p(\text{Data}) \rightarrow \text{Needs to be approximated}$$

$$p(\text{Data}) = \int d\vec{\theta} p(\text{Data} | \vec{\theta}) p(\vec{\theta}) \quad \left. \vphantom{\int d\vec{\theta} p(\text{Data} | \vec{\theta}) p(\vec{\theta})} \right\} \begin{array}{l} \text{typically} \\ \text{intractable} \end{array}$$

- Approximate methods

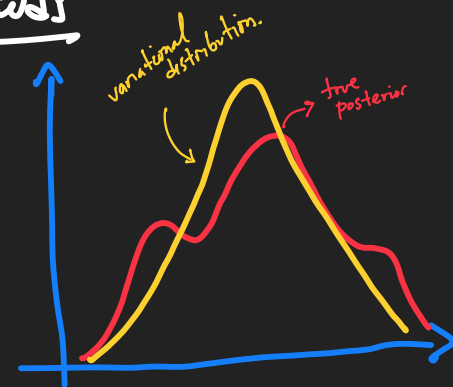
### ① Sampling methods

e.g. MCMC,  
ABC, ...



### ② Optimization methods

e.g. Laplace,  
Variational Inference,  
N. ratio-estimation, ...



## VARIATIONAL INFERENCE

---

## VARIATIONAL INFERENCE

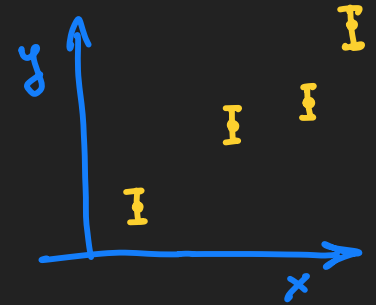
---

- Data  $\{x_i, y_i\}_{i=1}^N$

## VARIATIONAL INFERENCE

---

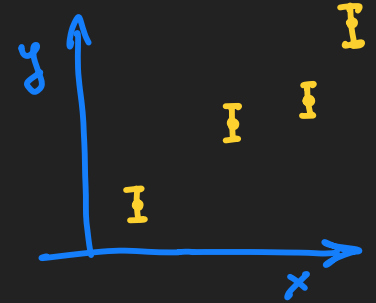
- Data  $\{x_i, y_i\}_{i=1}^N$



## VARIATIONAL INFERENCE

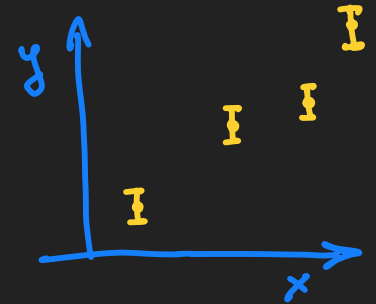
---

- Data  $\{x_i, y_i\}_{i=1}^N$
- Fitting function  $f(x; \vec{\theta})$



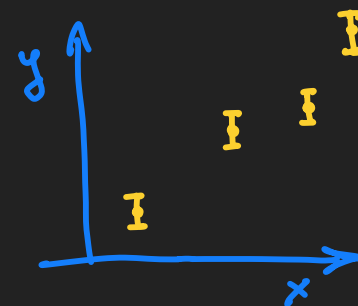
## VARIATIONAL INFERENCE

- Data  $\{x_i, y_i\}_{i=1}^N$
- Fitting function  $f(x; \vec{\theta})$
- Likelihood  $p(\vec{y} | \vec{\theta}) = \prod_{i=1}^N N(y_i | f(x_i, \vec{\theta}), \sigma_i)$   
↳ known uncertainty



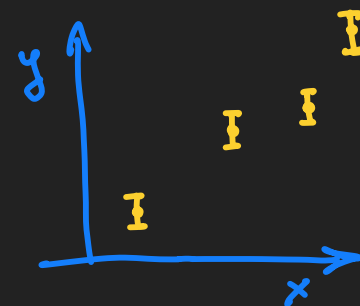
## VARIATIONAL INFERENCE

- Data  $\{x_i, y_i\}_{i=1}^N$
- Fitting function  $f(x; \vec{\theta})$
- Likelihood  $p(\vec{y} | \vec{\theta}) = \prod_{i=1}^N N(y_i | f(x_i, \vec{\theta}), \sigma_i)$   
 $\sigma_i$  known uncertainty
- Prior distrib. for  $\vec{\theta}$ :  $p(\vec{\theta}) \Rightarrow$  e.g. also Gaussian



## VARIATIONAL INFERENCE

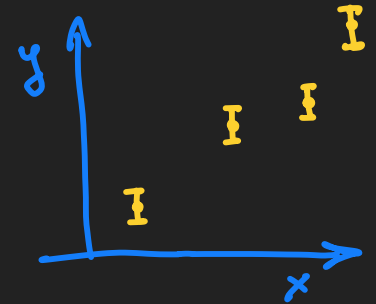
- Data  $\{x_i, y_i\}_{i=1}^N$
- Fitting function  $f(x; \vec{\theta})$
- Likelihood  $p(\vec{y} | \vec{\theta}) = \prod_{i=1}^N N(y_i | f(x_i, \vec{\theta}), \sigma_i)$   
 $\sigma_i$  known uncertainty
- Prior distrib. for  $\vec{\theta}$ :  $p(\vec{\theta}) \Rightarrow$  e.g. also Gaussian
- Posterior  $p(\vec{\theta} | \vec{y})$  not Gaussian in general  $\Rightarrow$  Typically intractable





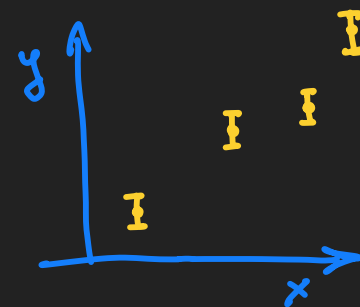
## VARIATIONAL INFERENCE

- Data  $\{x_i, y_i\}_{i=1}^N$
- Fitting function  $f(x; \vec{\theta})$
- Likelihood  $p(\vec{y} | \vec{\theta}) = \prod_{i=1}^N N(y_i | f(x_i, \vec{\theta}), \sigma_i)$   
 $\sigma_i$  known uncertainty
- Prior distrib. for  $\vec{\theta}$ :  $p(\vec{\theta}) \Rightarrow$  e.g. also Gaussian
- Posterior  $p(\vec{\theta} | \vec{y})$  not Gaussian in general  $\Rightarrow$  Typically intractable
- $p(\vec{\theta} | \vec{y}) \approx q(\vec{\theta} | \vec{\eta})$ 
  - Tractable
  - Expressive enough



## VARIATIONAL INFERENCE

- Data  $\{x_i, y_i\}_{i=1}^N$
  - Fitting function  $f(x; \vec{\theta})$
  - Likelihood  $p(\vec{y} | \vec{\theta}) = \prod_{i=1}^N N(y_i | f(x_i, \vec{\theta}), \sigma_i)$   
 $\sigma_i$  known uncertainty
  - Prior distrib. for  $\vec{\theta}$ :  $p(\vec{\theta}) \Rightarrow$  e.g. also Gaussian
  - Posterior  $p(\vec{\theta} | \vec{y})$  not Gaussian in general  $\Rightarrow$  Typically intractable
  - $p(\vec{\theta} | \vec{y}) \approx q(\vec{\theta} | \vec{\eta})$ 
    - Tractable
    - Expressive enough
- Need a procedure to optimize parameters  $\vec{\eta}$  s.t. this approx. is as good as possible



## VARIATIONAL INFERENCE

---

## VARIATIONAL INFERENCE

---

- Standard procedure: Minimize the KL divergence:

$$KL[q_{\vec{\eta}} | p] = \int d\vec{\theta} \cdot q_{\vec{\eta}} \cdot \ln\left(\frac{q_{\vec{\eta}}}{p}\right) \geq 0$$

## VARIATIONAL INFERENCE

---

- Standard procedure: Minimize the KL divergence:

$$KL[q_{\vec{\eta}} | p] = \int d\vec{\theta} \cdot q_{\vec{\eta}} \cdot \ln\left(\frac{q_{\vec{\eta}}}{p}\right) \geq 0$$

$\hookrightarrow$  true posterior

## VARIATIONAL INFERENCE

---

- Standard procedure: Minimize the KL divergence:

$$KL[q_{\vec{\eta}} | p] = \int d\vec{\theta} \cdot q_{\vec{\eta}} \cdot \ln\left(\frac{q_{\vec{\eta}}}{p}\right) \geq 0$$

proposed approximation  $\nwarrow$   $\searrow$  true posterior

## VARIATIONAL INFERENCE

- Standard procedure: Minimize the KL divergence:

$$KL[q_{\vec{\eta}} | p] = \int d\vec{\theta} \cdot q_{\vec{\eta}} \cdot \ln\left(\frac{q_{\vec{\eta}}}{p}\right) \geq 0$$

proposed approximation  $\nwarrow$   $\nearrow$  true posterior

useless since it depends on  $p$

## VARIATIONAL INFERENCE

- Standard procedure: Minimize the KL divergence:

$$KL[q_{\vec{\eta}} | p] = \int d\vec{\theta} \cdot q_{\vec{\eta}} \cdot \ln\left(\frac{q_{\vec{\eta}}}{p}\right) \geq 0$$

proposed approximation  $\nwarrow$   $\nearrow$  true posterior

useless since it depends on  $p$

- Intractability of  $p$  typically due to its normalization ("Evidence")



## VARIATIONAL INFERENCE

- Standard procedure: Minimize the KL divergence:

$$KL[q_{\vec{\eta}} | p] = \int d\vec{\theta} \cdot q_{\vec{\eta}} \cdot \ln\left(\frac{q_{\vec{\eta}}}{p}\right) \geq 0$$

proposed approximation  $\nwarrow$   $\nearrow$  true posterior

useless since it depends on  $p$

- Intractability of  $p$  typically due to its normalization ("Evidence")  
 $p = \tilde{p} / E$

## VARIATIONAL INFERENCE

- Standard procedure: Minimize the KL divergence:

$$KL[q_{\vec{\eta}} | p] = \int d\vec{\theta} \cdot q_{\vec{\eta}} \cdot \ln\left(\frac{q_{\vec{\eta}}}{p}\right) \geq 0$$

proposed approximation  $\nwarrow$   $\nearrow$  true posterior

useless since it depends on  $p$

- Intractability of  $p$  typically due to its normalization ("Evidence")

$$p = \tilde{p} / E$$

$$ELBO = \int d\vec{\theta} \cdot q_{\vec{\eta}} \ln\left(\frac{\tilde{p}}{q_{\vec{\eta}}}\right)$$

## VARIATIONAL INFERENCE

- Standard procedure: Minimize the KL divergence:

$$KL[q_{\vec{\eta}} | p] = \int d\vec{\theta} \cdot q_{\vec{\eta}} \cdot \ln\left(\frac{q_{\vec{\eta}}}{p}\right) \geq 0$$

proposed approximation  $\nwarrow$   $\nearrow$  true posterior

useless since it depends on  $p$

- Intractability of  $p$  typically due to its normalization ("Evidence")

$$p = \tilde{p} / E$$

$$ELBO = \int d\vec{\theta} \cdot q_{\vec{\eta}} \ln\left(\frac{\tilde{p}}{q_{\vec{\eta}}}\right)$$

Maximize wrt  $\vec{\eta}$

## VARIATIONAL INFERENCE

- Standard procedure: Minimize the KL divergence:

$$KL[q_{\vec{\eta}} | p] = \int d\vec{\theta} \cdot q_{\vec{\eta}} \cdot \ln\left(\frac{q_{\vec{\eta}}}{p}\right) \geq 0$$

proposed approximation  $\leftarrow$   $\rightarrow$  true posterior

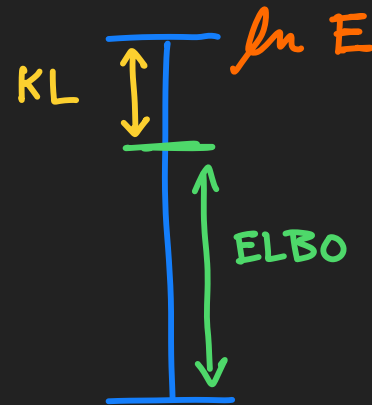
useless since it depends on  $p$

- Intractability of  $p$  typically due to its normalization ("Evidence")

$$p = \tilde{p} / E$$

$$ELBO = \int d\vec{\theta} \cdot q_{\vec{\eta}} \ln\left(\frac{\tilde{p}}{q_{\vec{\eta}}}\right)$$

Maximize wrt  $\vec{\eta}$



## VARIATIONAL INFERENCE

- Standard procedure: Minimize the KL divergence:

$$KL[q_{\vec{\eta}} | p] = \int d\vec{\theta} \cdot q_{\vec{\eta}} \cdot \ln\left(\frac{q_{\vec{\eta}}}{p}\right) \geq 0$$

proposed approximation  $\leftarrow$   $\rightarrow$  true posterior

useless since it depends on  $p$

- Intractability of  $p$  typically due to its normalization ("Evidence")

$$p = \tilde{p} / E$$

$$ELBO = \int d\vec{\theta} \cdot q_{\vec{\eta}} \ln\left(\frac{\tilde{p}}{q_{\vec{\eta}}}\right)$$

Maximize wrt  $\vec{\eta}$



- Bonus:

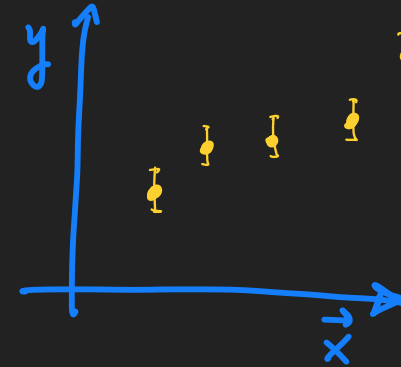
Get an approx. to  $E$ :  
 $\ln E \approx ELBO$

---

CASE # 1

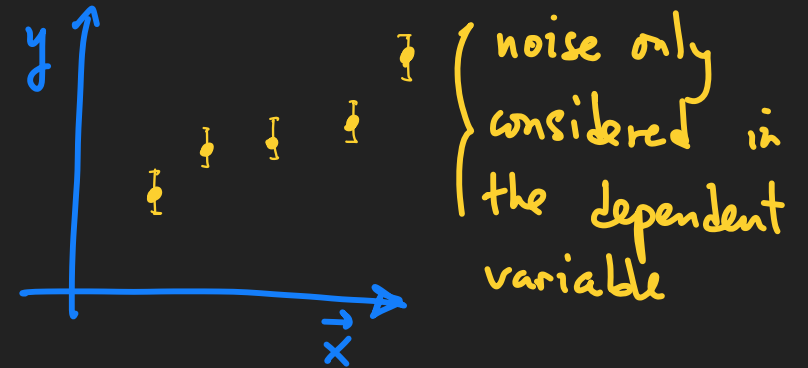


- In a typical data-science problem

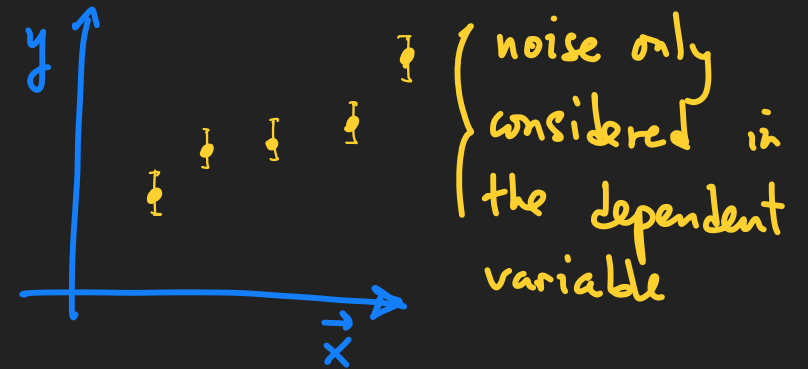




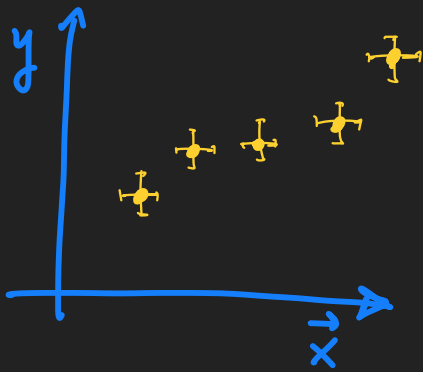
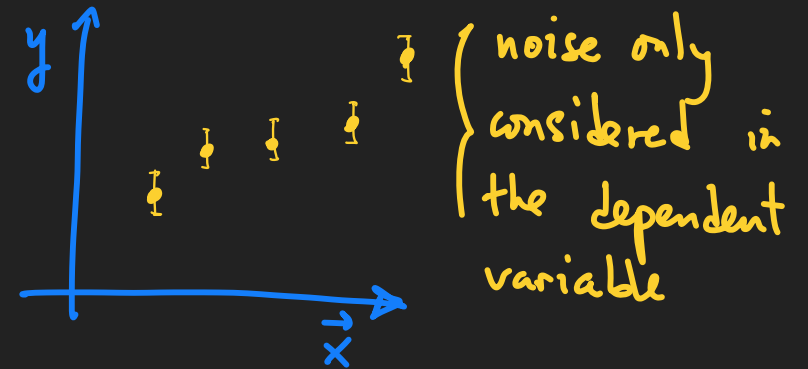
- In a typical data-science problem



- In a typical data-science problem
- In physics it is quite common to have uncertainties also in the  $\vec{x}$  (e.g. instrumental errors)

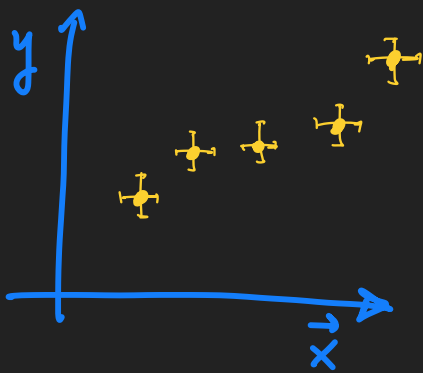
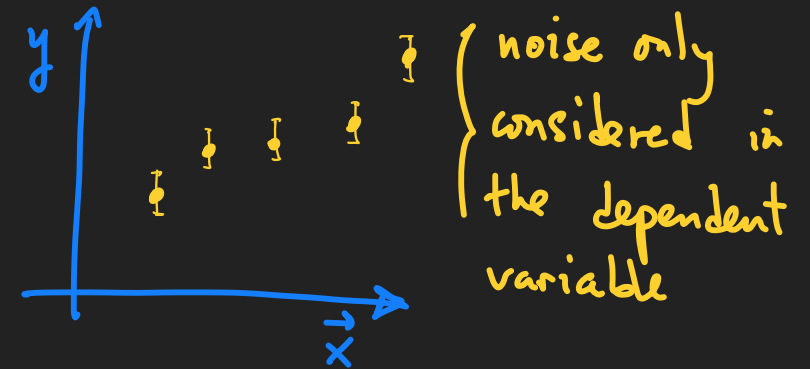


- In a typical data-science problem
- In physics it is quite common to have uncertainties also in the  $\vec{x}$  (e.g. instrumental errors)



What is the impact  
on the fitting function?

- In a typical data-science problem
- In physics it is quite common to have uncertainties also in the  $\vec{x}$  (e.g. instrumental errors)



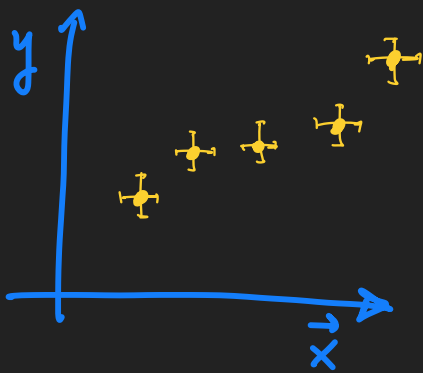
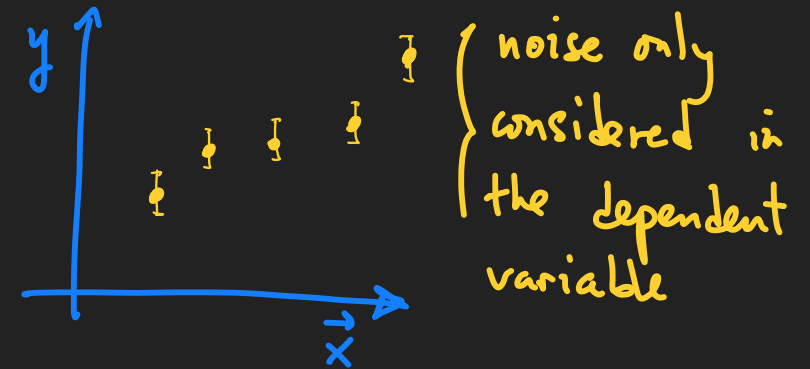
[Numerical Recipes]

e.g.  $f(x) = a + bx$

$$\chi^2(a, b) = \sum_{i=1}^N \frac{(y_i - f(x_i))^2}{\sigma_{y_i}^2 + b^2 \sigma_{x_i}^2}$$

What is the impact on the fitting function?

- In a typical data-science problem
- In physics it is quite common to have uncertainties also in the  $\vec{x}$  (e.g. instrumental errors)



[Numerical Recipes]

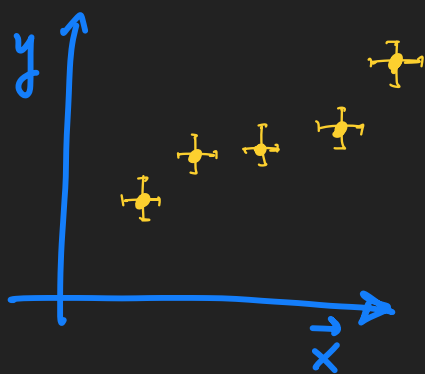
e.g.  $f(x) = a + bx$

$$\chi^2(a, b) = \sum_{i=1}^N \frac{(y_i - f(x_i))^2}{\sigma_{y_i}^2 + b^2 \sigma_{x_i}^2}$$

What is the impact on the fitting function?

+ Many refinements in the literature

- In a typical data-science problem
- In physics it is quite common to have uncertainties also in the  $\vec{x}$  (e.g. instrumental errors)



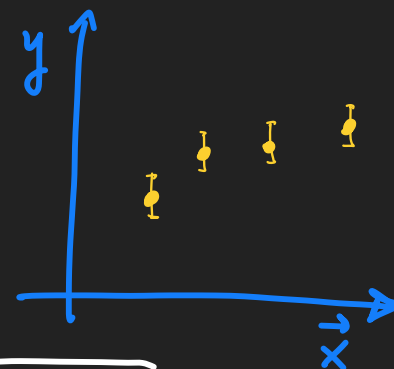
[Numerical Recipes]

e.g.  $f(x) = a + bx$

$$\chi^2(a, b) = \sum_{i=1}^N \frac{(y_i - f(x_i))^2}{\sigma_{y_i}^2 + b^2 \sigma_{x_i}^2}$$

+ Many refinements in the literature

What is the impact on the fitting function?

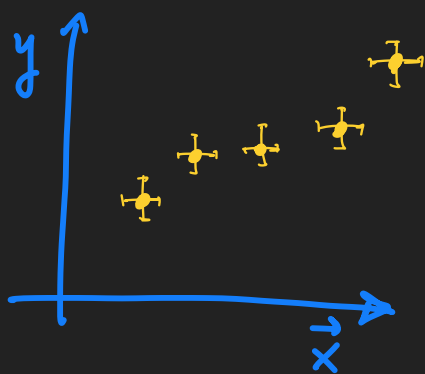


noise only considered in the dependent variable

What about classification problems?

(e.g. signal-vs-bckg, cats-vs-dogs-vs-donkeys)

- In a typical data-science problem
- In physics it is quite common to have uncertainties also in the  $\vec{x}$  (e.g. instrumental errors)



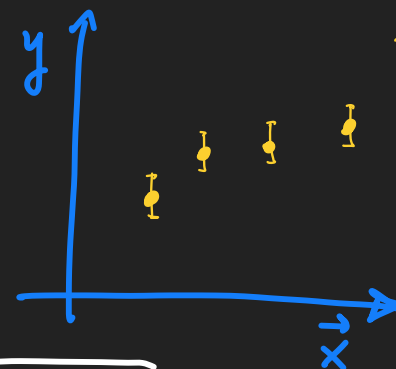
[Numerical Recipes]

e.g.  $f(x) = a + bx$

$$\chi^2(a, b) = \sum_{i=1}^N \frac{(y_i - f(x_i))^2}{\sigma_{y_i}^2 + b^2 \sigma_{x_i}^2}$$

What is the impact on the fitting function?

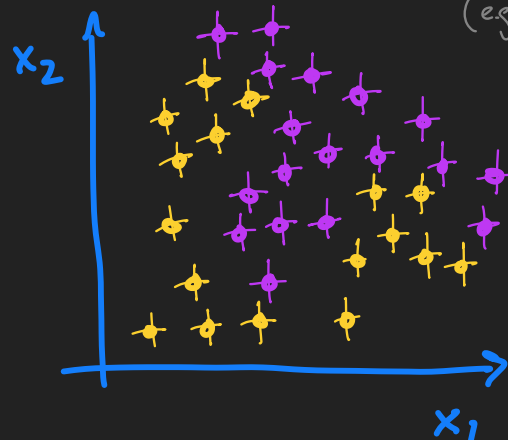
+ Many refinements in the literature



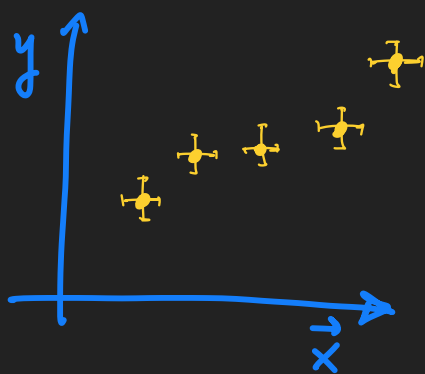
noise only considered in the dependent variable

What about classification problems?

(e.g. signal-vs-bckg, cats-vs-dogs-vs-donkeys)



- In a typical data-science problem
- In physics it is quite common to have uncertainties also in the  $\vec{x}$  (e.g. instrumental errors)



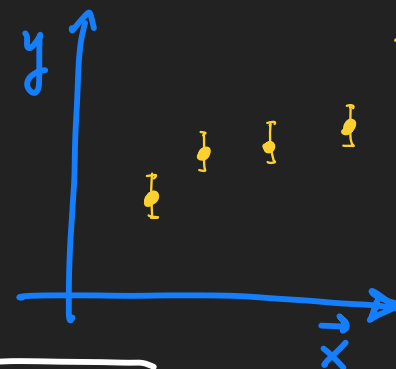
[Numerical Recipes]

e.g.  $f(x) = a + bx$

$$\chi^2(a,b) = \sum_{i=1}^N \frac{(y_i - f(x_i))^2}{\sigma_{y_i}^2 + b^2 \sigma_{x_i}^2}$$

+ Many refinements in the literature

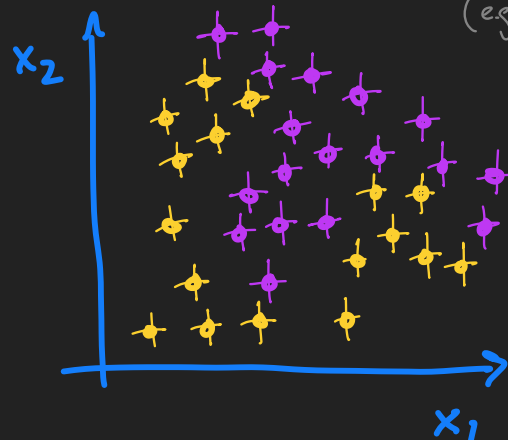
What is the impact on the fitting function?



noise only considered in the dependent variable

What about classification problems?

(e.g. signal-vs-bckg, cats-vs-dogs-vs-donkeys)



Not formally addressed before 2019



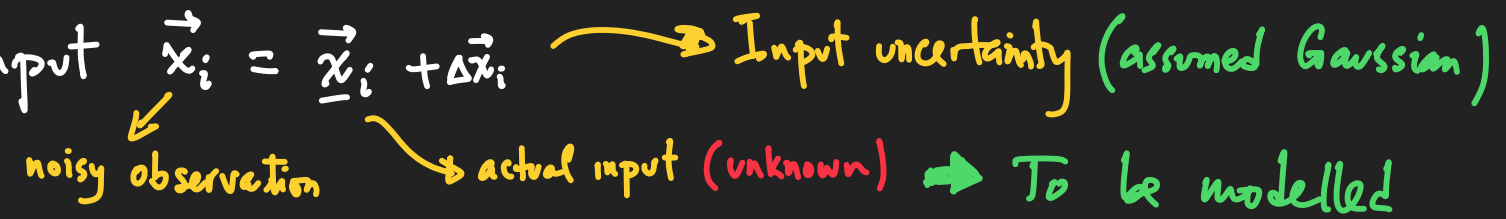


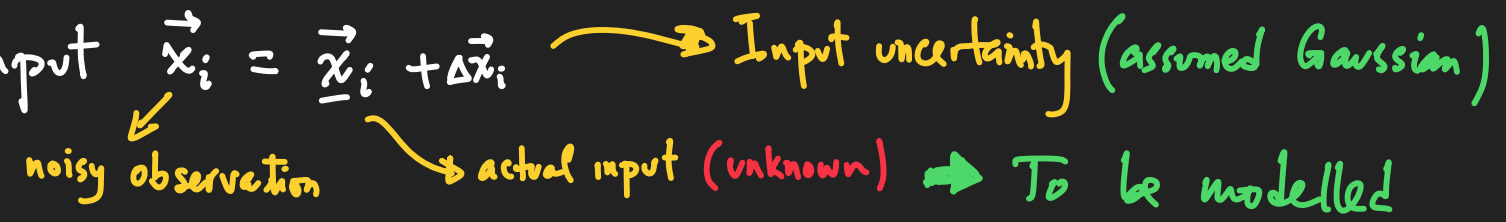
- Idea : Each input  $\vec{x}_i = \vec{\underline{x}}_i + \Delta\vec{x}_i$

- Idea : Each input  $\vec{x}_i = \vec{\bar{x}}_i + \Delta\vec{x}_i$   
noisy observation

- Idea : Each input  $\vec{x}_i = \underline{\vec{x}}_i + \Delta\vec{x}_i$   
noisy observation  $\swarrow$   $\searrow$  actual input (unknown)

- Idea : Each input  $\vec{x}_i = \underline{\vec{x}}_i + \Delta\vec{x}_i$   
noisy observation  $\swarrow$  actual input (unknown)  $\rightarrow$  To be modelled

- Idea : Each input  $\vec{x}_i = \vec{\bar{x}}_i + \Delta\vec{x}_i$     
noisy observation  $\rightarrow$  actual input (unknown)  $\rightarrow$  To be modelled   
Input uncertainty (assumed Gaussian)

- Idea : Each input  $\vec{x}_i = \vec{\bar{x}}_i + \Delta\vec{x}_i$    
noisy observation  $\rightarrow$  actual input (unknown)  $\rightarrow$  To be modelled  
(in analogy to what we typically do for the dependent variable  $y$ )

- Idea: Each input  $\vec{x}_i = \vec{\bar{x}}_i + \Delta\vec{x}_i$ 
  - $\vec{\bar{x}}_i$ : actual input (unknown)  $\rightarrow$  To be modelled
  - $\Delta\vec{x}_i$ : Input uncertainty (assumed Gaussian)
  - $\vec{x}_i$ : noisy observation

(in analogy to what we typically do for the dependent variable  $y$ )
- Likelihood of data  $p(Y, X | F, \underline{X}) = \prod_i^N p(y_i | f(\vec{x}_i)) \cdot \mathcal{N}(\vec{x}_i | \vec{\bar{x}}_i, \Delta\vec{x}_i)$ 
  - $\underline{X}$ : latent variables for the output  $Y$



- Idea: Each input  $\vec{x}_i = \vec{\bar{x}}_i + \Delta\vec{x}_i$    
↙ noisy observation ↘ actual input (unknown) → Input uncertainty (assumed Gaussian) → To be modelled

(in analogy to what we typically do for the dependent variable  $y$ )

- Likelihood of data  $p(Y, X | F, \underline{X}) = \prod_i^N p(y_i | f(\vec{x}_i)) \cdot \mathcal{N}(\vec{x}_i | \vec{\bar{x}}_i, \Delta\vec{x}_i)$    
↙ Latent variables for the output  $Y$

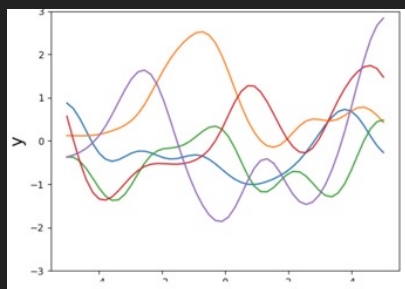
- $Y$  modelled as a Gaussian Process (GP)

- Idea: Each input  $\vec{x}_i = \vec{\bar{x}}_i + \Delta\vec{x}_i$    
noisy observation actual input (unknown) Input uncertainty (assumed Gaussian)   
To be modelled   
 (in analogy to what we typically do for the dependent variable  $y$ )
- Likelihood of data  $p(Y, X | F, \underline{X}) = \prod_i^N p(y_i | f(\vec{x}_i)) \cdot \mathcal{N}(\vec{x}_i | \vec{\bar{x}}_i, \Delta\vec{x}_i)$    
latent variables for the output  $Y$
- $Y$  modelled as a Gaussian Process (GP)  $\Rightarrow$  Very popular Stochastic Process in ML, based on Gaussian distrib. [over functions]

- Idea: Each input  $\vec{x}_i = \vec{\bar{x}}_i + \Delta\vec{x}_i$   $\Rightarrow$  Input uncertainty (assumed Gaussian)  
 $\swarrow$  noisy observation  $\searrow$  actual input (unknown)  $\Rightarrow$  To be modelled  
 (in analogy to what we typically do for the dependent variable  $y$ )

- Likelihood of data  $p(Y, X | F, \underline{X}) = \prod_i^N p(y_i | f(\vec{x}_i)) \cdot \mathcal{N}(\vec{x}_i | \vec{\bar{x}}_i, \Delta\vec{x}_i)$   
 $\swarrow$  latent variables for the output  $Y$

- $Y$  modelled as a Gaussian Process (GP)  $\Rightarrow$  Very popular Stochastic Process in ML, based on Gaussian distrib. [over functions]  
 GP prior

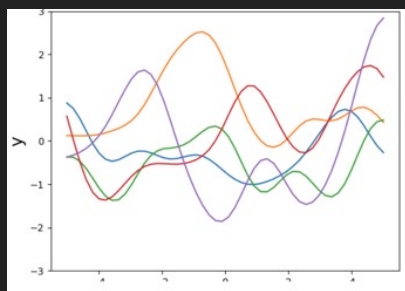


- Idea: Each input  $\vec{x}_i = \vec{\bar{x}}_i + \Delta\vec{x}_i$   $\rightarrow$  Input uncertainty (assumed Gaussian)  
 $\swarrow$  noisy observation  $\searrow$  actual input (unknown)  $\rightarrow$  To be modelled  
 (in analogy to what we typically do for the dependent variable  $y$ )

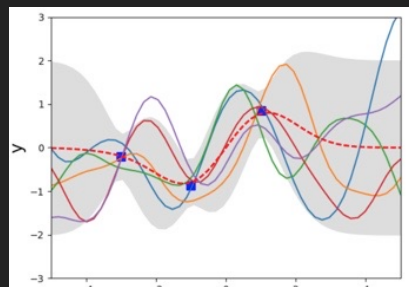
- Likelihood of data  $p(Y, X | F, \underline{X}) = \prod_i^N p(y_i | f(\vec{x}_i)) \cdot \mathcal{N}(\vec{x}_i | \vec{\bar{x}}_i, \Delta\vec{x}_i)$   
 $\swarrow$  Latent variables for the output  $Y$

- $Y$  modelled as a Gaussian Process (GP)  $\Rightarrow$  Very popular Stochastic Process in ML, based on Gaussian distrib. [over functions]

GP prior



GP posterior after 3 obs.

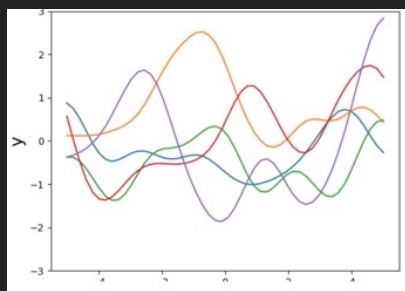


- Idea: Each input  $\vec{x}_i = \vec{\bar{x}}_i + \Delta\vec{x}_i$   $\rightarrow$  Input uncertainty (assumed Gaussian)  
 $\swarrow$  noisy observation  $\searrow$  actual input (unknown)  $\rightarrow$  To be modelled  
 (in analogy to what we typically do for the dependent variable  $y$ )

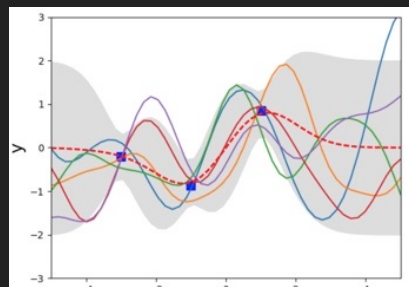
- Likelihood of data  $p(Y, X | F, \underline{X}) = \prod_i^N p(y_i | f(\vec{x}_i)) \cdot \mathcal{N}(\vec{x}_i | \vec{\bar{x}}_i, \Delta\vec{x}_i)$   
 $\swarrow$  latent variables for the output  $Y$

- $Y$  modelled as a Gaussian Process (GP)  $\Rightarrow$  Very popular Stochastic Process in ML, based on Gaussian distrib. [over functions]

GP prior



GP posterior after 3 obs.



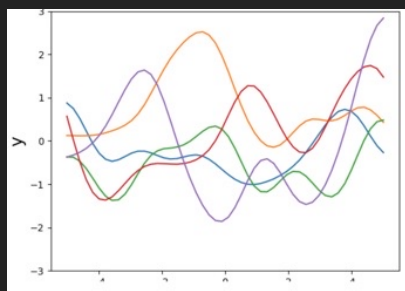
- GP give analytical predictions in regression problems

- Idea: Each input  $\vec{x}_i = \vec{\bar{x}}_i + \Delta\vec{x}_i$   $\rightarrow$  Input uncertainty (assumed Gaussian)  
 $\swarrow$  noisy observation  $\searrow$  actual input (unknown)  $\rightarrow$  To be modelled  
 (in analogy to what we typically do for the dependent variable  $y$ )

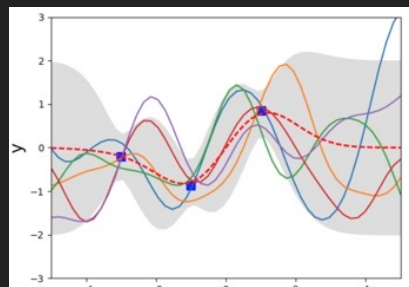
- Likelihood of data  $p(Y, X | F, \underline{X}) = \prod_i^N p(y_i | f(\vec{x}_i)) \cdot \mathcal{N}(\vec{x}_i | \vec{\bar{x}}_i, \Delta\vec{x}_i)$   
 $\swarrow$  latent variables for the output  $Y$

- $Y$  modelled as a Gaussian Process (GP)  $\Rightarrow$  Very popular Stochastic Process in ML, based on Gaussian distrib. [over functions]

GP prior



GP posterior after 3 obs.




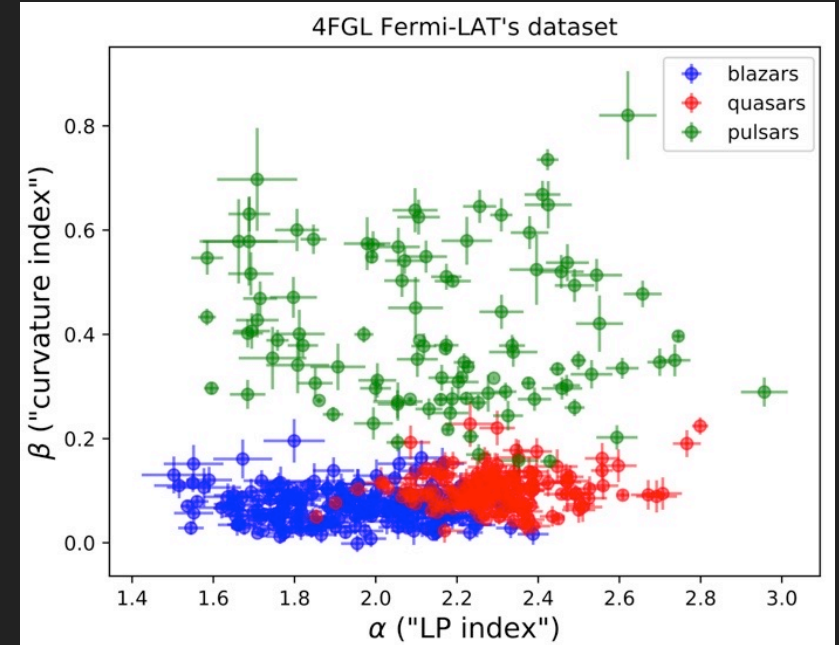
- GP give analytical predictions in regression problems
- For classification the posterior should be approximated  
 $\rightarrow$  nowadays typically using Variational Inference




- Tested the benefits of This approach in several synthetic and real datasets,

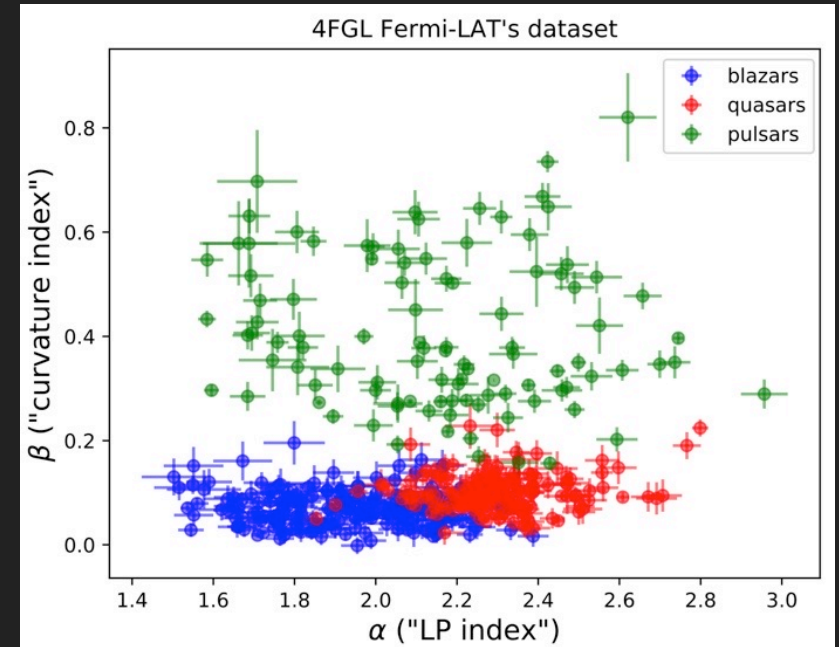



- Tested the benefits of This approach in several synthetic and real datasets, e.g. 



- Tested the benefits of This approach in several synthetic and real datasets, e.g. 

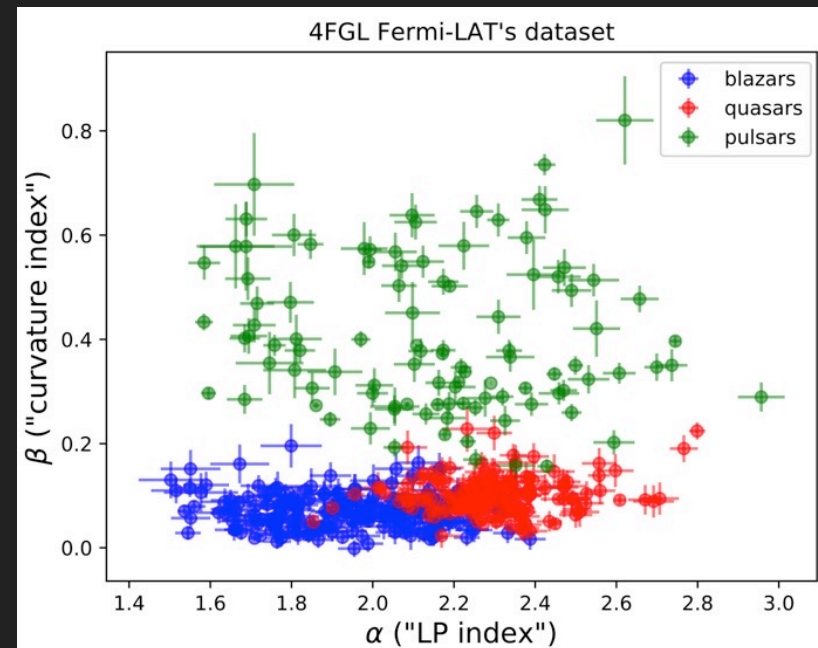
Classifying among 3 types of astro. sources from latest Fermi Catalog




- Tested the benefits of This approach in several synthetic and real datasets, e.g. 

Classifying among 3 types of astro. sources from latest Fermi Catalog

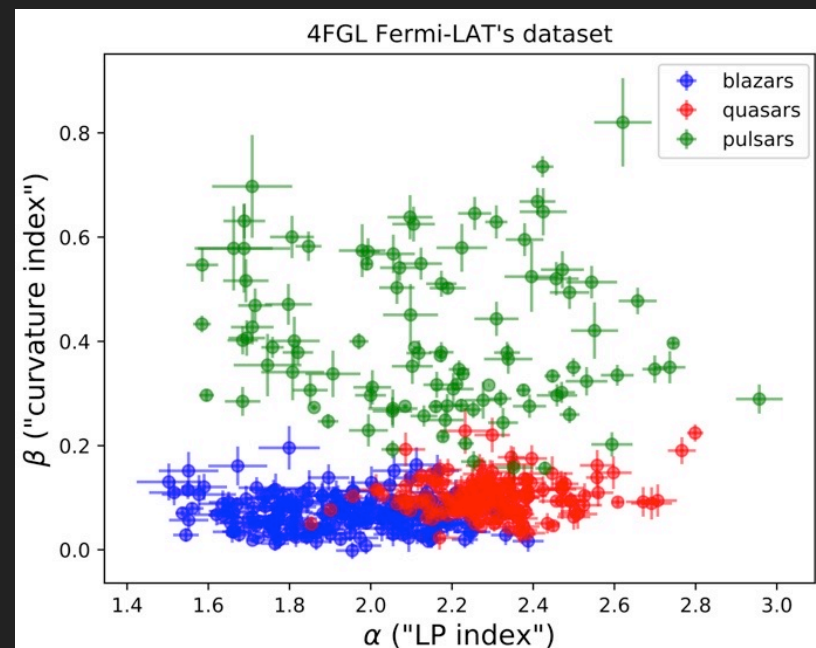
- Considering 6 representative attributes of the point-like sources



- Tested the benefits of This approach in several synthetic and real datasets, e.g. 


Classifying among 3 types of astro. sources from latest Fermi Catalog

- Considering 6 representative attributes of the point-like sources



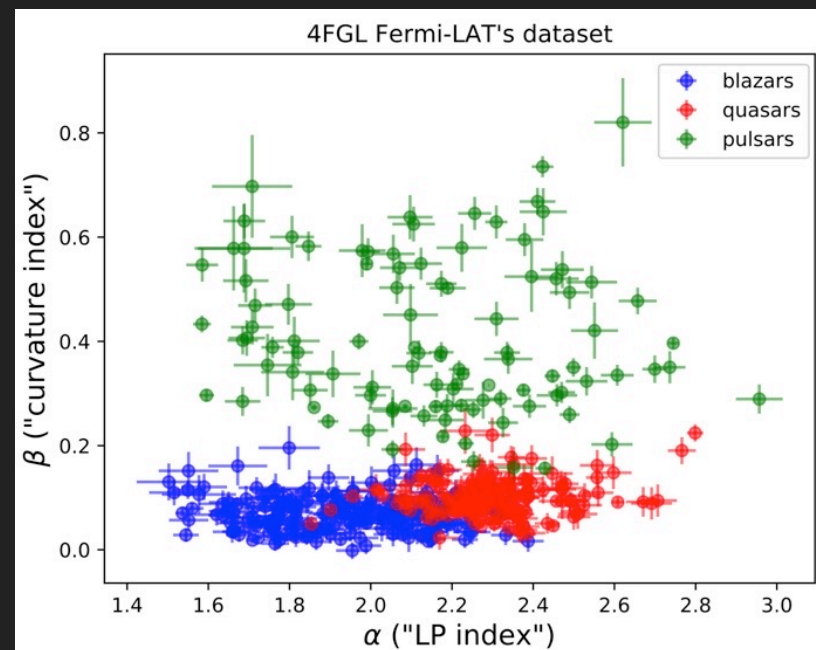
NOISELESS      VARIANT #1      VARIANT #2      VARIANT #3

	MGP	NIMGP	NIMGP <sub>NN</sub>	NIMGP <sub>FO</sub>
NLL	0.377±0.0194	<b>0.246±0.0097</b>	0.261±0.011	0.292±0.0158
Test error	0.075±0.0043	0.088±0.0038	0.082±0.0041	<b>0.071±0.0038</b>

- Tested the benefits of This approach in several synthetic and real datasets, e.g. 

Classifying among 3 types of astro. sources from latest Fermi Catalog

- Considering 6 representative attributes of the point-like sources



Take-home message:

- No clear benefit in Accuracy
- Clear advantage for predictive distribution!

NOISELESS      VARIANT #1      VARIANT #2      VARIANT #3

	MGP	NIMGP	NIMGP <sub>NN</sub>	NIMGP <sub>FO</sub>
NLL	0.377±0.0194	<b>0.246±0.0097</b>	0.261±0.011	0.292±0.0158
Test error	0.075±0.0043	0.088±0.0038	0.082±0.0041	<b>0.071±0.0038</b>

TEXT

---

CASE #2

SANTANA, BZ, HERNANDEZ, INTERNATIONALCONFERENCE ON ML 2022, (2110.07618)

---

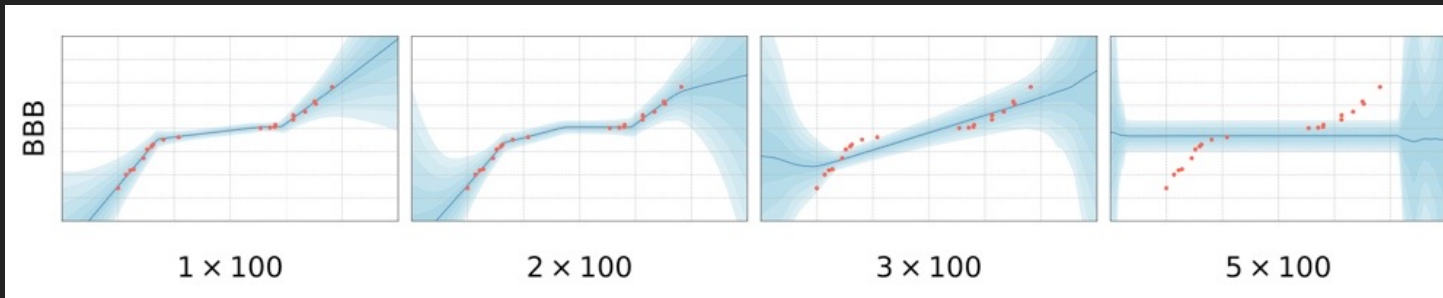
- What happens in typical Variational Inference for very complicated models? (e.g. deep neural networks)



- What happens in typical Variational Inference for very complicated models? (e.g. deep neural networks)  $\left\{ \begin{array}{l} \text{Difficult to keep} \\ \text{trade-off } \mathcal{F}(\vec{\theta}) \\ \text{Tractability} \\ \text{vs.} \\ \text{Expressivity} \end{array} \right.$

- What happens in typical Variational Inference for very complicated models? (e.g. deep neural networks)

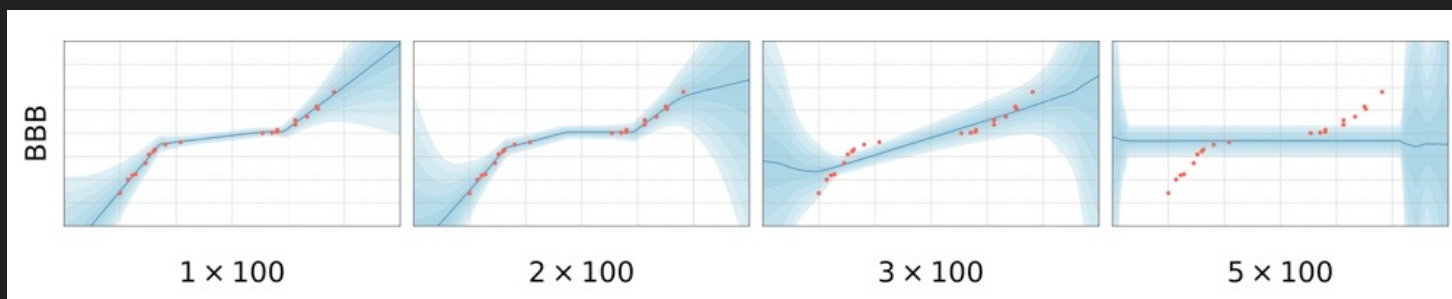
} Difficult to keep  
trade-off  $\mathcal{F}(\vec{\theta})$   
Tractability  
vs.  
Expressivity



1903.05779

- What happens in typical Variational Inference for very complicated models? (e.g. deep neural networks)

} Difficult to keep  
trade-off  $\mathcal{F}(\vec{\theta})$   
Tractability  
vs.  
Expressivity

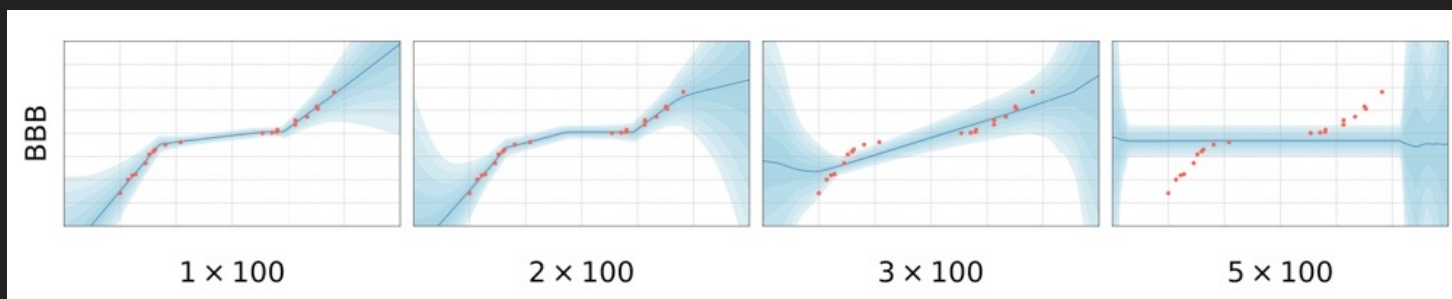


1903.05779

- ML community now considering doing inference in "function space" instead of parameter space

- What happens in typical Variational Inference for very complicated models? (e.g. deep neural networks)

} Difficult to keep  
trade-off  $\mathcal{F}(\vec{\theta})$   
Tractability  
vs.  
Expressivity

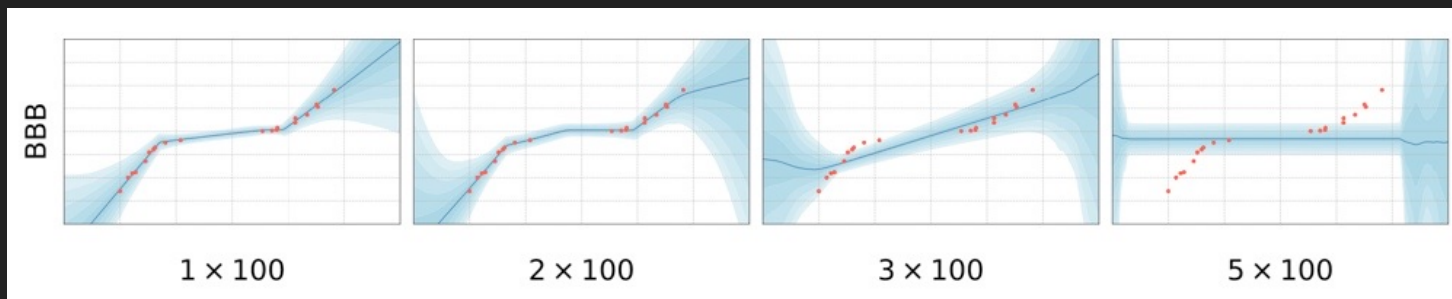


1903.05779

- ML community now considering doing inference in "function space" instead of parameter space
  - Priors over functions (GP, Implicit Processes)
- $\{f(x_1), f(x_2), \dots, f(x_n)\}$

- What happens in typical Variational Inference for very complicated models? (e.g. deep neural networks)

} Difficult to keep  
trade-off  $\mathcal{F}(\vec{\theta})$   
Tractability  
vs.  
Expressivity



1903.05779

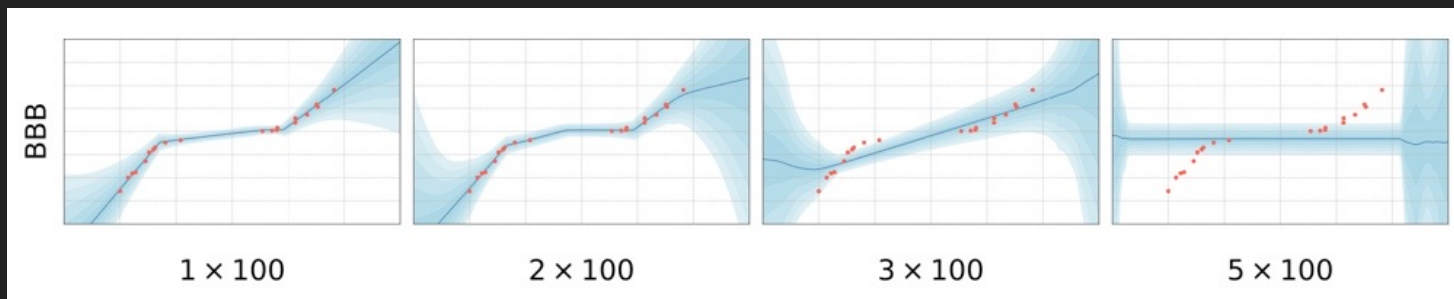
- ML community now considering doing inference in "function space" instead of parameter space

- Priors over functions (GP, Implicit Processes)

$\{f(x_1), f(x_2), \dots, f(x_n)\}$

IP defined by "Implicit distributions"  
[cannot be evaluated, but can be sampled]

- What happens in typical Variational Inference for very complicated models? (e.g. deep neural networks)
- Difficult to keep trade-off  $\mathcal{F}(\vec{\theta})$
- Tractability vs. Expressivity



1903.05779

- ML community now considering doing inference in "function space" instead of parameter space
  - Priors over functions (GP, Implicit Processes)
- $\{f(x_1), f(x_2), \dots, f(x_n)\}$

IP defined by "Implicit distributions"  
[cannot be evaluated, but can be sampled]

e.g. Physics simulators:

$$f(x) = g_{\theta}(x, z), \quad z \sim p(z)$$

$\theta$  encodes physics       $z$  remaining randomness (Latent variables)

SANTANA, BZ, HERNANDEZ, INTERNATIONALCONFERENCE ON ML 2022, (2110.07618)

---

- V.I. procedures in function space are computationally challenging  
ELBO  $[q | \tilde{p}]$  now for two implicit distributions over functions!



- V.I. procedures in function space are computationally challenging  
 $\text{ELBO}[q|\tilde{p}]$  now for two implicit distributions over functions!

IP priors given by  
neural networks with  
 $(x, z)$  inputs  
 $\uparrow$   
 $\hookrightarrow$  stochastic

- V.I. procedures in function space are computationally challenging  
 $\text{ELBO}[q|\tilde{p}]$  now for two implicit distributions over functions!
- Useful intuition from standard VI for Gaussian Processes

IP priors given by  
neural networks with  
 $(x, z)$  inputs  
 $\uparrow$   
 $\rightarrow$  Stochastic

- V.I. procedures in function space are computationally challenging  
 $\text{ELBO}[q | \tilde{p}]$  now for two implicit distributions over functions!
- Useful intuition from standard VI for Gaussian Processes
- $\text{ELBO} \supset \mathbb{E}_q \left[ \ln \frac{q}{\text{prior}} \right] = ?$

IP priors given by  
neural networks with  
 $(x, z)$  inputs  
 $\uparrow$   
Stochastic

- V.I. procedures in function space are computationally challenging  
 $\text{ELBO}[q | \tilde{p}]$  now for two implicit distributions over functions!

IP priors given by  
 neural networks with  
 $(x, z)$  inputs  
 $\uparrow$  stochastic

- Useful intuition from standard VI for Gaussian Processes

- $\text{ELBO} \supset \mathbb{E}_q \left[ \ln \frac{q}{\text{prior}} \right] = ?$   $\ln \frac{q}{\text{prior}} \approx \mathcal{D}_{w*}$  } a "discriminator" (e.g. neural net)  
 trained to distinguish samples  
 from  $q$  and prior

- V.I. procedures in function space are computationally challenging  
 $\text{ELBO}[q | \tilde{p}]$  now for two implicit distributions over functions!

IP priors given by  
 neural networks with  
 $(x, z)$  inputs  
 $\uparrow$  stochastic

- Useful intuition from standard VI for Gaussian Processes

- $\text{ELBO} \supset \mathbb{E}_q \left[ \ln \frac{q}{\text{prior}} \right] = ?$   $\ln \frac{q}{\text{prior}} \approx \mathcal{D}_{w*}$  } a "discriminator" (e.g. neural net)  
 trained to distinguish samples  
 from  $q$  and prior

• Results

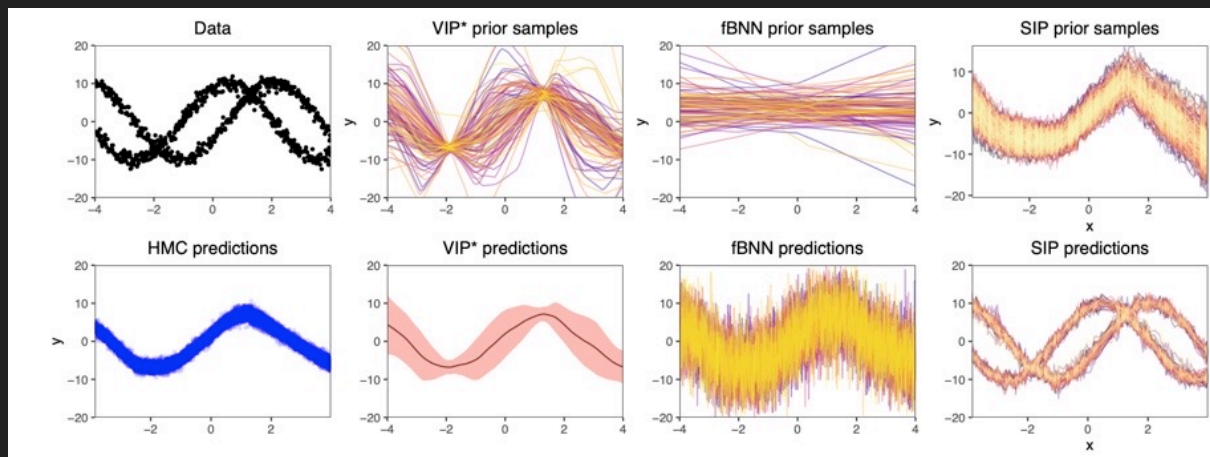
- V.I. procedures in function space are computationally challenging  
 $\text{ELBO}[q | \tilde{p}]$  now for two implicit distributions over functions!

IP priors given by  
 neural networks with  
 $(x, z)$  inputs  
 $\uparrow$  stochastic

- Useful intuition from standard VI for Gaussian Processes

- $\text{ELBO} \supset \mathbb{E}_q \left[ \ln \frac{q}{\text{prior}} \right] = ?$   $\ln \frac{q}{\text{prior}} \approx \mathcal{D}_{w*}$  } a "discriminator" (e.g. neural net)  
 trained to distinguish samples  
 from  $q$  and prior

Results



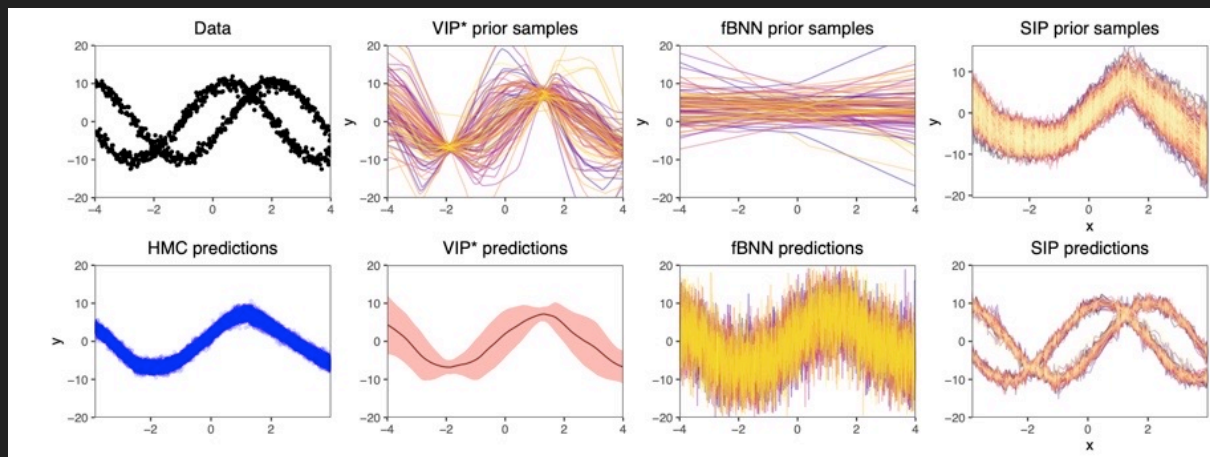
- V.I. procedures in function space are computationally challenging  
 $\text{ELBO}[q | \tilde{p}]$  now for two implicit distributions over functions!

IP priors given by  
 neural networks with  
 $(x, z)$  inputs  
 $\uparrow$  stochastic

- Useful intuition from standard VI for Gaussian Processes

- $\text{ELBO} \supset \mathbb{E}_q \left[ \ln \frac{q}{\text{prior}} \right] = ?$   $\ln \frac{q}{\text{prior}} \approx \mathcal{D}_{w*}$  } a "discriminator" (e.g. neural net)  
 trained to distinguish samples  
 from  $q$  and prior

Results



"Golden truth"

Existing method #1

Existing method #2

Our method

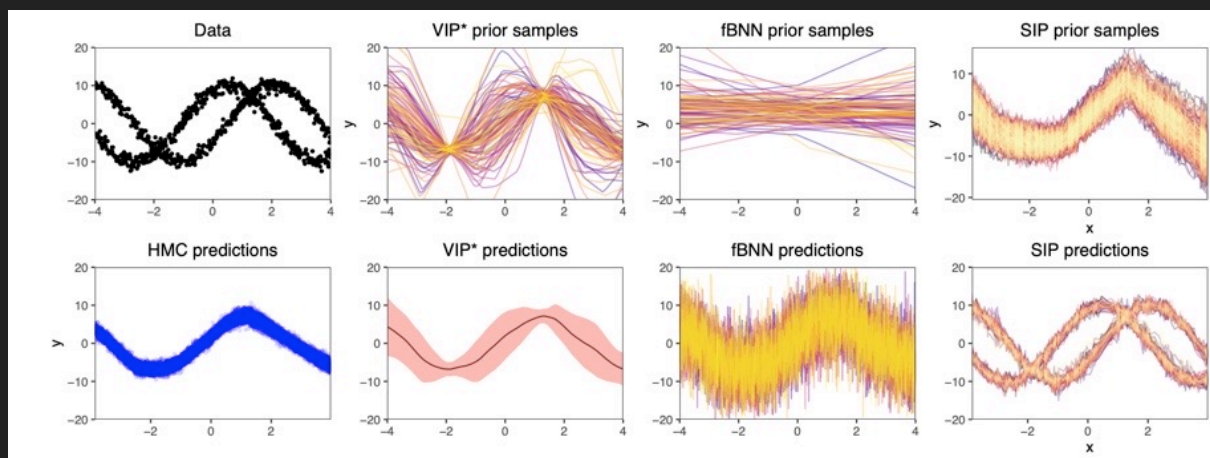
- V.I. procedures in function space are computationally challenging  
 $\text{ELBO}[q | \tilde{p}]$  now for two implicit distributions over functions!

IP priors given by  
 neural networks with  
 $(x, z)$  inputs  
 $\uparrow$  Stochastic

- Useful intuition from standard VI for Gaussian Processes

- $\text{ELBO} \supset \mathbb{E}_q \left[ \ln \frac{q}{\text{prior}} \right] = ?$   $\ln \frac{q}{\text{prior}} \approx \mathcal{D}_{w*}$  } a "discriminator" (e.g. neural net)  
 trained to distinguish samples  
 from  $q$  and prior

Results



"Golden truth"

Existing method #1

Existing method #2

Our method

+ Our method performs  
 better in general for a  
 collection of other  
 real datasets



---

CASE #3

## OTHER ACTIVE RESEARCH BRANCH IN VARIATIONAL INFERENCE

---

## OTHER ACTIVE RESEARCH BRANCH IN VARIATIONAL INFERENCE

---

Consider the following mathematical identity ("reweighting" - Lattece QCD jargon)

$$\mathbb{E}_p[f(\theta)] = \int d\theta p(\theta) f(\theta) = \int d\theta q(\theta) \frac{p(\theta)}{q(\theta)} f(\theta) = \mathbb{E}_q \left[ \frac{p}{q} f(\theta) \right]$$

↘ "weights"

## OTHER ACTIVE RESEARCH BRANCH IN VARIATIONAL INFERENCE

---

Consider the following mathematical identity ("reweighting" - Lattece QCD jargon)

$$\mathbb{E}_p[f(\theta)] = \int d\theta p(\theta) f(\theta) = \int d\theta q(\theta) \frac{p(\theta)}{q(\theta)} f(\theta) = \mathbb{E}_q \left[ \frac{p}{q} f(\theta) \right]$$

↘ "weights"

- This is the basis of "Importance Sampling"

## OTHER ACTIVE RESEARCH BRANCH IN VARIATIONAL INFERENCE

---

Consider the following mathematical identity ("reweighting" - Latent QCD jargon)

$$\mathbb{E}_p[f(\theta)] = \int d\theta p(\theta) f(\theta) = \int d\theta q(\theta) \frac{p(\theta)}{q(\theta)} f(\theta) = \mathbb{E}_q \left[ \frac{p}{q} f(\theta) \right]$$

↗ "weights"

- This is the basis of "Importance Sampling"  $\Rightarrow$  computing  $\mathbb{E}$  under complicated  $p$ .

## OTHER ACTIVE RESEARCH BRANCH IN VARIATIONAL INFERENCE

Consider the following mathematical identity ("reweighting" - Latent QCD jargon)

$$\mathbb{E}_p[f(\theta)] = \int d\theta p(\theta) f(\theta) = \int d\theta q(\theta) \frac{p(\theta)}{q(\theta)} f(\theta) = \mathbb{E}_q \left[ \frac{p}{q} f(\theta) \right]$$

↘ "weights"

- This is the basis of "Importance Sampling"  $\Rightarrow$  computing  $\mathbb{E}$  under complicated  $p$ .
- Useful for V.I.  $\Rightarrow$

## OTHER ACTIVE RESEARCH BRANCH IN VARIATIONAL INFERENCE

Consider the following mathematical identity ("reweighting" - Latte QCD jargon)

$$\mathbb{E}_p[f(\theta)] = \int d\theta p(\theta) f(\theta) = \int d\theta q(\theta) \frac{p(\theta)}{q(\theta)} f(\theta) = \mathbb{E}_q \left[ \underbrace{\frac{p}{q}}_{\text{"weights"}} f(\theta) \right]$$

- This is the basis of "Importance Sampling"  $\Rightarrow$  computing  $\mathbb{E}$  under complicated  $p$ .
- Useful for V.I.  $\Rightarrow$   $\text{posterior} \approx q(\theta)$  s.t.  $p(y_* | x_*) = \mathbb{E}_p[p(y_* | x_*, \theta)] \approx \mathbb{E}_q[ \quad ]$

## OTHER ACTIVE RESEARCH BRANCH IN VARIATIONAL INFERENCE

Consider the following mathematical identity ("reweighting" - Latent QCD jargon)

$$\mathbb{E}_p[f(\theta)] = \int d\theta p(\theta) f(\theta) = \int d\theta q(\theta) \frac{p(\theta)}{q(\theta)} f(\theta) = \mathbb{E}_q \left[ \underbrace{\frac{p}{q}}_{\text{"weights"}} f(\theta) \right]$$

- This is the basis of "Importance Sampling"  $\Rightarrow$  computing  $\mathbb{E}$  under complicated  $p$ .
- Useful for V.I.  $\Rightarrow$  posterior  $\approx q(\theta)$  s.t.  $p(y_* | x_*) = \mathbb{E}_p[p(y_* | x_*, \theta)] \approx \mathbb{E}_q[ \quad ]$

V.I. is inherently biased



## OTHER ACTIVE RESEARCH BRANCH IN VARIATIONAL INFERENCE

Consider the following mathematical identity ("reweighting" - Latte QCD jargon)

$$\mathbb{E}_p[f(\theta)] = \int d\theta p(\theta) f(\theta) = \int d\theta q(\theta) \frac{p(\theta)}{q(\theta)} f(\theta) = \mathbb{E}_q \left[ \underbrace{\frac{p}{q}}_{\text{"weights"}} f(\theta) \right]$$

- This is the basis of "Importance Sampling"  $\Rightarrow$  computing  $\mathbb{E}$  under complicated  $p$ .
- Useful for V.I.  $\Rightarrow$   $\text{posterior} \approx q(\theta)$  s.t.  $p(y_* | x_*) = \mathbb{E}_p[p(y_* | x_*, \theta)] \approx \mathbb{E}_q[ \quad ]$

V.I. is inherently biased

$\Rightarrow$  Reweighting helps "de-biasing" V.I

(see Jordan et al, 2106.15980  
and refs. therein)

## OTHER ACTIVE RESEARCH BRANCH IN VARIATIONAL INFERENCE

Consider the following mathematical identity ("reweighting" - Lattice QCD jargon)

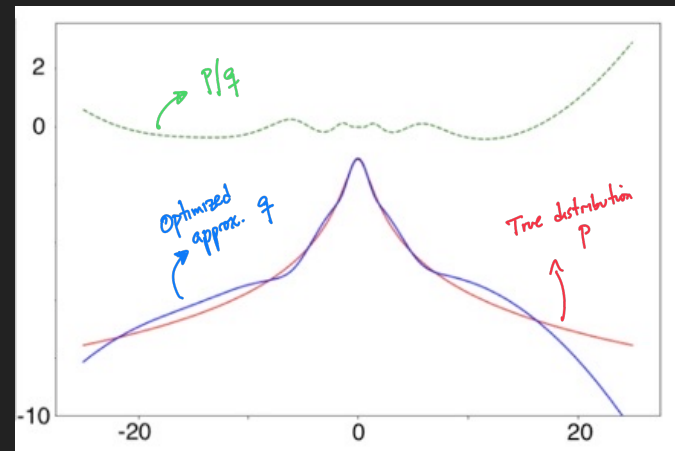
$$\mathbb{E}_p[f(\theta)] = \int d\theta p(\theta) f(\theta) = \int d\theta q(\theta) \frac{p(\theta)}{q(\theta)} f(\theta) = \mathbb{E}_q \left[ \underbrace{\frac{p}{q}}_{\text{"weights"}} f(\theta) \right]$$

- This is the basis of "Importance Sampling"  $\Rightarrow$  computing  $\mathbb{E}$  under complicated  $p$ .
- Useful for V.I.  $\Rightarrow$   $\text{posterior} \approx q(\theta)$  s.t.  $p(y_* | x_*) = \mathbb{E}_p[p(y_* | x_*, \theta)] \approx \mathbb{E}_q[ \quad ]$

V.I. is inherently biased

$\Rightarrow$  Reweighting helps "de-biasing" V.I

(see Jordan et al, 2106.15980 and refs. therein)



## OTHER ACTIVE RESEARCH BRANCH IN VARIATIONAL INFERENCE

---

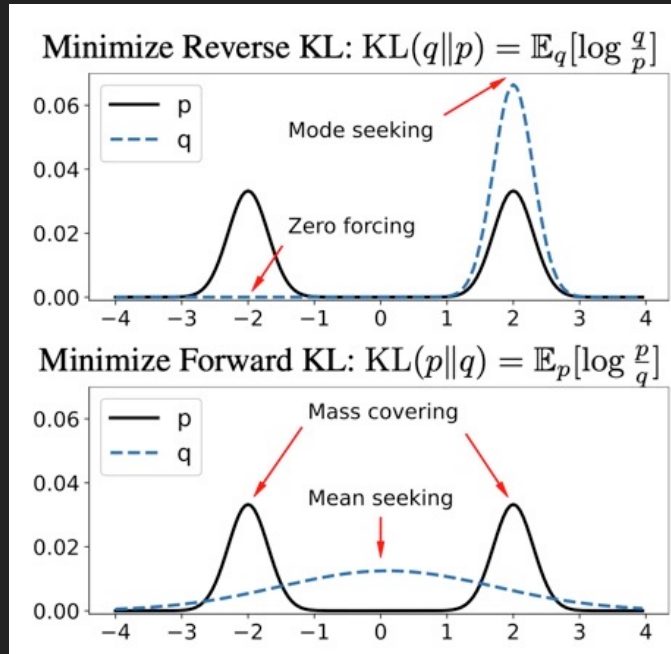
## OTHER ACTIVE RESEARCH BRANCH IN VARIATIONAL INFERENCE

---

- Challenges here :

## OTHER ACTIVE RESEARCH BRANCH IN VARIATIONAL INFERENCE

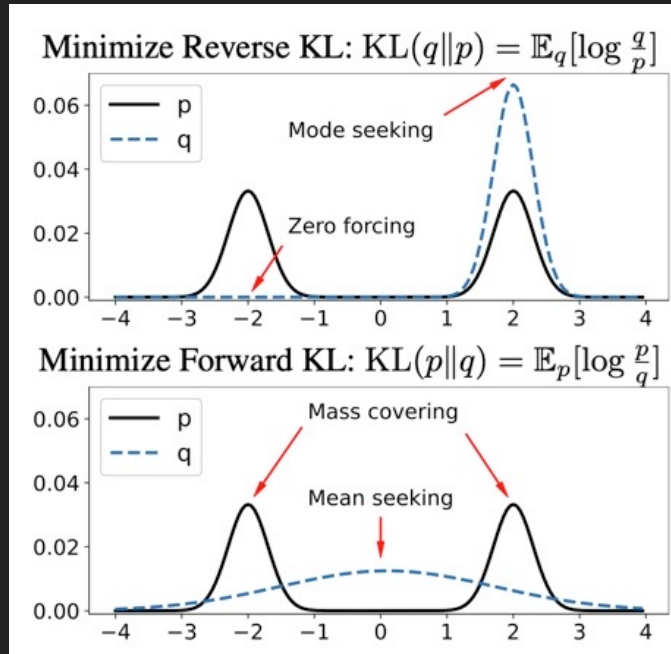
- Challenges here :



## OTHER ACTIVE RESEARCH BRANCH IN VARIATIONAL INFERENCE

- Challenges here :

weights have  
typically very  
large variance

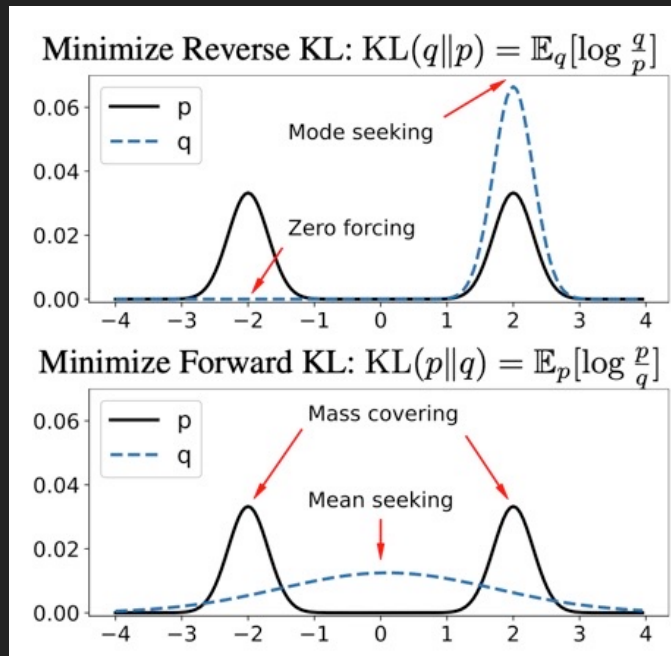


## OTHER ACTIVE RESEARCH BRANCH IN VARIATIONAL INFERENCE

- Challenges here :

weights have  
typically very  
large variance

smaller  
variance, but  
computationally  
very costly

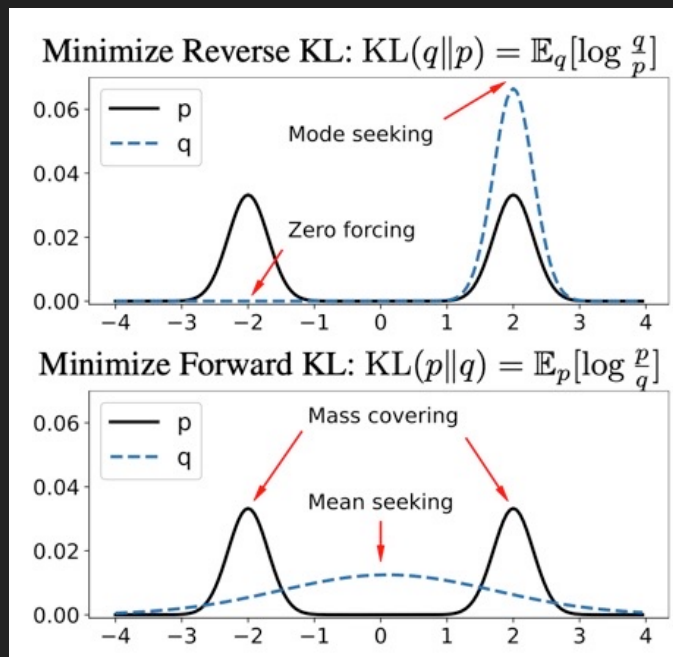


## OTHER ACTIVE RESEARCH BRANCH IN VARIATIONAL INFERENCE

- Challenges here:

weights have  
typically very  
large variance

smaller  
variance, but  
computationally  
very costly



\* Computation - Quality trade-off

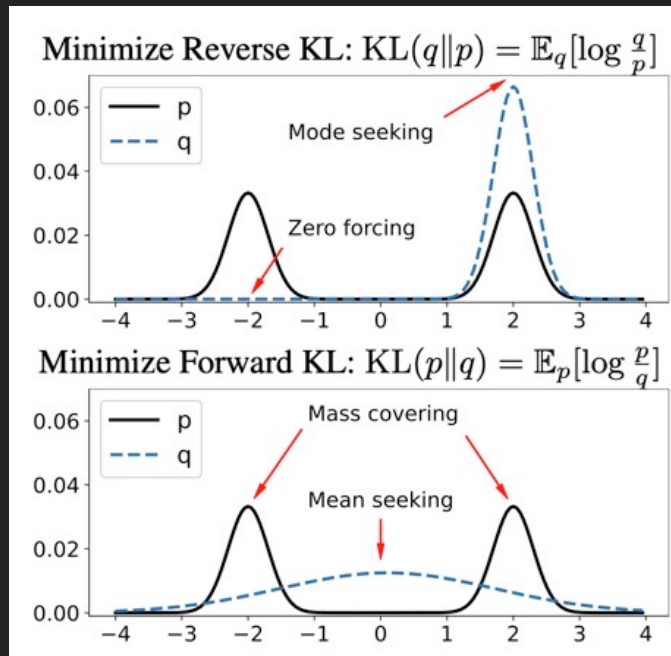


## OTHER ACTIVE RESEARCH BRANCH IN VARIATIONAL INFERENCE

- Challenges here:

weights have typically very large variance

smaller variance, but computationally very costly



\* Computation-Quality trade-off

\* Other "distance measures" are definitely worth exploring!

[A. Dimitriou, A. Ramos, G. Teb and BZ work in progress]

---

CASE # 4

## MARKOV CHAIN MONTE CARLO

---

## MARKOV CHAIN MONTE CARLO

---

- Idea of MCMC

### • Idea of MCMC

\* Get samples  $\vec{\theta}_1, \vec{\theta}_2, \dots, \vec{\theta}_t, \vec{\theta}_{t+1}, \dots$   
from proposal dist.  $q(\vec{\theta}_{t+1} | \vec{\theta}_t)$

### • Idea of MCMC

\* Get samples  $\vec{\theta}_1, \vec{\theta}_2, \dots, \vec{\theta}_t, \vec{\theta}_{t+1}, \dots$  from proposal dist.  $q(\vec{\theta}_{t+1} | \vec{\theta}_t)$  } Markov chain

### • Idea of MCMC

\* Get samples  $\vec{\theta}_1, \vec{\theta}_2, \dots, \vec{\theta}_t, \vec{\theta}_{t+1}, \dots$  from proposal dist.  $q(\vec{\theta}_{t+1} | \vec{\theta}_t)$  { Markov chain

\* Accept/Reject  $\vec{\theta}_{t+1}$  according to some criterion  
(typically ensuring good properties of the Markov chain)

## • Idea of MCMC

\* Get samples  $\vec{\theta}_1, \vec{\theta}_2, \dots, \vec{\theta}_t, \vec{\theta}_{t+1}, \dots$  from proposal dist.  $q(\vec{\theta}_{t+1} | \vec{\theta}_t)$  { Markov chain

\* Accept/Reject  $\vec{\theta}_{t+1}$  according to some criterion  
(typically ensuring good properties of the Markov chain)

$$p(\vec{\theta} | \text{Data}) = \frac{f(\vec{\theta})}{p(\text{Data})}$$



Idea of MCMC

\* Get samples  $\vec{\theta}_1, \vec{\theta}_2, \dots, \vec{\theta}_t, \vec{\theta}_{t+1}, \dots$  from proposal dist.  $q(\vec{\theta}_{t+1} | \vec{\theta}_t)$  } Markov chain

\* Accept/Reject  $\vec{\theta}_{t+1}$  according to some criterion  
(typically ensuring good properties of the Markov chain)

$$p(\vec{\theta} | \text{Data}) \propto \frac{f(\vec{\theta})}{p(\text{Data})} \quad \left| \quad \begin{array}{l} r = \frac{f(\vec{\theta}_{t+1})}{f(\vec{\theta}_t)} \\ u \sim \text{unif}[0,1] \end{array} \right.$$

## • Idea of MCMC

\* Get samples  $\vec{\theta}_1, \vec{\theta}_2, \dots, \vec{\theta}_t, \vec{\theta}_{t+1}, \dots$  from proposal dist.  $q(\vec{\theta}_{t+1} | \vec{\theta}_t)$  } Markov chain

\* Accept/Reject  $\vec{\theta}_{t+1}$  according to some criterion  
(typically ensuring good properties of the Markov chain)

$$p(\vec{\theta} | \text{Data}) \propto \frac{f(\vec{\theta})}{p(\text{Data})} \left| \begin{array}{l} r = \frac{f(\vec{\theta}_{t+1})}{f(\vec{\theta}_t)} \\ u \sim \text{unif}[0,1] \end{array} \right. \left. \begin{array}{l} u \leq r \Rightarrow \text{Accept } \vec{\theta}_{t+1} \\ u > r \Rightarrow \text{Reject it} \end{array} \right.$$

## HYBRID (A.K.A. HAMILTON) MONTE CARLO

---

## HYBRID (A.K.A. HAMILTON) MONTE CARLO

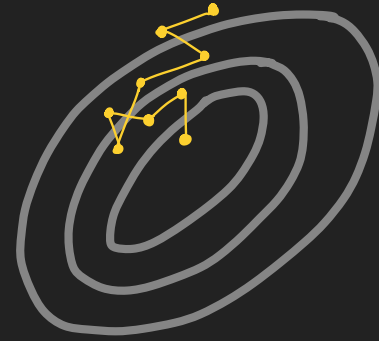
---

- Random walk behaviour  
(e.g. Metropolis-Hastings)

## HYBRID (A.K.A. HAMILTON) MONTE CARLO

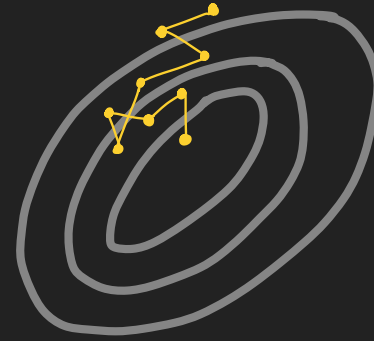
---

- Random walk behaviour  
(e.g. Metropolis-Hastings)



## HYBRID (A.K.A. HAMILTON) MONTE CARLO

- Random walk behaviour  
(e.g. Metropolis-Hastings)



Would be  
convenient to have  
guidance on where  
to move next!

## HYBRID (A.K.A. HAMILTON) MONTE CARLO

- Random walk behaviour  
(e.g. Metropolis-Hastings)
- Hamiltonian dynamics at work



Would be  
convenient to have  
guidance on where  
to move next!

## HYBRID (A.K.A. HAMILTON) MONTE CARLO

- Random walk behaviour  
(e.g. Metropolis-Hastings)
- Hamiltonian dynamics at work  
\*  $\vec{\theta} \equiv \vec{q} \rightarrow (\vec{q}, \vec{p})$

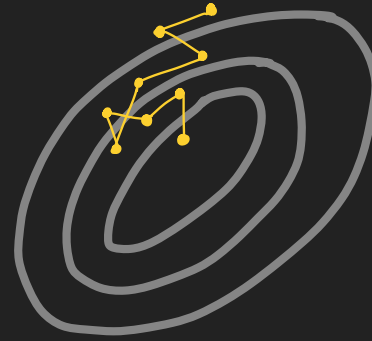


Would be  
convenient to have  
guidance on where  
to move next!



## HYBRID (A.K.A. HAMILTON) MONTE CARLO

- Random walk behaviour  
(e.g. Metropolis-Hastings)
- Hamiltonian dynamics at work
  - \*  $\vec{\theta} \equiv \vec{q} \rightarrow (\vec{q}, \vec{p})$
  - \*  $H(\vec{q}, \vec{p}) = K(\vec{p}) + V(\vec{q})$



Would be convenient to have guidance on where to move next!

## HYBRID (A.K.A. HAMILTON) MONTE CARLO

- Random walk behaviour  
(e.g. Metropolis-Hastings)



Would be convenient to have guidance on where to move next!

- Hamiltonian dynamics at work

$$* \vec{\theta} \equiv \vec{q} \rightarrow (\vec{q}, \vec{p}) \quad \rightarrow = \ln p(\vec{q}, \text{Data})$$

$$* H(\vec{q}, \vec{p}) = K(\vec{p}) + V(\vec{q})$$

## HYBRID (A.K.A. HAMILTON) MONTE CARLO

- Random walk behaviour  
(e.g. Metropolis-Hastings)



Would be convenient to have guidance on where to move next!

- Hamiltonian dynamics at work

$$* \vec{\theta} \equiv \vec{q} \rightarrow (\vec{q}, \vec{p}) \quad \underbrace{\hspace{1cm}}_{= \ln p(\vec{q}, \text{Data})}$$

$$* H(\vec{q}, \vec{p}) = K(\vec{p}) + V(\vec{q})$$

$$\frac{d\vec{q}}{dt} = \frac{\partial H}{\partial \vec{p}}, \quad \frac{d\vec{p}}{dt} = -\frac{\partial H}{\partial \vec{q}}$$

## HYBRID (A.K.A. HAMILTON) MONTE CARLO

- Random walk behaviour  
(e.g. Metropolis-Hastings)



Would be convenient to have guidance on where to move next!

- Hamiltonian dynamics at work

$$* \vec{\theta} \equiv \vec{q} \rightarrow (\vec{q}, \vec{p}) \quad \underbrace{\hspace{1cm}}_{= \ln p(\vec{q}, \text{Data})}$$

$$* H(\vec{q}, \vec{p}) = K(\vec{p}) + V(\vec{q})$$

$$\left. \begin{aligned} \frac{d\vec{q}}{dt} &= \frac{\partial H}{\partial \vec{p}} \\ \frac{d\vec{p}}{dt} &= -\frac{\partial H}{\partial \vec{q}} \end{aligned} \right\} \begin{array}{l} \text{to be solved} \\ \text{numerically by} \\ \text{a good algorithm} \end{array}$$

## HYBRID (A.K.A. HAMILTON) MONTE CARLO

- Random walk behaviour  
(e.g. Metropolis-Hastings)



Would be convenient to have guidance on where to move next!

- Hamiltonian dynamics at work

$$* \vec{\theta} \equiv \vec{q} \rightarrow (\vec{q}, \vec{p}) \quad \underbrace{\hspace{1cm}}_{= \ln p(\vec{q}, \text{Data})}$$

$$* H(\vec{q}, \vec{p}) = K(\vec{p}) + V(\vec{q})$$

$$\left. \begin{aligned} \frac{d\vec{q}}{dt} &= \frac{\partial H}{\partial \vec{p}}, & \frac{d\vec{p}}{dt} &= -\frac{\partial H}{\partial \vec{q}} \end{aligned} \right\} \begin{array}{l} \text{to be solved} \\ \text{numerically by} \\ \text{a good algorithm} \end{array}$$

$$\begin{array}{c} \text{deterministic} \\ \swarrow \\ x(\vec{q}_0, \vec{p}_0) \end{array} \quad \begin{array}{c} \searrow \\ x(\vec{q}', \vec{p}') \end{array}$$

## HYBRID (A.K.A. HAMILTON) MONTE CARLO

- Random walk behaviour  
(e.g. Metropolis-Hastings)



Would be convenient to have guidance on where to move next!

- Hamiltonian dynamics at work

$$* \vec{\theta} \equiv \vec{q} \rightarrow (\vec{q}, \vec{p}) \quad \rightarrow = \ln p(\vec{q}, \text{Data})$$

$$* H(\vec{q}, \vec{p}) = K(\vec{p}) + V(\vec{q})$$

$$\left. \begin{aligned} \frac{d\vec{q}}{dt} &= \frac{\partial H}{\partial \vec{p}}, & \frac{d\vec{p}}{dt} &= -\frac{\partial H}{\partial \vec{q}} \end{aligned} \right\} \begin{array}{l} \text{to be solved} \\ \text{numerically by} \\ \text{a good algorithm} \end{array}$$

Probability of accepting

$$\min\left(1, \frac{e^{-H'}}{e^{-H_0}}\right)$$

$\vec{q}_0, \vec{p}_0$   $\xrightarrow{\text{deterministic}}$   $\vec{q}', \vec{p}'$

DIMITRIOU, RAMOS, TELO AND BZ (WORK IN PROGRESS)

---

- The predictive distribution contain "hyper-parameters"



- The predictive distribution contain "hyper-parameters"

$$P(y_* | x_*, \text{Data}) = \int d\vec{\theta} \left( \begin{array}{c} \text{Likelihood} \\ \text{of } y_* | x_*, \vec{\theta} \end{array} \right) \cdot \left( \begin{array}{c} \text{Posterior of} \\ \vec{\theta} | \text{Data} \end{array} \right)$$

*model parameters*

- The predictive distribution contain "hyper-parameters"

$$P(y_* | x_*, \text{Data}) = \int d\vec{\theta} \left( \text{Likelihood of } y_* | x_*, \vec{\theta} \right) \cdot \left( \text{Posterior of } \vec{\theta} | \text{Data} \right) \propto p(\text{Data} | \vec{\theta}, \vec{\alpha}) p(\vec{\theta} | \vec{\beta})$$

The equation is annotated with handwritten notes:
 

- A blue arrow points from "model parameters" to  $d\vec{\theta}$ .
- Red arrows point from "Likelihood hyper-parameters" to  $\vec{\alpha}$  and from "prior hyper-parameters" to  $\vec{\beta}$ .

- The predictive distribution contain "hyper-parameters"

$$P(y_* | x_*, \text{Data}) = \int d\vec{\theta} \left( \text{Likelihood of } y_* | x_*, \vec{\theta} \right) \cdot \left( \text{Posterior of } \vec{\theta} | \text{Data} \right) \propto p(\text{Data} | \vec{\theta}, \vec{\alpha}) p(\vec{\theta} | \vec{\beta})$$

↗ model parameters  
↖ Likelihood hyper-parameters  
↖ Prior hyper-parameters

e.g. Prior =  $N(\vec{\theta} | \vec{\mu}, \sigma \cdot \mathbb{I})$

- The predictive distribution contain "hyper-parameters"

$$P(y_* | x_*, \text{Data}) = \int d\vec{\theta} \underbrace{\left( \text{Likelihood of } y_* | x_*, \vec{\theta} \right)}_{P(y_* | x_*, \vec{\theta}, \vec{\alpha})} \cdot \left( \text{Posterior of } \vec{\theta} | \text{Data} \right) \propto p(\text{Data} | \vec{\theta}, \vec{\alpha}) p(\vec{\theta} | \vec{\beta})$$

Likelihood hyper-parameters prior hyper-parameters

e.g. prior =  $N(\vec{\theta} | \vec{\mu}, \sigma \cdot \mathbb{I})$

- The predictive distribution contain "hyper-parameters"

$$P(y_* | x_*, \text{Data}) = \int d\vec{\theta} \underbrace{\left( \text{Likelihood of } y_* | x_*, \vec{\theta} \right)}_{P(y_* | x_*, \vec{\theta}, \vec{\alpha})} \cdot \left( \text{Posterior of } \vec{\theta} | \text{Data} \right) \propto p(\text{Data} | \vec{\theta}, \vec{\alpha}) p(\vec{\theta} | \vec{\beta})$$

Likelihood hyper-parameters prior hyper-parameters

e.g. prior =  $N(\vec{\theta} | \vec{\mu}, \sigma \cdot \mathbb{I})$

- How to optimize these hyper-parameters?

- The predictive distribution contain "hyper-parameters"

$$P(y_* | x_*, \text{Data}) = \int d\vec{\theta} \underbrace{\left( \text{Likelihood of } y_* | x_*, \vec{\theta} \right)}_{P(y_* | x_*, \vec{\theta}, \vec{\alpha})} \cdot \left( \text{Posterior of } \vec{\theta} | \text{Data} \right) \propto p(\text{Data} | \vec{\theta}, \vec{\alpha}) p(\vec{\theta} | \vec{\beta})$$

Likelihood hyper-parameters prior hyper-parameters

e.g. prior =  $N(\vec{\theta} | \vec{\mu}, \sigma \cdot \mathbb{I})$

- How to optimize these hyper-parameters?

1) Grid/random search

- The predictive distribution contain "hyper-parameters"

$$P(y_* | x_*, \text{Data}) = \int d\vec{\theta} \underbrace{\left( \text{Likelihood of } y_* | x_*, \vec{\theta} \right)}_{P(y_* | x_*, \vec{\theta}, \vec{\alpha})} \cdot \left( \text{Posterior of } \vec{\theta} | \text{Data} \right) \propto p(\text{Data} | \vec{\theta}, \vec{\alpha}) p(\vec{\theta} | \vec{\beta})$$

Likelihood hyper-parameters  $\rightarrow \vec{\alpha}$ 
prior hyper-parameters  $\rightarrow \vec{\beta}$

e.g. prior =  $N(\vec{\theta} | \vec{\mu}, \sigma \cdot \mathbb{I})$

- How to optimize these hyper-parameters?

1) Grid/random search  $\leftarrow$  curse of dimensionality

- The predictive distribution contain "hyper-parameters"

$$P(y_* | x_*, \text{Data}) = \int d\vec{\theta} \underbrace{\left( \text{Likelihood of } y_* | x_*, \vec{\theta} \right)}_{P(y_* | x_*, \vec{\theta}, \vec{\alpha})} \cdot \left( \text{Posterior of } \vec{\theta} | \text{Data} \right) \propto p(\text{Data} | \vec{\theta}, \vec{\alpha}) p(\vec{\theta} | \vec{\beta})$$

Likelihood hyper-parameters  $\rightarrow \vec{\alpha}$ 
prior hyper-parameters  $\rightarrow \vec{\beta}$

e.g. prior =  $N(\vec{\theta} | \vec{\mu}, \sigma \cdot \mathbb{I})$

- How to optimize these hyper-parameters?

- 1) Grid/random search ↔ curse of dimensionality
- 2) "Bayesian optimisation" (new proposal based on previous evaluations)



- The predictive distribution contain "hyper-parameters"

$$P(y_* | x_*, \text{Data}) = \int d\vec{\theta} \underbrace{\left( \text{Likelihood of } y_* | x_*, \vec{\theta} \right) \cdot \left( \text{Posterior of } \vec{\theta} | \text{Data} \right)}_{P(y_* | x_*, \vec{\theta}, \vec{\alpha})} \propto p(\text{Data} | \vec{\theta}, \vec{\alpha}) p(\vec{\theta} | \vec{\beta})$$

Likelihood hyper-parameters  $\rightarrow \vec{\alpha}$ 
prior hyper-parameters  $\rightarrow \vec{\beta}$

e.g. prior =  $N(\vec{\theta} | \vec{\mu}, \sigma \cdot \mathbb{I})$

- How to optimize these hyper-parameters?

- Grid/random search  $\leftarrow$  curse of dimensionality
- "Bayesian optimisation" (new proposal based on previous evaluations)  $\leftarrow$  Depends on extra tuning parameters

- The predictive distribution contain "hyper-parameters"

$$P(y_* | x_*, \text{Data}) = \int d\vec{\theta} \underbrace{\left( \text{Likelihood of } y_* | x_*, \vec{\theta} \right)}_{P(y_* | x_*, \vec{\theta}, \vec{\alpha})} \cdot \left( \text{Posterior of } \vec{\theta} | \text{Data} \right) \propto p(\text{Data} | \vec{\theta}, \vec{\alpha}) p(\vec{\theta} | \vec{\beta})$$

Likelihood hyper-parameters  $\rightarrow \vec{\alpha}$ 
prior hyper-parameters  $\rightarrow \vec{\beta}$

e.g. prior =  $N(\vec{\theta} | \vec{\mu}, \sigma \cdot \mathbb{I})$

- How to optimize these hyper-parameters?

- Grid/random search  $\leftarrow$  curse of dimensionality
- "Bayesian optimisation" (new proposal based on previous evaluations)  $\leftarrow$  Depends on extra tuning parameters
- Maximizing the model Evidence  
an approximation of

- The predictive distribution contain "hyper-parameters"

$$P(y_* | x_*, \text{Data}) = \int d\vec{\theta} \underbrace{\left( \text{Likelihood of } y_* | x_*, \vec{\theta} \right)}_{P(y_* | x_*, \vec{\theta}, \vec{\alpha})} \cdot \left( \text{Posterior of } \vec{\theta} | \text{Data} \right) \propto p(\text{Data} | \vec{\theta}, \vec{\alpha}) p(\vec{\theta} | \vec{\beta})$$

Likelihood hyper-parameters prior hyper-parameters

e.g. prior =  $N(\vec{\theta} | \vec{\mu}, \sigma \cdot \mathbb{I})$

- How to optimize these hyper-parameters?

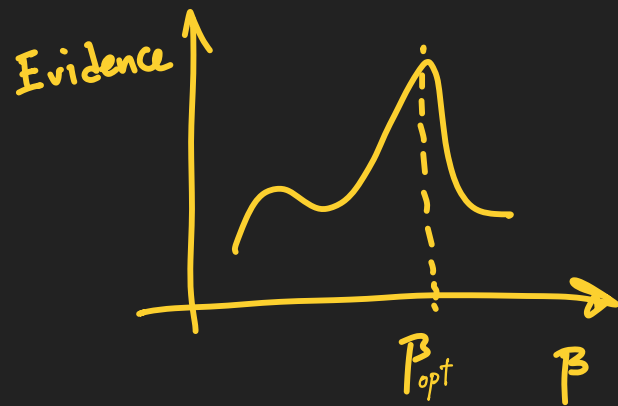
- 1) Grid/random search ← Curse of dimensionality
- 2) "Bayesian optimisation" (new proposal based on previous evaluations) ← Depends on extra tuning parameters
- 3) Maximizing the model Evidence an approximation of ←
  - No guarantees for robust predictions
  - Contribution to the predictive variance?

DIMITRIOU, RAMOS, TELO AND BZ (WORK IN PROGRESS)

---

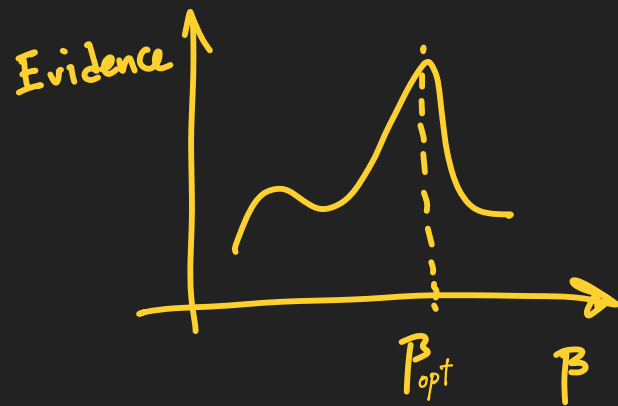
- What does it mean a good value for the hyper-parameters?

- What does it mean a good value for the hyper-parameters?
  - The one maximizing the model evidence?

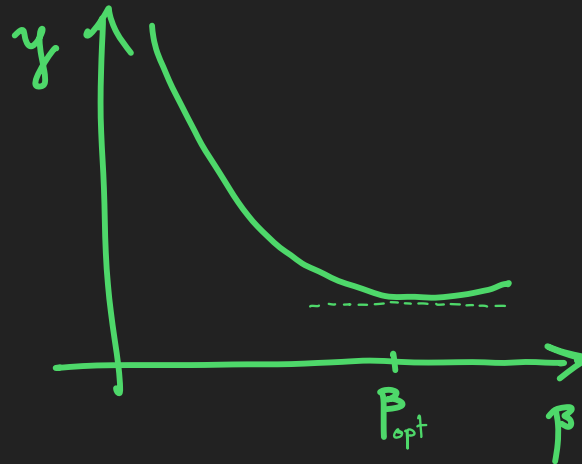


- What does it mean a good value for the hyper-parameters?

- The one maximizing the model evidence?

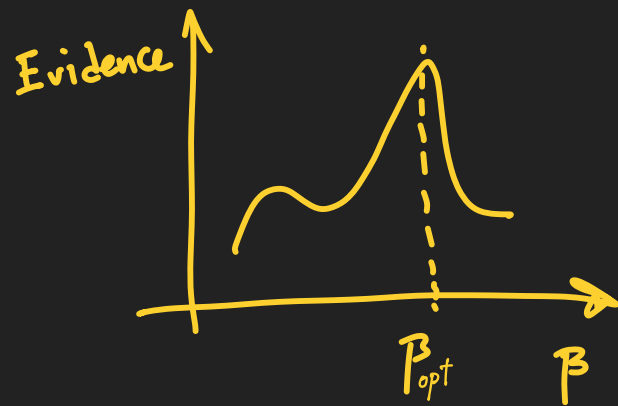


- The one ensuring robust predictions?

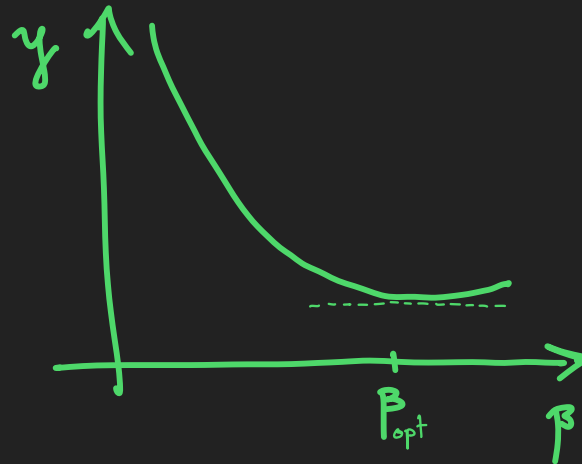


- What does it mean a good value for the hyper-parameters?

- The one maximizing the model evidence?



- The one ensuring robust predictions?

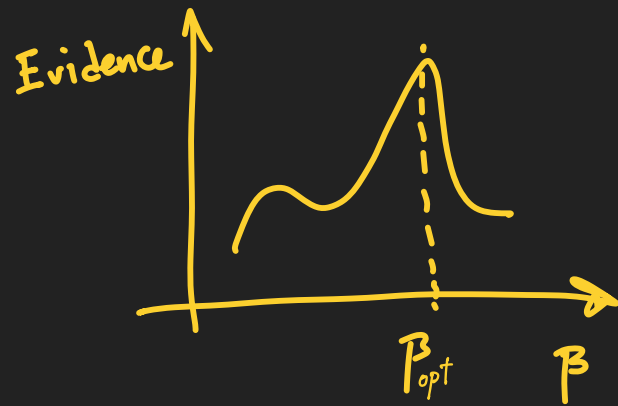


} Depends on your scope!

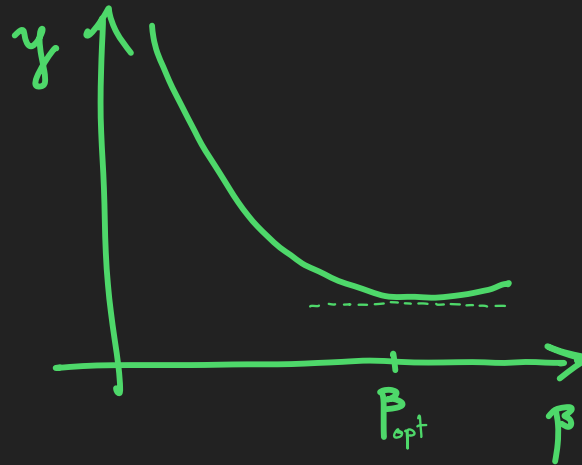


- What does it mean a good value for the hyper-parameters?

- The one maximizing the model evidence?



- The one ensuring robust predictions?

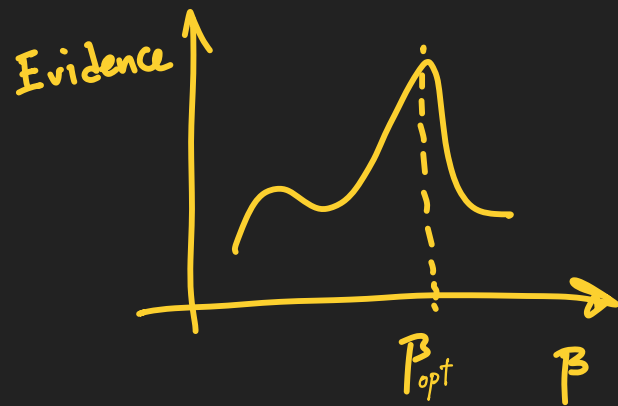


} Depends on your scope!

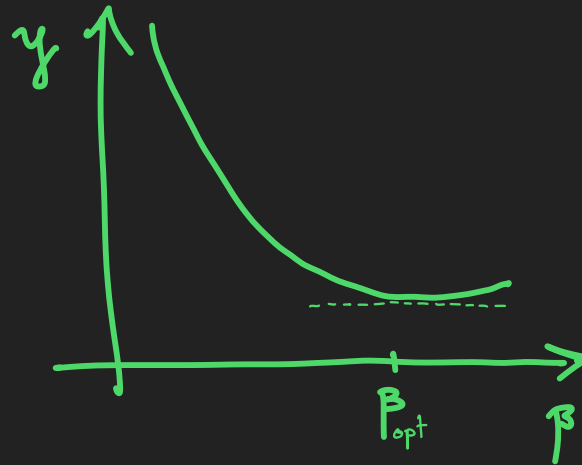
- Here we follow the second approach

- What does it mean a good value for the hyper-parameters?

- The one maximizing the model evidence?



- The one ensuring robust predictions?



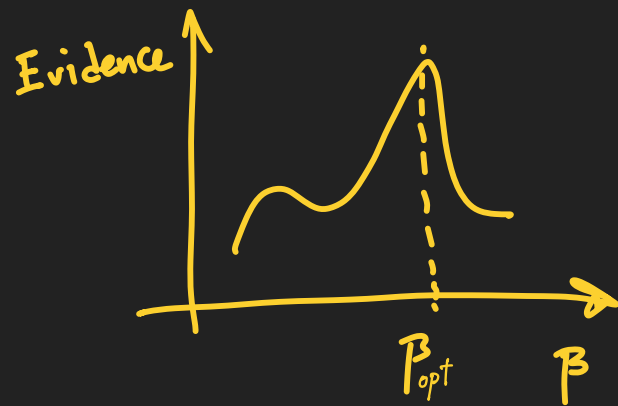
} Depends on your scope!

- Here we follow the **second approach**

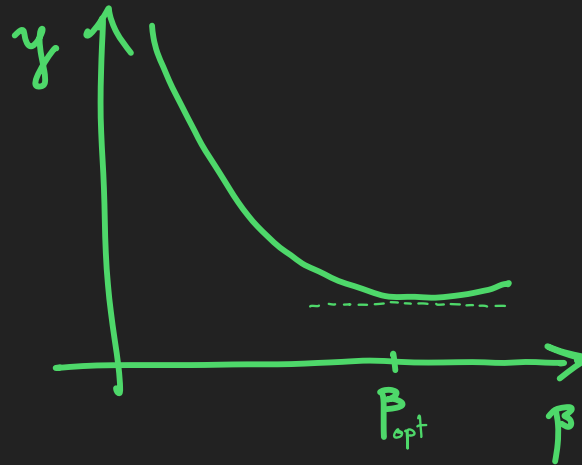
- Develop a methodology for estimating uncertainties associated with our ignorance about the hyper- $p$  values

- What does it mean a good value for the hyper-parameters?

- The one maximizing the model evidence?



- The one ensuring robust predictions?



Depends on your scope!

- Here we follow the **second approach**

- Develop a methodology for estimating uncertainties associated with our ignorance about the hyper- $p$  values

Not addressed so far in the literature

DEMTRIIOU, RAMOS, TELO AND BZ (WORK IN PROGRESS)

---

- Idea: Compute the derivative of the predictive distribution w.r.t. the hyper-parameters  $\rightarrow$  if very small, then predictions are robust

- Idea: Compute the derivative of the predictive distribution w.r.t. the hyper-parameters  $\rightarrow$  if very small, then predictions are robust

$$P(y_* | x_*, \text{Data}) \stackrel{\lim N_s \rightarrow \infty}{=} \frac{1}{N_s} \sum_{s=1}^{N_s} P(y_* | x_*, \vec{\theta}_s, \vec{\alpha})$$

- Idea: Compute the derivative of the predictive distribution w.r.t. the hyper-parameters  $\rightarrow$  if very small, then predictions are robust

$$P(y_* | x_*, \text{Data}) \stackrel{\lim N_s \rightarrow \infty}{=} \frac{1}{N_s} \sum_{s=1}^{N_s} P(y_* | x_*, \vec{\theta}_s, \vec{\alpha}) \quad \text{with } \vec{\theta}_s \sim \text{Posterior}(\vec{\theta} | \text{Data}, \vec{\alpha}, \vec{\beta})$$

$\nwarrow$  MCMC  $\swarrow$

- Idea: Compute the derivative of the predictive distribution w.r.t. the hyper-parameters  $\rightarrow$  if very small, then predictions are robust

$$P(y_* | x_*, \text{Data}) \stackrel{\lim N_s \rightarrow \infty}{=} \frac{1}{N_s} \sum_{s=1}^{N_s} P(y_* | x_*, \vec{\theta}_s, \vec{\alpha}) \quad \text{with } \vec{\theta}_s \sim \text{Posterior}(\vec{\theta} | \text{Data}, \vec{\alpha}, \vec{\beta})$$

$\nwarrow$  MCMC  $\swarrow$

$\rightarrow$  explicit dependence on hyper-parameters is lost!



- Idea: Compute the derivative of the predictive distribution w.r.t. the hyper-parameters  $\rightarrow$  if very small, then predictions are robust

$$P(y_* | x_*, \text{Data}) \stackrel{\lim N_s \rightarrow \infty}{=} \frac{1}{N_s} \sum_{s=1}^{N_s} P(y_* | x_*, \vec{\theta}_s, \vec{\alpha}) \quad \text{with } \vec{\theta}_s \sim \text{Posterior}(\vec{\theta} | \text{Data}, \vec{\alpha}, \vec{\beta})$$

$\nwarrow$  MCMC  $\nwarrow$

$\rightarrow$  explicit dependence on hyper-parameters is lost!

$$\frac{\partial P(y_* | x_*, \text{Data})}{\partial \vec{\beta}} = ?$$

DIMITRIOU, RAMOS, TELO AND BZ (WORK IN PROGRESS)

---

One possible way: Numerical Stochastic Perturbation Theory  
using HMC (borrowed from Lattice QCD)

One possible way: Numerical Stochastic Perturbation Theory  
using HMC (borrowed from Lattice QCD)

- Expand in Taylor series
- $$\vec{\theta} = \vec{\theta}^{(0)} + \vec{\theta}^{(1)} \cdot \delta \beta + \dots$$
- $$\vec{p} = \vec{p}^{(0)} + \vec{p}^{(1)} \cdot \delta \beta + \dots$$
- hyper-parameter*

One possible way: Numerical Stochastic Perturbation Theory  
using HMC (borrowed from Lattice QCD)

- Expand in Taylor series

$$\begin{aligned}\vec{\theta} &= \vec{\theta}^{(0)} + \vec{\theta}^{(1)} \cdot \delta \beta + \dots \\ \vec{p} &= \vec{p}^{(0)} + \vec{p}^{(1)} \cdot \delta \beta + \dots\end{aligned}$$

↗ hyper-parameter

$$\left\{ \begin{aligned}\dot{\vec{\theta}} &= \frac{\partial H}{\partial \vec{p}} \\ \dot{\vec{p}} &= -\frac{\partial H}{\partial \vec{\theta}}\end{aligned}\right.$$

One possible way: Numerical Stochastic Perturbation Theory  
using HMC (borrowed from Lattice QCD)

- Expand in Taylor series
 
$$\vec{\theta} = \vec{\theta}^{(0)} + \vec{\theta}^{(1)} \cdot \delta \beta + \dots$$

$$\vec{p} = \vec{p}^{(0)} + \vec{p}^{(1)} \cdot \delta \beta + \dots$$

$$\frac{\partial H}{\partial \vec{p}} = \left( \frac{\partial H}{\partial \vec{p}} \right)^{(0)} + \left( \frac{\partial H}{\partial \vec{p}} \right)^{(1)} \cdot \delta \beta + \dots$$

$$\left\{ \begin{array}{l} \dot{\vec{\theta}} = \frac{\partial H}{\partial \vec{p}} \\ \dot{\vec{p}} = - \frac{\partial H}{\partial \vec{\theta}} \end{array} \right.$$

One possible way: Numerical Stochastic Perturbation Theory  
using HMC (borrowed from Lattice QCD)

• Expand in Taylor series

$$\vec{\theta} = \vec{\theta}^{(0)} + \vec{\theta}^{(1)} \cdot \delta \beta + \dots$$

$$\vec{p} = \vec{p}^{(0)} + \vec{p}^{(1)} \cdot \delta \beta + \dots$$

$\rightarrow$  hyper-parameter

$$\left\{ \begin{array}{l} \dot{\vec{\theta}} = \frac{\partial H}{\partial \vec{p}} \\ \dot{\vec{p}} = -\frac{\partial H}{\partial \vec{\theta}} \end{array} \right.$$

$$\frac{\partial H}{\partial \vec{p}} = \left( \frac{\partial H}{\partial \vec{p}} \right)^{(0)} + \left( \frac{\partial H}{\partial \vec{p}} \right)^{(1)} \cdot \delta \beta + \dots \quad (\text{similarly with } \frac{\partial H}{\partial \vec{\theta}})$$

One possible way: Numerical Stochastic Perturbation Theory  
using HMC (borrowed from Lattice QCD)

- Expand in Taylor series

$$\begin{aligned}\vec{\theta} &= \vec{\theta}^{(0)} + \vec{\theta}^{(1)} \cdot \delta \beta + \dots \\ \vec{p} &= \vec{p}^{(0)} + \vec{p}^{(1)} \cdot \delta \beta + \dots\end{aligned}$$

↗ hyper-parameter

$$\left\{ \begin{aligned}\dot{\vec{\theta}} &= \frac{\partial H}{\partial \vec{p}} \\ \dot{\vec{p}} &= -\frac{\partial H}{\partial \vec{\theta}}\end{aligned} \right.$$

$$\frac{\partial H}{\partial \vec{p}} = \left( \frac{\partial H}{\partial \vec{p}} \right)^{(0)} + \left( \frac{\partial H}{\partial \vec{p}} \right)^{(1)} \cdot \delta \beta + \dots \quad (\text{similarly with } \frac{\partial H}{\partial \vec{\theta}})$$

- Results in a tower of Hamilton eqs. up to given truncation order



One possible way: Numerical Stochastic Perturbation Theory  
using HMC (borrowed from Lattice QCD)

- Expand in Taylor series

$$\begin{aligned}\vec{\theta} &= \vec{\theta}^{(0)} + \vec{\theta}^{(1)} \cdot \delta \beta + \dots \\ \vec{p} &= \vec{p}^{(0)} + \vec{p}^{(1)} \cdot \delta \beta + \dots\end{aligned}$$

hyper-parameter

$$\left\{ \begin{aligned}\dot{\vec{\theta}} &= \frac{\partial H}{\partial \vec{p}} \\ \dot{\vec{p}} &= -\frac{\partial H}{\partial \vec{\theta}}\end{aligned} \right.$$

$$\frac{\partial H}{\partial \vec{p}} = \left( \frac{\partial H}{\partial \vec{p}} \right)^{(0)} + \left( \frac{\partial H}{\partial \vec{p}} \right)^{(1)} \cdot \delta \beta + \dots \quad (\text{similarly with } \frac{\partial H}{\partial \vec{\theta}})$$

- Results in a tower of Hamilton eqs. up to given truncation order

$$\left\{ \begin{aligned}\dot{\vec{\theta}}^{(n)} &= \left( \frac{\partial H}{\partial \vec{p}} \right)^{(n)} \\ \dot{\vec{p}}^{(n)} &= -\left( \frac{\partial H}{\partial \vec{\theta}} \right)^{(n)}\end{aligned} \right. \quad n=1, \dots, N_{\text{trunc}}$$

One possible way: Numerical Stochastic Perturbation Theory  
using HMC (borrowed from Lattice QCD)

- Expand in Taylor series

$$\begin{aligned}\vec{\theta} &= \vec{\theta}^{(0)} + \vec{\theta}^{(1)} \cdot \delta \beta + \dots \\ \vec{p} &= \vec{p}^{(0)} + \vec{p}^{(1)} \cdot \delta \beta + \dots\end{aligned}$$

hyper-parameter

$$\begin{cases} \dot{\vec{\theta}} = \frac{\partial H}{\partial \vec{p}} \\ \dot{\vec{p}} = -\frac{\partial H}{\partial \vec{\theta}} \end{cases}$$

$$\frac{\partial H}{\partial \vec{p}} = \left( \frac{\partial H}{\partial \vec{p}} \right)^{(0)} + \left( \frac{\partial H}{\partial \vec{p}} \right)^{(1)} \cdot \delta \beta + \dots \quad (\text{similarly with } \frac{\partial H}{\partial \vec{\theta}})$$

- Results in a tower of Hamilton eqs. up to given truncation order

$$\begin{cases} \dot{\vec{\theta}}^{(n)} = \left( \frac{\partial H}{\partial \vec{p}} \right)^{(n)} \\ \dot{\vec{p}}^{(n)} = -\left( \frac{\partial H}{\partial \vec{\theta}} \right)^{(n)} \end{cases} \quad n=1, \dots, N_{\text{trunc}}$$



$$P(y_* | x_*, \text{Data}) \equiv \mathcal{P}_*$$

$$\mathcal{P}_* = \mathcal{P}_*^{(0)} + \mathcal{P}_*^{(1)} \cdot \delta \beta + \dots$$

One possible way: Numerical Stochastic Perturbation Theory  
using HMC (borrowed from Lattice QCD)

- Expand in Taylor series

$$\begin{aligned}\vec{\theta} &= \vec{\theta}^{(0)} + \vec{\theta}^{(1)} \cdot \delta \beta + \dots \\ \vec{p} &= \vec{p}^{(0)} + \vec{p}^{(1)} \cdot \delta \beta + \dots\end{aligned}$$

$\delta \beta$  hyper-parameter

$$\begin{cases} \dot{\vec{\theta}} = \frac{\partial H}{\partial \vec{p}} \\ \dot{\vec{p}} = -\frac{\partial H}{\partial \vec{\theta}} \end{cases}$$

$$\frac{\partial H}{\partial \vec{p}} = \left( \frac{\partial H}{\partial \vec{p}} \right)^{(0)} + \left( \frac{\partial H}{\partial \vec{p}} \right)^{(1)} \cdot \delta \beta + \dots \quad (\text{similarly with } \frac{\partial H}{\partial \vec{\theta}})$$

- Results in a tower of Hamilton eqs. up to given truncation order

$$\begin{cases} \dot{\vec{\theta}}^{(n)} = \left( \frac{\partial H}{\partial \vec{p}} \right)^{(n)} \\ \dot{\vec{p}}^{(n)} = -\left( \frac{\partial H}{\partial \vec{\theta}} \right)^{(n)} \end{cases} \quad n=1, \dots, N_{\text{trunc}}$$



$$\begin{aligned} P(y_* | x_*, \text{Data}) &\equiv \mathcal{P}_* \\ \mathcal{P}_* &= \mathcal{P}_*^{(0)} + \mathcal{P}_*^{(1)} \cdot \delta \beta + \dots \end{aligned}$$

this is what we are searching for

DIMITRIOU, RAMOS, TELO AND BZ (WORK IN PROGRESS)

---

This can be implemented using a  
single HMC simulation!

This can be implemented using a  
single HMC simulation!  
[with  $N_{\text{part}}$  coupled Hamilton eqs.]

This can be implemented using a single HMC simulation!

[with  $N_{\text{pert}}$  coupled Hamilton eqs.]

the magic is possible thanks to the concept of "Automatic Differentiation" and "Hyper-dual  $\mathbb{H}_s$ "  
(see A. Ramos, 1809.01289)

This can be implemented using a single HMC simulation!

[with  $N_{\text{part}}$  coupled Hamilton eqs.]

the magic is possible thanks to the concept of "Automatic Differentiation" and "Hyper-dual  $\mathbb{H}_s$ "  
(see A. Ramos, 1809.01289)

- Toy example :



This can be implemented using a single HMC simulation!

[with  $N_{\text{part}}$  coupled Hamilton eqs.]

the magic is possible thanks to the concept of "Automatic Differentiation" and "Hyper-dual  $\mathbb{H}_s$ "  
(see A. Ramos, 1809.01289)

- Toy example :

- Gaussian data generated by

$$f(x) = x^3 + x^2 + x + 1$$

This can be implemented using a single HMC simulation!

[with  $N_{\text{pert}}$  coupled Hamilton eqs.]

the magic is possible thanks to the concept of "Automatic Differentiation" and "Hyper-dual  $\mathbb{H}_s$ "  
(see A. Ramos, 1809.01289)

- Toy example :

- Gaussian data generated by

$$f(x) = x^3 + x^2 + x + 1$$

- Model is a  
1-hidden layer  
Bayesian Neural net  
with prior variance =  $\sigma_p^2$

This can be implemented using a single HMC simulation!

[with  $N_{\text{pert}}$  coupled Hamilton eqs.]

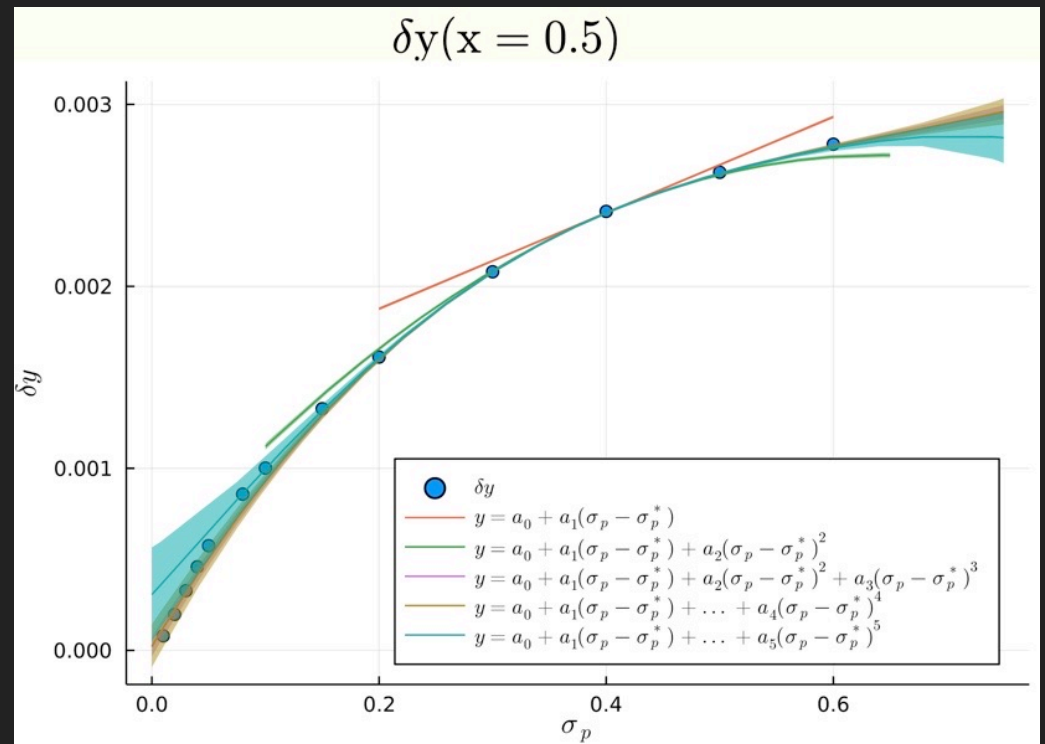
the magic is possible thanks to the concept of "Automatic Differentiation" and "Hyper-dual  $\#_s$ " (see A. Ramos, 1809.01289)

### • Toy example :

- Gaussian data generated by

$$f(x) = x^3 + x^2 + x + 1$$

- Model is a  
1-hidden layer  
Bayesian Neural net  
with prior variance =  $\sigma_p^2$



This can be implemented using a single HMC simulation!

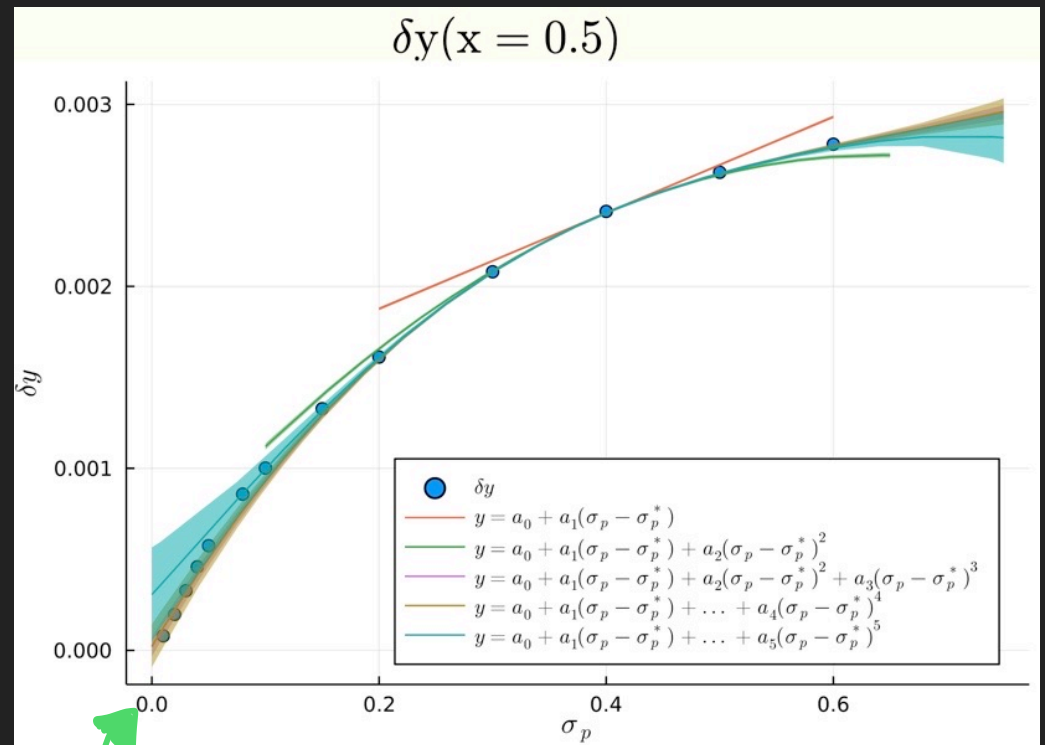
[with  $N_{\text{pert}}$  coupled Hamilton eqs.]

the magic is possible thanks to the concept of "Automatic Differentiation" and "Hyper-dual  $\#_s$ " (see A. Ramos, 1809.01289)

- Toy example :

- Gaussian data generated by  $f(x) = x^3 + x^2 + x + 1$

- Model is a 1-hidden layer Bayesian Neural net with prior variance  $= \sigma_p^2$



Results: The expansion captures successfully the actual behavior!

This can be implemented using a single HMC simulation!

[with  $N_{\text{pert}}$  coupled Hamilton eqs.]

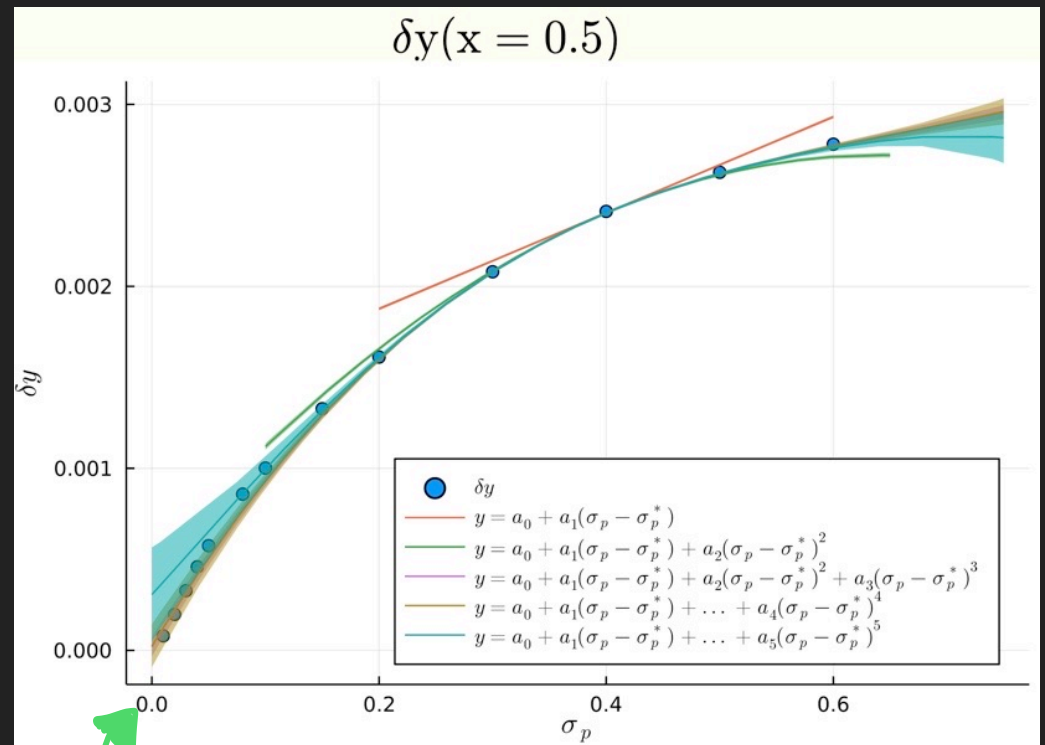
the magic is possible thanks to the concept of "Automatic Differentiation" and "Hyper-dual  $\#_s$ " (see A. Ramos, 1809.01289)

- Toy example :

- Gaussian data generated by

$$f(x) = x^3 + x^2 + x + 1$$

- Model is a 1-hidden layer Bayesian Neural net with prior variance =  $\sigma_p^2$



Results: The expansion captures successfully the actual behavior! Clearly this works

## FINAL TAKE-HOME MESSAGE

---

## FINAL TAKE-HOME MESSAGE

---

Bayesian Inference is in a golden era nowadays

## FINAL TAKE-HOME MESSAGE

---

Bayesian Inference is in a golden era nowadays  
Plenty of opportunities for two-way contributions



## FINAL TAKE-HOME MESSAGE

---

Bayesian Inference is in a golden era nowadays

Plenty of opportunities for two-way contributions

Physics,  
Biology, etc

Bayesian Inference,  
Machine Learning

## FINAL TAKE-HOME MESSAGE

---

Bayesian Inference is in a golden era nowadays

Plenty of opportunities for two-way contributions



## FINAL TAKE-HOME MESSAGE

---

Bayesian Inference is in a golden era nowadays

Plenty of opportunities for two-way contributions

