# Bayesian machine learning: a tutorial

#### Rémi Bardenet

CNRS & CRIStAL, Univ. Lille, France





#### **Outline**

#### The what

Typical statistical problems Statistical decision theory Posterior expected utility and Bayes rules

## The why

The philosophical why
The practical why

#### The how

Conjugacy Monte Carlo methods Metropolis-Hastings

# Variational approximations

# In depth with Gaussian processes in ML

From linear regression to GPs Modeling and learning More applications References and open issues

#### **Outline**

#### The what

Typical statistical problems Statistical decision theory Posterior expected utility and Bayes rules

### The why

The philosophical why The practical why

#### The how

Monte Carlo methods
Metropolis-Hastings

## In depth with Gaussian processes in ML

From linear regression to G Modeling and learning More applications References and open issues

## Typical jobs for statisticians

#### **Estimation**

- You have data  $x_1, \ldots, x_n$  that you assume drawn from  $p(x_1, \ldots, x_n | \theta^*)$ , with  $\theta^* \in \mathbb{R}^d$ .
- ▶ You want an estimate  $\hat{\theta}(x_1, ..., x_n)$  of  $\theta^* \in \mathbb{R}^d$ .

#### Confidence regions

- You have data  $x_1, \ldots, x_n$  that you assume drawn from  $p(x_1, \ldots, x_n \cdot | \theta^*)$ , with  $\theta^* \in \mathbb{R}^d$ .
- You want a region  $A(x_1, ..., x_n) \subset \mathbb{R}^d$  and make a statement that  $\theta \in A(x_1, ..., x_n)$  with some certainty.

## Typical jobs for statisticians

#### **Estimation**

- You have data  $x_1, \ldots, x_n$  that you assume drawn from  $p(x_1, \ldots, x_n | \theta^*)$ , with  $\theta^* \in \mathbb{R}^d$ .
- ▶ You want an estimate  $\hat{\theta}(x_1, ..., x_n)$  of  $\theta^* \in \mathbb{R}^d$ .

## **Confidence regions**

- You have data  $x_1, \ldots, x_n$  that you assume drawn from  $p(x_1, \ldots, x_n \cdot | \theta^*)$ , with  $\theta^* \in \mathbb{R}^d$ .
- You want a region  $A(x_1,...,x_n) \subset \mathbb{R}^d$  and make a statement that  $\theta \in A(x_1,...,x_n)$  with some certainty.

## Statistical decision theory<sup>1</sup>



Figure: Abraham Wald (1902–1950)

<sup>&</sup>lt;sup>1</sup>A. Wald. Statistical decision functions. Wiley, 1950.

- Let Θ be the "states of the world", typically the space of parameters of interest.
- ▶ Decisions are functions  $d(x_1, ..., x_n) \in \mathcal{D}$ .
- Let  $L(d, \theta)$  denote the loss of making decision d when the state of the world is  $\theta$ .
- Wald defines the risk of a decision as

$$R(d,\theta) = \int L(d,\theta) p(x_{1:n}|\theta) dx_{1:n}.$$

▶ Wald says  $d_1$  is a better decision than  $d_2$  if

$$\forall \theta \in \Theta, \quad L(d_1, \theta) \leqslant L(d_2, \theta).$$
 (1)

 $\triangleright$  d is called admissible if there is no better decision than d.

- Let Θ be the "states of the world", typically the space of parameters of interest.
- ▶ Decisions are functions  $d(x_1, ..., x_n) \in \mathcal{D}$ .
- Let  $L(d, \theta)$  denote the loss of making decision d when the state of the world is  $\theta$ .
- Wald defines the risk of a decision as

$$R(d,\theta) = \int L(d,\theta) p(x_{1:n}|\theta) dx_{1:n}.$$

▶ Wald says  $d_1$  is a better decision than  $d_2$  if

$$\forall \theta \in \Theta, \quad L(d_1, \theta) \leqslant L(d_2, \theta).$$
 (1)

 $\triangleright$  d is called admissible if there is no better decision than d.

- Let Θ be the "states of the world", typically the space of parameters of interest.
- ▶ Decisions are functions  $d(x_1, ..., x_n) \in \mathcal{D}$ .
- Let  $L(d, \theta)$  denote the loss of making decision d when the state of the world is  $\theta$ .
- Wald defines the risk of a decision as

$$R(d,\theta) = \int L(d,\theta) p(x_{1:n}|\theta) dx_{1:n}.$$

▶ Wald says  $d_1$  is a better decision than  $d_2$  if

$$\forall \theta \in \Theta, \quad L(d_1, \theta) \leqslant L(d_2, \theta).$$
 (1)

d is called admissible if there is no better decision than d.

- Let Θ be the "states of the world", typically the space of parameters of interest.
- ▶ Decisions are functions  $d(x_1, ..., x_n) \in \mathcal{D}$ .
- Let  $L(d, \theta)$  denote the loss of making decision d when the state of the world is  $\theta$ .
- Wald defines the risk of a decision as

$$R(d,\theta) = \int L(d,\theta)p(x_{1:n}|\theta)dx_{1:n}.$$

▶ Wald says  $d_1$  is a better decision than  $d_2$  if

$$\forall \theta \in \Theta, \quad L(d_1, \theta) \leqslant L(d_2, \theta).$$
 (1)

d is called admissible if there is no better decision than d.

- Let Θ be the "states of the world", typically the space of parameters of interest.
- ▶ Decisions are functions  $d(x_1, ..., x_n) \in \mathcal{D}$ .
- Let  $L(d, \theta)$  denote the loss of making decision d when the state of the world is  $\theta$ .
- Wald defines the risk of a decision as

$$R(d,\theta) = \int L(d,\theta)p(x_{1:n}|\theta)dx_{1:n}.$$

▶ Wald says  $d_1$  is a better decision than  $d_2$  if

$$\forall \theta \in \Theta, \quad L(d_1, \theta) \leqslant L(d_2, \theta).$$
 (1)

 $\triangleright$  d is called admissible if there is no better decision than d.

- Let Θ be the "states of the world", typically the space of parameters of interest.
- ▶ Decisions are functions  $d(x_1, ..., x_n) \in \mathcal{D}$ .
- Let  $L(d, \theta)$  denote the loss of making decision d when the state of the world is  $\theta$ .
- Wald defines the risk of a decision as

$$R(d,\theta) = \int L(d,\theta)p(x_{1:n}|\theta)dx_{1:n}.$$

ightharpoonup Wald says  $d_1$  is a better decision than  $d_2$  if

$$\forall \theta \in \Theta, \quad L(d_1, \theta) \leqslant L(d_2, \theta).$$
 (1)

d is called admissible if there is no better decision than d.

### Illustration with a simple estimation problem

You have data  $x_1, \ldots, x_n$  that you assume drawn from

$$p(x_1,\ldots,x_n|\theta^*)=\prod_{i=1}^n \mathcal{N}(x_i|\theta^*,\sigma^2),$$

and you know  $\sigma^2$ .

- ▶ You choose a loss function, say  $L(\hat{\theta}, \theta) = ||\hat{\theta} \theta||^2$ .
- You restrict your decision space to unbiased estimators.
- ▶ The sample mean  $\tilde{\theta} := n^{-1} \sum_{i=1}^{n} x_i$  is unbiased, and has minimum variance among unbiased estimators.
- Since

$$R(\tilde{\theta}, \theta) = Var\tilde{\theta},$$

 $ilde{ heta}$  is the best decision you can make in Wald's framework.

## Wald's view of frequentist estimation

#### **Estimation**

- You have data  $x_1, \ldots, x_n$  that you assume drawn from  $p(x_1, \ldots, x_n | \theta^*)$ , with  $\theta^* \in \mathbb{R}^d$ .
- ▶ You want an estimate  $\hat{\theta}(x_1, ..., x_n)$  of  $\theta^* \in \mathbb{R}^d$ .

#### A Waldian answer

- Our decisions are estimates  $d(x_1, \ldots, x_n) = \hat{\theta}(x_1, \ldots, x_n)$ .
- ▶ We pick a loss, say  $L(d, \theta) = L(\hat{\theta}, \theta) = ||\hat{\theta} \theta||^2$ .
- ► If you have an unbiased estimator with minimum variance, then this is the best decision among unbiased estimators.

### Wald's view of frequentist estimation

#### **Estimation**

- You have data  $x_1, \ldots, x_n$  that you assume drawn from  $p(x_1, \ldots, x_n | \theta^*)$ , with  $\theta^* \in \mathbb{R}^d$ .
- ▶ You want an estimate  $\hat{\theta}(x_1, ..., x_n)$  of  $\theta^* \in \mathbb{R}^d$ .

#### A Waldian answer

- Our decisions are estimates  $d(x_1, \ldots, x_n) = \hat{\theta}(x_1, \ldots, x_n)$ .
- ► In general, the loss can be more complex and unbiased estimors unknown/irrelevant.
- ▶ In these cases, you may settle for a minimax estimator

$$\hat{\theta}(x_1,\ldots,x_n) = \underset{d}{\operatorname{arg\,min\,sup}} R(d,\theta).$$

## Wald's is only one view of frequentist statistics...

► On estimation, some would argue in favour of the maximum likelihood<sup>2</sup>.



Figure: Ronald Fisher (1890–1962)

<sup>&</sup>lt;sup>2</sup>S. M. Stigler. "The epic story of maximum likelihood". In: *Statistical Science* (2007), pp. 598–620.

## ... but bear with me, since it is predominant in machine learning

For instance, supervised learning is usually formalized as

$$g^* = \arg\min_{g} \mathbb{E}L(y, g(x)).$$
 (2)

which you approximate by

$$\hat{g} = \underset{g}{\operatorname{arg \, min}} \sum_{i=1}^{n} L(y_i, g(x_i)) + \operatorname{penalty}(g),$$

while trying to control the excess risk

$$\mathbb{E}L(y,\hat{g}(x)) - \mathbb{E}L(y,g^*(x)).$$

## Wald's view of frequentist confidence regions

#### **Confidence regions**

- You have data  $x_1, \ldots, x_n$  that you assume drawn from  $p(x_1, \ldots, x_n \cdot | \theta^*)$ , with  $\theta^* \in \mathbb{R}^d$ .
- You want a region  $A(x_1,...,x_n) \subset \mathbb{R}^d$  and make a statement that  $\theta \in A(x_1,...,x_n)$  with some certainty.

#### A Waldian answer

- ▶ Our decisions are subsets of  $\mathbb{R}^d$ :  $d(x_{1:n}) = A(x_{1:n})$ .
- ▶ A common loss is  $L(d, \theta) = L(A, \theta) = 1_{\theta \notin A} + \gamma |A|$ .
- ▶ So you want to find  $A(x_{1:n})$  that minimizes

$$L(A,\theta) = \int \left[1_{\theta^{\star} \notin A} p(x_{1:n} | \theta^{\star}) + \gamma |A|\right] dx_{1:n}.$$

## Illustration with a simple confidence interval problem

 $\blacktriangleright$  You have data  $x_1, \ldots, x_n$  that you assume drawn from

$$p(x_1,\ldots,x_n|\theta^*)=\prod_{i=1}^n \mathcal{N}(x_i|\theta^*,\sigma^2).$$

- ▶ You choose a loss function, say  $L(A, \theta) = 1_{\theta \notin A} + \gamma |A|$ .
- You restrict your decisions to intervals centered around the sample mean  $\tilde{\theta}$ .
- ► Since  $\frac{\theta \tilde{\theta}}{\sigma / \sqrt{n}} \sim \mathcal{N}(0, 1)$ , we know (exercise) that for

$$\tilde{A} := [\tilde{\theta} - k\sigma/\sqrt{n}, \tilde{\theta} + k\sigma/\sqrt{n}],$$

it comes

$$R(A, \theta) = \mathbb{P}(|\mathcal{N}(0, 1)| \geqslant k) + \frac{2\gamma k\sigma}{\sqrt{n}}.$$

- ► All is left to do is choose k.
- ► Textbook examples bypass the need for  $\gamma$ : they fix  $\alpha > 0$  and find the smallest k such that  $\mathbb{P}(|\mathcal{N}(0,1)| \ge k) \le \alpha$ .

#### **Summary so far**

Waldian frequentists measure risks as expectations w.r.t. the data-generating process.

$$R(d,\theta) = \int L(d(x_{1:n}),\theta) p(x_{1:n}|\theta) dx_{1:n}$$

- ▶ One major difficulty is that the risk remains a function of  $\theta$ .
- Without additional structure (unbiasedness, Gaussianity, etc.), it is difficult to go beyond minimax rules.

#### Idea

What if we introduced a distribution on  $\Theta$ , and tried to minimize

$$r(d) = \int R(d,\theta)p(\theta)d\theta$$
$$= \int \left[\int L(d(x_{1:n}),\theta)p(x_{1:n}|\theta)dx_{1:n}\right]p(\theta)d\theta$$

## From expected frequentist loss to posterior expected loss

$$r(d) = \int R(d,\theta)p(\theta)d\theta$$

$$= \int \left[ \int L(d(x_{1:n}),\theta)p(x_{1:n}|\theta)dx_{1:n} \right] p(\theta)d\theta$$

$$= \int \left[ \int L(d(x_{1:n}),\theta)p(x_{1:n}|\theta)p(\theta)d\theta \right] dx_{1:n}$$

$$= \int \left[ \int L(d(x_{1:n}),\theta)\frac{p(x_{1:n}|\theta)p(\theta)}{p(x_{1:n})}d\theta \right] p(x_{1:n})dx_{1:n}$$

$$= \int \left[ \int L(d(x_{1:n}),\theta)p(\theta|x_{1:n})d\theta \right] p(x_{1:n})dx_{1:n}$$

## Bayesians minimize posterior expected utility

#### The posterior expected utility paradigm: Bayes rules

Pick d to solve

$$\arg\min_{d} \int L(d(x_{1:n}), \theta) p(\theta|x_{1:n}) d\theta.$$

## Bayes rules have good frequentist properties<sup>3</sup>

- ▶ Under general conditions, Bayes decision rules are admissible, all admissible rules are limits of Bayes rules.
- ▶ Bayes rules with "least favourable priors" are minimax.

<sup>&</sup>lt;sup>3</sup>G. Parmigiani and L. Inoue. *Decision theory: principles and approaches*. Vol. 812. John Wiley & Sons, 2009.

## Illustration with a simple estimation problem

You have data  $x_1, \ldots, x_n$  that you assume drawn from

$$p(x_1,\ldots,x_n|\theta^*)=\prod_{i=1}^n \mathcal{N}(x_i|\theta^*,\sigma^2),$$

and you know  $\sigma^2$ .

- ▶ You choose a loss function, say  $L(\hat{\theta}, \theta) = \|\hat{\theta} \theta\|^2$ .
- You choose a prior p over  $\theta$ .
- ► Your Bayes decision minimizes

$$\int \|\hat{\theta} - \theta\|^2 p(\theta|x_{1:n}) d\theta,$$

so you pick

$$\hat{\theta} = \int \theta p(\theta|x_{1:n})d\theta.$$

► Conceptually, it is simpler. In practice, you need to compute an integral.

#### Illustration with a simple confidence interval problem

You have data  $x_1, \ldots, x_n$  that you assume drawn from

$$p(x_1,\ldots,x_n|\theta^*)=\prod_{i=1}^n \mathcal{N}(x_i|\theta^*,\sigma^2).$$

- ▶ You choose a loss function, say  $L(A, \theta) = 1_{\theta \notin A} + \gamma |A|$ .
- ightharpoonup You choose a prior p over  $\theta$ .
- Your Bayes decision minimizes

$$\int 1_{\theta \notin A} p(\theta|x_{1:n}) d\theta + \gamma |A|,$$

► Conceptually, it is simpler. In practice, you need to carefully pick your prior and/or restrict the decision space and/or compute many integrals.

### Summary so far

- ▶ Bayes rules fit into Wald's framework.
- ► For a fixed prior, the Bayesian risk completely orders decision rules.
- ► The key idea is posterior expected utility.
- You can answer most basic statistical questions using this principle: [more examples].

#### A recent motivating success

PRL 119, 141101 (2017)

PHYSICAL REVIEW LETTERS

week ending 6 OCTOBER 2017

ဏ္

#### GW170814: A Three-Detector Observation of Gravitational Waves from a Binary Black Hole Coalescence

B. P. Abbott et al.\*

(LIGO Scientific Collaboration and Virgo Collaboration) (Received 23 September 2017; published 6 October 2017)

On August 14, 2017 at 10:30:43 UTC, the Advanced Virgo detector and the two Advanced LIGO detectors coherently observed a transient gravitational-wave signal produced by the coalescence of two stellar mass black holes, with a false-alarm rate of  $\lesssim$ 1 in 27 000 years. The signal was observed with a three-detector network matched-filter signal-to-noise ratio of 18. The inferred masses of the initial black holes are  $30.5^{\circ}_{23}M_{\odot}$  and  $45.5^{\circ}_{23}M_{\odot}$  and 45

DOI: 10.1103/PhysRevLett.119.141101

#### I. INTRODUCTION

The era of gravitational-wave (GW) astronomy began with the detection of binary black hole (BBH) mergers, by the Advanced Laser Interferometer Gravitational-Wave Observatory (LIGO) detectors [11] during the first of the

waveform obtained from analysis of the LIGO detectors' data alone, we find that the probability, in 5000 s of data around the event, of a peak in SNR from Virgo data due to noise and as large as the one observed, within a time window determined by the maximum possible time of

#### **Outline**

#### The what

Typical statistical problems Statistical decision theory Posterior expected utility and Bayes rules

## The why

The philosophical why The practical why

#### The how

Monte Carlo methods
Metropolis-Hastings

## In depth with Gaussian processes in ML

From linear regression to G Modeling and learning More applications References and open issues

## The subjectivistic viewpoint

- ► Top requirement is internal coherence of decisions.
- ► Favourizes interpreting probability distributions as personal beliefs.



**Figure:** Bruno de Finetti (1906–1985) and L. Jimmie Savage (1917–1971)

## The logical justification

- ► Top requirement is to find a version of propositional logic that allows taking into account uncertainty.
- ► Also favourizes interpreting probability distributions as beliefs, but aims for objective priors.



**Figure:** Richart T. Cox (1898–1991), Edwin T. Jaynes (1917–1971), and Harold Jeffreys (1891–1989)

## The hybrid view<sup>4</sup>

- The starting point is posterior expected utility, loosely justified by Wald's theory.
- ▶ It is simple, widely applicable, has good frequentist properties.
- It satisfies the likelihood principle.
- ▶ It is easy to interpret: beliefs are
  - represented by probabilities,
  - updated using Bayes' rule,
  - integrated when making decisions.
- It is easy to communicate your uncertainty
  - Simply give your posterior.
  - When making a decision, make sure that the priors of everyone involved would yield the same decision.

<sup>&</sup>lt;sup>4</sup>C. P. Robert. The Bayesian choice: from decision-theoretic foundations to computational implementation. Springer Science & Business Media, 2007.

#### Practical advantages of posterior expected utility

- Conceptually answers all ML problems.
- Suits all applications where quantifying uncertainty is vital vs computational complexity: all basic sciences, health, even one-shot commercial decisions.
- ► We never invoked any large-sample argument, so suits all sizes of datasets.

#### **Outline**

#### The what

Typical statistical problems Statistical decision theory Posterior expected utility and Bayes rules

## The why

The philosophical why The practical why

#### The how

Conjugacy Monte Carlo methods Metropolis-Hastings Variational approximations

## In depth with Gaussian processes in ML

From linear regression to GPs
Modeling and learning
More applications
References and open issues

Say we have a linear regression problem

$$y_i = f(x_i) + \varepsilon_i,$$

$$f(x) = \theta^T x$$
,  $\varepsilon_i$  i.i.d. Gaussians  $\mathcal{N}(0, \sigma^2)$ .

▶ If we choose  $p(\theta) \sim \mathcal{N}(0, \Sigma)$ , then (exercise)

$$p(\theta|(x,y)_{1:n}) \propto p((x,y)_{1:n}|\theta)p(\theta)$$
  
=  $\mathcal{N}(\sigma^{-2}A^{-1}X\mathbf{y},A^{-1})$ .

where 
$$A = \sigma^{-2}X^TX + \Sigma^{-1}$$
.

▶ If the loss is not too complicated, then integrals are easy. For instance, prediction with  $L^2$  loss is simple:

Say we have a linear regression problem

$$y_i = f(x_i) + \varepsilon_i,$$

$$f(x) = \theta^T x$$
,  $\varepsilon_i$  i.i.d. Gaussians  $\mathcal{N}(0, \sigma^2)$ .

▶ If we choose  $p(\theta) \sim \mathcal{N}(0, \Sigma)$ , then (exercise)

$$p(\theta|(x,y)_{1:n}) \propto p((x,y)_{1:n}|\theta)p(\theta)$$
  
=  $\mathcal{N}(\sigma^{-2}A^{-1}X\mathbf{y},A^{-1}),$ 

where 
$$A = \sigma^{-2}X^TX + \Sigma^{-1}$$
.

▶ If the loss is not too complicated, then integrals are easy. For instance, prediction with L² loss is simple:

Say we have a linear regression problem

$$y_i = f(x_i) + \varepsilon_i,$$

$$f(x) = \theta^T x$$
,  $\varepsilon_i$  i.i.d. Gaussians  $\mathcal{N}(0, \sigma^2)$ .

▶ If we choose  $p(\theta) \sim \mathcal{N}(0, \Sigma)$ , then (exercise)

$$p(\theta|(x,y)_{1:n}) \propto p((x,y)_{1:n}|\theta)p(\theta)$$
  
=  $\mathcal{N}(\sigma^{-2}A^{-1}X\mathbf{y},A^{-1}),$ 

where 
$$A = \sigma^{-2}X^TX + \Sigma^{-1}$$
.

▶ If the loss is not too complicated, then integrals are easy. For instance, prediction with  $L^2$  loss is simple:

$$\arg\min_{\mathbf{y}_{\star}} \int \|\mathbf{y}_{\star} - f(\mathbf{x}_{\star})\|^2 p(\theta|(\mathbf{x}, \mathbf{y})_{1:n}) d\theta$$

Say we have a linear regression problem

$$y_i = f(x_i) + \varepsilon_i,$$

$$f(x) = \theta^T x$$
,  $\varepsilon_i$  i.i.d. Gaussians  $\mathcal{N}(0, \sigma^2)$ .

▶ If we choose  $p(\theta) \sim \mathcal{N}(0, \Sigma)$ , then (exercise)

$$p(\theta|(x,y)_{1:n}) \propto p((x,y)_{1:n}|\theta)p(\theta)$$
  
=  $\mathcal{N}(\sigma^{-2}A^{-1}X\mathbf{y},A^{-1}),$ 

where 
$$A = \sigma^{-2}X^TX + \Sigma^{-1}$$
.

▶ If the loss is not too complicated, then integrals are easy. For instance, prediction with  $L^2$  loss is simple:

$$\arg\min_{y_{\star}} \int \|y_{\star} - \theta^{\mathsf{T}} x_{\star}\|^{2} p(\theta|(x, y)_{1:n}) d\theta$$

## **Conjugacy**

Say we have a linear regression problem

$$y_i = f(x_i) + \varepsilon_i,$$

$$f(x) = \theta^T x$$
,  $\varepsilon_i$  i.i.d. Gaussians  $\mathcal{N}(0, \sigma^2)$ .

▶ If we choose  $p(\theta) \sim \mathcal{N}(0, \Sigma)$ , then (exercise)

$$p(\theta|(x,y)_{1:n}) \propto p((x,y)_{1:n}|\theta)p(\theta)$$
  
=  $\mathcal{N}(\sigma^{-2}A^{-1}X\mathbf{y},A^{-1}),$ 

where  $A = \sigma^{-2}X^TX + \Sigma^{-1}$ .

If the loss is not too complicated, then integrals are easy. For instance, prediction with  $L^2$  loss is simple:

$$\hat{\mathbf{y}}_{\star} := \sigma^{-2} \mathbf{x}_{\star}^{\mathsf{T}} \mathbf{A}^{-1} \mathbf{X} \mathbf{y}.$$

#### Monte Carlo methods

Sometimes, you're less lucky. Say we're doing logistic regression.

$$y_i = \text{Bernoulli}\left[\sigma(f(x_i))\right],$$

with 
$$f(x) = \theta^T x$$
,  $\sigma(x) = 1/(1 + e^{-x})$ .

▶ Even if we choose  $p(\theta) \sim \mathcal{N}(0, \Sigma)$ 

$$p(\theta|(x,y)_{1:n}) \propto p((x,y)_{1:n}|\theta)p(\theta)$$

does not have a simple closed form.

We need powerful numerical integration methods, that is, constructions of nodes  $(\theta_i)$  and weights  $w_i$  such that

$$\int h d\pi \approx \sum_{i=1}^{N} w_i h(\theta_i)$$

#### Monte Carlo methods

Sometimes, you're less lucky. Say we're doing logistic regression.

$$y_i = \text{Bernoulli}\left[\sigma(f(x_i))\right],$$

with 
$$f(x) = \theta^T x$$
,  $\sigma(x) = 1/(1 + e^{-x})$ .

▶ Even if we choose  $p(\theta) \sim \mathcal{N}(0, \Sigma)$ ,

$$p(\theta|(x,y)_{1:n}) \propto p((x,y)_{1:n}|\theta)p(\theta)$$

does not have a simple closed form.

We need powerful numerical integration methods, that is, constructions of nodes  $(\theta_i)$  and weights  $w_i$  such that

$$\int h d\pi \approx \sum_{i=1}^{N} w_i h(\theta_i)$$

#### Monte Carlo methods

Sometimes, you're less lucky. Say we're doing logistic regression.

$$y_i = \text{Bernoulli}\left[\sigma(f(x_i))\right],$$

with 
$$f(x) = \theta^T x$$
,  $\sigma(x) = 1/(1 + e^{-x})$ .

▶ Even if we choose  $p(\theta) \sim \mathcal{N}(0, \Sigma)$ ,

$$p(\theta|(x,y)_{1:n}) \propto p((x,y)_{1:n}|\theta)p(\theta)$$

does not have a simple closed form.

▶ We need powerful numerical integration methods, that is, constructions of nodes  $(\theta_i)$  and weights  $w_i$  such that

$$\int hd\pi \approx \sum_{i=1}^{N} w_i h(\theta_i).$$

```
MH(\pi(\theta), q(\theta'|\theta), \theta_0, N_{iter})
                   for k \leftarrow 1 to N_{\text{iter}}
```

```
MH(\pi(\theta), q(\theta'|\theta), \theta_0, N_{iter})
                      for k \leftarrow 1 to N_{\text{iter}}
                                   \theta \leftarrow \theta_{k-1}
                                  \theta' \sim q(.|\theta), \ u \sim \mathcal{U}_{(0.1)},
```

```
MH(\pi(\theta), q(\theta'|\theta), \theta_0, N_{iter})
                            for k \leftarrow 1 to N_{\text{iter}}
                                           \theta \leftarrow \theta_{k-1}
                                           \theta' \sim q(.|\theta), \ u \sim \mathcal{U}_{(0.1)},
                                          \alpha = \frac{\pi(\theta')}{\pi(\theta)} \frac{q(\theta|\theta')}{q(\theta'|\theta)}
```

```
MH(\pi(\theta), q(\theta'|\theta), \theta_0, N_{iter})
                         for k \leftarrow 1 to N_{\text{iter}}
                                       \theta \leftarrow \theta_{k-1}
                                      \theta' \sim q(.|\theta), \ u \sim \mathcal{U}_{(0,1)},
                                      \alpha = \frac{\pi(\theta')}{\pi(\theta)} \frac{q(\theta|\theta')}{q(\theta'|\theta)}
      5
                                      if u < \alpha
                                                   \theta_k \leftarrow \theta' \qquad \triangleright Accept
      6
                                       else \theta_k \leftarrow \theta \triangleright Reject
```

```
MH(\pi(\theta), q(\theta'|\theta), \theta_0, N_{iter})
                        for k \leftarrow 1 to N_{\text{iter}}
                                      \theta \leftarrow \theta_{k-1}
                                     \theta' \sim q(.|\theta), \ u \sim \mathcal{U}_{(0,1)},
                                     \alpha = \frac{\pi(\theta')}{\pi(\theta)} \frac{q(\theta|\theta')}{q(\theta'|\theta)}
     5
                                     if u < \alpha
                                                  \theta_k \leftarrow \theta' \qquad \triangleright Accept
     6
                                      else \theta_k \leftarrow \theta \triangleright Reject
                        return (\theta_k)_{k=1,\dots,N_{\text{iter}}}
     8
```

## The MCMC magic

Under assumptions<sup>5</sup>,

$$\sqrt{\textit{N}_{\text{iter}}} \left[ \frac{1}{\textit{N}_{\text{iter}}} \sum_{k=0}^{\textit{N}_{\text{iter}}} \textit{h}(\theta_k) - \int \textit{h}\left(\theta\right) \pi(\theta) d\theta \right] \rightarrow \mathcal{N}(0, \sigma_{\text{lim}}^2(\textit{h})),$$

- If you choose q carefully, you can hope for a polynomial increase of the mixing time and  $\sigma_{\lim}^2(h)$  with d.
- Most MCMC algorithms are instances of Metropolis-Hastings with clever choices of proposal<sup>6</sup>, even the NUTS HMC of Stan and PyMC3.
- ► For nice illustrations, check out https://chi-feng.github.io/mcmc-demo/

<sup>&</sup>lt;sup>5</sup>R. Douc et al. *Nonlinear time series*. Chapman-Hall, 2014.

<sup>&</sup>lt;sup>6</sup>C. P. Robert and G. Casella. *Monte Carlo Statistical Methods*. New York: Springer-Verlag, 2004.

## **Variational approximations**

When in a hurry, you can settle for a good approximation to your posterior

$$\pi(\theta) = p(\theta|\mathbf{x}) \propto p(\mathbf{x}|\theta)p(\theta),$$

say minimize in q

$$\begin{aligned} \mathsf{KL}(q,\pi) &= \mathbb{E}_q \log q - \mathbb{E}_q \log p(\theta|\mathbf{x}) \\ &= -[-\mathbb{E}_q \log q + \mathbb{E}_q \log p(\mathbf{x},\theta)] + \log p(\mathbf{x}). \end{aligned}$$

- Equivalently, we can maximize the evidence lower bound (ELBO)<sup>7</sup>.
- ▶ Ideally, I would rather cast the choice of *q* into a Wald-like problem.

<sup>&</sup>lt;sup>7</sup>D. M. Blei et al. "Variational inference: A review for statisticians". In: Journal of the American Statistical Association 112.518 (2017), pp. 859–877.

#### **Outline**

#### The what

Typical statistical problems Statistical decision theory Posterior expected utility and Bayes rules

## The why

The philosophical why The practical why

#### The how

Conjugacy Monte Carlo methods Metropolis-Hastings Variational approximations

# In depth with Gaussian processes in ML

From linear regression to GPs Modeling and learning More applications References and open issues

- $ightharpoonup y_i = f(x_i) + \varepsilon_i, \ f(x) = \theta^T x, \ \varepsilon_i \ \text{i.i.d.} \ \text{Gaussians} \ \mathcal{N}(0, \sigma^2).$
- ▶ If we choose  $p(\theta) \sim \mathcal{N}(0, \Sigma^2)$ , then

$$p(\theta|(x,y)_{1:n}) \propto p((x,y)_{1:n}|\theta)p(\theta)$$
= [Exercise]

where 
$$A = \sigma^{-2}X^TX + \Sigma^{-1}$$
.

- $y_i = f(x_i) + \varepsilon_i$ ,  $f(x) = \theta^T x$ ,  $\varepsilon_i$  i.i.d. Gaussians  $\mathcal{N}(0, \sigma^2)$ .
- ▶ If we choose  $p(\theta) \sim \mathcal{N}(0, \Sigma^2)$ , then

$$p(\theta|(x,y)_{1:n}) \propto p((x,y)_{1:n}|\theta)p(\theta)$$

$$\propto \exp\left[-\frac{\|\mathbf{y} - X\theta\|^2}{2\sigma^2} - \frac{\theta^T \Sigma^{-1}\theta}{2}\right]$$

- $y_i = f(x_i) + \varepsilon_i$ ,  $f(x) = \theta^T x$ ,  $\varepsilon_i$  i.i.d. Gaussians  $\mathcal{N}(0, \sigma^2)$ .
- ▶ If we choose  $p(\theta) \sim \mathcal{N}(0, \Sigma^2)$ , then

$$p(\theta|(x,y)_{1:n}) \propto p((x,y)_{1:n}|\theta)p(\theta)$$

$$\propto \exp\left[-\frac{\|\mathbf{y} - X\theta\|^2}{2\sigma^2} - \frac{\theta^T \Sigma^{-1}\theta}{2}\right]$$

$$\propto \exp\left[\frac{\mathbf{y}^T X\theta}{\sigma^2} - \frac{\theta^T X^T X\theta}{2\sigma^2} - \frac{\theta^T \Sigma^{-1}\theta}{2}\right]$$

- $ightharpoonup y_i = f(x_i) + \varepsilon_i, \ f(x) = \theta^T x, \ \varepsilon_i \ \text{i.i.d.} \ \text{Gaussians} \ \mathcal{N}(0, \sigma^2).$
- ▶ If we choose  $p(\theta) \sim \mathcal{N}(0, \Sigma^2)$ , then

$$p(\theta|(x,y)_{1:n}) \propto p((x,y)_{1:n}|\theta)p(\theta)$$

$$\propto \exp\left[-\frac{\|\mathbf{y} - X\theta\|^2}{2\sigma^2} - \frac{\theta^T \Sigma^{-1}\theta}{2}\right]$$

$$\propto \exp\left[\frac{\mathbf{y}^T X\theta}{\sigma^2} - \frac{\theta^T X^T X\theta}{2\sigma^2} - \frac{\theta^T \Sigma^{-1}\theta}{2}\right]$$

$$\propto \exp\left[-\frac{1}{2}\left(\theta - \sigma^{-2}A^{-1}X\mathbf{y}\right)^T A\left(\theta - \sigma^{-2}A^{-1}X\mathbf{y}\right)\right]$$

- $y_i = f(x_i) + \varepsilon_i$ ,  $f(x) = \theta^T x$ ,  $\varepsilon_i$  i.i.d. Gaussians  $\mathcal{N}(0, \sigma^2)$ .
- ▶ If we choose  $p(\theta) \sim \mathcal{N}(0, \Sigma^2)$ , then

$$p(\theta|(x,y)_{1:n}) \propto p((x,y)_{1:n}|\theta)p(\theta)$$

$$\propto \exp\left[-\frac{\|\mathbf{y} - X\theta\|^2}{2\sigma^2} - \frac{\theta^T \Sigma^{-1}\theta}{2}\right]$$

$$\propto \exp\left[\frac{\mathbf{y}^T X\theta}{\sigma^2} - \frac{\theta^T X^T X\theta}{2\sigma^2} - \frac{\theta^T \Sigma^{-1}\theta}{2}\right]$$

$$\propto \exp\left[-\frac{1}{2}\left(\theta - \sigma^{-2}A^{-1}X\mathbf{y}\right)^T A\left(\theta - \sigma^{-2}A^{-1}X\mathbf{y}\right)\right]$$

$$= \mathcal{N}(\sigma^{-2}A^{-1}X\mathbf{y}, A^{-1}),$$

- $y_i = f(x_i) + \varepsilon_i$ ,  $f(x) = \theta^T x$ ,  $\varepsilon_i$  i.i.d. Gaussians  $\mathcal{N}(0, \sigma^2)$ .
- ▶ If we choose  $p(\theta) \sim \mathcal{N}(0, \Sigma^2)$ , then

$$\begin{split} p(\theta|(x,y)_{1:n}) &\propto p((x,y)_{1:n}|\theta)p(\theta) \\ &\propto \exp\left[-\frac{\|\mathbf{y}-X\theta\|^2}{2\sigma^2} - \frac{\theta^T \Sigma^{-1}\theta}{2}\right] \\ &\propto \exp\left[\frac{\mathbf{y}^T X \theta}{\sigma^2} - \frac{\theta^T X^T X \theta}{2\sigma^2} - \frac{\theta^T \Sigma^{-1}\theta}{2}\right] \\ &\propto \exp\left[-\frac{1}{2}\left(\theta - \sigma^{-2}A^{-1}X\mathbf{y}\right)^T A\left(\theta - \sigma^{-2}A^{-1}X\mathbf{y}\right)\right] \\ &= \mathcal{N}(\sigma^{-2}A^{-1}X\mathbf{y}, A^{-1}), \end{split}$$

where  $A = \sigma^{-2}X^TX + \Sigma^{-1}$ .

▶ Remember prediction with  $L^2$  loss is simple:

$$\arg \min_{y_{\star}} \int \|y_{\star} - f(x_{\star})\|^{2} p(\theta|(x, y)_{1:n}) d\theta = \sigma^{-2} x_{\star}^{T} A^{-1} X \mathbf{y}.$$

- $y_i = f(x_i) + \varepsilon_i$ ,  $f(x) = \theta^T x$ ,  $\varepsilon_i$  i.i.d. Gaussians  $\mathcal{N}(0, \sigma^2)$ .
- ▶ If we choose  $p(\theta) \sim \mathcal{N}(0, \Sigma^2)$ , then

$$\begin{split} \rho(\theta|(x,y)_{1:n}) &\propto \rho((x,y)_{1:n}|\theta)\rho(\theta) \\ &\propto \exp\left[-\frac{\|\mathbf{y}-X\theta\|^2}{2\sigma^2} - \frac{\theta^T\Sigma^{-1}\theta}{2}\right] \\ &\propto \exp\left[\frac{\mathbf{y}^TX\theta}{\sigma^2} - \frac{\theta^TX^TX\theta}{2\sigma^2} - \frac{\theta^T\Sigma^{-1}\theta}{2}\right] \\ &\propto \exp\left[-\frac{1}{2}\left(\theta - \sigma^{-2}A^{-1}X\mathbf{y}\right)^TA\left(\theta - \sigma^{-2}A^{-1}X\mathbf{y}\right)\right] \\ &= \mathcal{N}(\sigma^{-2}A^{-1}X\mathbf{y}, A^{-1}), \end{split}$$

where  $A = \sigma^{-2}X^TX + \Sigma^{-1}$ .

► Actually, we can even check that

$$p(f(x_{\star})|x_{\star},(x,y)_{1:n}) = \mathcal{N}(\sigma^{-2}x_{\star}^{T}A^{-1}X\mathbf{y},x_{\star}^{T}A^{-1}x_{\star}^{\star}).$$

# Linear regression with nonlinear features

▶ Replace each x by a vector of features  $\varphi(x) \in \mathbb{R}^p$ :

$$y_i = f(x_i) + \varepsilon_i, \quad i = 1, \ldots, n,$$

$$f(x) = \theta^T \varphi(x)$$
,  $\varepsilon_i$  i.i.d. Gaussians  $\mathcal{N}(0, \sigma^2)$ ,  $\theta \sim \mathcal{N}(0, \Sigma)$ .

- ► Think  $\varphi(x) = (1, x^1, x^2, x^1x^2, ...)$
- Recall

$$p(f(x_{\star})|x_{\star},(x,y)_{1:n}) = \mathcal{N}(\sigma^{-2}\Phi_{\star}^{\mathsf{T}}A^{-1}\mathbf{\Phi}\mathbf{y},\Phi^{\star}A^{-1}\Phi^{\star})$$

where  $A = \sigma^{-2} \mathbf{\Phi}^T \mathbf{\Phi} + \Sigma^{-1}$ .

- ▶ Requires  $p \times p$  inversion.
- ▶ But let  $K = ΦΣΦ^T$ , then can rewrite (Exercise)

$$p(f(x_{\star})|(x,y)_{1:n}) = \mathcal{N}(\mu_{\star},\sigma_{\star}^{2}),$$

where

$$\mu_{\star} = \Phi_{\star}^{T} \Sigma \mathbf{\Phi}^{T} (\mathbf{K} + \sigma^{2} I)^{-1} \mathbf{y},$$
  
$$\sigma_{\star}^{2} = \varphi_{\star} \Sigma \varphi_{\star} - \varphi_{\star}^{T} \Sigma \mathbf{\Phi}^{T} (\mathbf{K} + \sigma^{2} I)^{-1} \mathbf{\Phi} \Sigma \varphi_{\star}.$$

## Gaussian processes

▶ A distribution over a space of functions  $f : \mathbb{R}^d \to \mathbb{R}$ .

### **Gaussian processes**

If  $\forall p \in \mathbb{N}, \forall x_1, \dots, x_p \in \mathbb{R}^d$ 

$$[f(x_1),\ldots,f(x_p)]^T\sim \mathcal{N}(\mathbf{m},\mathbf{K}),$$

where  $\mathbf{m} = [\mu(x_1), \dots, \mu(x_p)]$  and

$$\mathbf{K} = ((K(x_i, x_j))),$$

then we say  $f \sim GP(\mu, K)$ .

- Unicity is usually easy, existence is tricky.
- ▶ Necessary condition is that all matrices **K** are positive definite.

# Sampling, conditioning and predicting

See notebook 01 on https://github.com/rbardenet/bnp-course

## Commonly-used kernels

covariance function	expression	S	ND
constant	$\sigma_0^2$		
linear	$egin{array}{l} \sum_{d=1}^D \sigma_d^2 x_d x_d' \ (\mathbf{x} \cdot \mathbf{x}' + \sigma_0^2)^p \end{array}$		
polynomial	$(\mathbf{x} \cdot \mathbf{x}' + \sigma_0^2)^p$		
squared exponential	$\exp(-\frac{r^2}{2\ell^2})$		$\checkmark$
Matérn	$\left  \begin{array}{c} rac{1}{2^{ u-1}\Gamma( u)} \left(rac{\sqrt{2 u}}{\ell}r ight)^ u K_ u \left(rac{\sqrt{2 u}}{\ell}r ight) \end{array}  ight $		$\checkmark$
exponential	$\exp(-\frac{r}{\ell})$		$\checkmark$
$\gamma$ -exponential	$\exp\left(-\left(\frac{r}{\ell}\right)^{\gamma}\right)$		$\checkmark$
rational quadratic	$\left(1 + \frac{r^2}{2\alpha\ell^2}\right)^{-\alpha}$		$\checkmark$
neural network	$\sin^{-1}\left(\frac{2\tilde{\mathbf{x}}^{\top}\Sigma\tilde{\mathbf{x}}'}{\sqrt{(1+2\tilde{\mathbf{x}}^{\top}\Sigma\tilde{\mathbf{x}})(1+2\tilde{\mathbf{x}}'^{\top}\Sigma\tilde{\mathbf{x}}')}}\right)$		$\checkmark$

Table 4.1: Summary of several commonly-used covariance functions. The covariances are written either as a function of  $\mathbf{x}$  and  $\mathbf{x}'$ , or as a function of  $r = |\mathbf{x} - \mathbf{x}'|$ . Two columns marked 'S' and 'ND' indicate whether the covariance functions are stationary and nondegenerate respectively. Degenerate covariance functions have finite rank, see section 4.3 for more discussion of this issue.

# Learning

In regression,

$$p(\mathbf{y}|\mathbf{x}) = \int p(\mathbf{y}|\mathbf{f})p(\mathbf{f}|\mathbf{x})d\mathbf{f}$$
$$= \mathcal{N}(\mathbf{y}|0, \mathbf{K} + \sigma^2 I_n).$$

- So simply put a prior over  $\eta = (\sigma, \theta)$  and integrate
- Prediction becomes

$$f_{\star} \sim \int p(f_{\star}|\mathbf{x},\eta)p(\eta)d\eta$$

► Alternately, lots of people maximize the marginal likelihood.

# Learning

In regression,

$$p(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta}) = \int p(\mathbf{y}|\mathbf{f})p(\mathbf{f}|\mathbf{x})d\mathbf{f}$$
$$= \mathcal{N}(\mathbf{y}|0, \mathbf{K}_{\boldsymbol{\theta}} + \sigma^2 I_n).$$

- ▶ So simply put a prior over  $\eta = (\sigma, \theta)$  and integrate.
- Prediction becomes

$$f_{\star} \sim \int p(f_{\star}|\mathbf{x},\eta)p(\eta)d\eta.$$

▶ Alternately, lots of people maximize the marginal likelihood.

# Beyond regression: classification<sup>8</sup>

- (Exercise) Find a simple classification model with GPs.
- ► Take for instance

$$p(y = +1|x, f) = Bernoulli(\sigma(f(x))), \quad f \sim GP(0, K).$$

Problem: prediction is not easy anymore

$$p(f_*|X,\mathbf{y},x_*) = \int p(f_*|X,\mathbf{f},x_*)p(\mathbf{f}|X,\mathbf{y})d\mathbf{f}$$

<sup>&</sup>lt;sup>8</sup>C. E. Rasmussen and C. K. I. Williams. *Gaussian Processes for Machine Learning*. MIT Press, 2006.

# Beyond regression: classification<sup>8</sup>

- ► (Exercise) Find a simple classification model with GPs.
- Take for instance

$$p(y = +1|x, f) = Bernoulli(\sigma(f(x))), \quad f \sim GP(0, K).$$

Problem: prediction is not easy anymore

$$p(f_*|X,\mathbf{y},x_*) = \int p(f_*|X,\mathbf{f},x_*)p(\mathbf{f}|X,\mathbf{y})d\mathbf{f}$$

<sup>&</sup>lt;sup>8</sup>Rasmussen and Williams, Gaussian Processes for Machine Learning.

# Beyond regression: classification<sup>8</sup>

- ► (Exercise) Find a simple classification model with GPs.
- ▶ Take for instance

$$p(y = +1|x, f) = Bernoulli(\sigma(f(x))), \quad f \sim GP(0, K).$$

▶ Problem: prediction is not easy anymore

$$p(f_{\star}|X,\mathbf{y},x_{\star}) = \int p(f_{\star}|X,\mathbf{f},x_{\star})p(\mathbf{f}|X,\mathbf{y})d\mathbf{f}$$

<sup>&</sup>lt;sup>8</sup>Rasmussen and Williams, Gaussian Processes for Machine Learning.

# Beyond regression: ranking<sup>9</sup>

- ▶ (Exercise) Find a simple ranking model with GPs: your data is  $(u, v)_{1:n}$  where  $\forall i, u_i \prec v_i$ . Your user wants to know whether a new  $u_{\star} \prec v_{\star}$ .
- ► Take for instance

$$p(u \prec v | u, v, f) = \varphi(f(v) - f(u)), \quad f \sim GP(0, K), \varphi \text{ increasing}$$

Same difficulties with learning.

<sup>&</sup>lt;sup>9</sup>W. Chu and Z. Ghahramani. "Preference learning with Gaussian processes". In: *Proceedings of the 22nd International Conference on Machine Learning*. 2005, pp. 137–144.

# Beyond regression: ranking<sup>9</sup>

- ▶ (Exercise) Find a simple ranking model with GPs: your data is  $(u, v)_{1:n}$  where  $\forall i, u_i \prec v_i$ . Your user wants to know whether a new  $u_{\star} \prec v_{\star}$ .
- Take for instance

$$p(u \prec v | u, v, f) = \varphi(f(v) - f(u)), \quad f \sim GP(0, K), \varphi \text{ increasing.}$$

► Same difficulties with learning.

<sup>&</sup>lt;sup>9</sup>Chu and Ghahramani, "Preference learning with Gaussian processes".

# Beyond regression: ranking<sup>9</sup>

- ▶ (Exercise) Find a simple ranking model with GPs: your data is  $(u, v)_{1:n}$  where  $\forall i, u_i \prec v_i$ . Your user wants to know whether a new  $u_{\star} \prec v_{\star}$ .
- ► Take for instance

$$p(u \prec v | u, v, f) = \varphi(f(v) - f(u)), \quad f \sim GP(0, K), \varphi \text{ increasing.}$$

► Same difficulties with learning.

<sup>&</sup>lt;sup>9</sup>Chu and Ghahramani, "Preference learning with Gaussian processes".

#### **Emulators of expensive models**



RESEARCH ARTICLE

## Bayesian Sensitivity Analysis of a Cardiac Cell Model Using a Gaussian Process Emulator

Eugene TY Chang<sup>1,2</sup>, Mark Strong<sup>3</sup>, Richard H Clayton<sup>1,2</sup>\*

1 Insigneo Institute for in-silico Medicine, University of Sheffield, Sheffield, United Kingdom, 2 Department of Computer Science University of Sheffield, Sheffield, United Kingdom, 3 School of Health and Related Research, University of Sheffield, Sheffield, United Kingdom

\* r.h.clayton@sheffield.ac.uk



OPEN ACCESS

Citation: Chang ETY, Strong M, Clayton RH (2015)

#### Abstract

Models of electrical activity in cardiac cells have become important research tools as they can provide a quantitative description of detailed and integrative physiology. However, cardiac cell models have many parameters, and how uncertainties in these parameters affect the model output is difficult to assess without undertaking large numbers of model runs. In this study we show that a surrogate statistical model of a cardiac cell model (the Luo-Rudy 1991 model) can be half under Caussian process (GP) emiliary. Islan this paromach we

### Nonparametric fits

Arman Shafieloo<sup>1</sup>, Alex G. Kim<sup>2</sup>, Eric V. Linder<sup>1,2,3</sup> Institute for the Early Universe WCU, Ewha Womans University, Seoul, Korea <sup>2</sup> Lawrence Berkeley National Laboratory, Berkeley, CA 94720, USA and <sup>3</sup> University of California, Berkeley, CA 94720, USA (Dated: July 11, 2012).

Gaussian Process Cosmography

Gaussian processes provide a method for extracting cosmological information from observations without assuming a cosmological model. We carry out cosmography - mapping the time evolution of the cosmic expansion - in a model-independent manner using kinematic variables and a geometric probe of cosmology. Using the state of the art supernova distance data from the Union 2.1 compilation, we constrain, without any assumptions about dark energy parametrization or matter density, the Hubble parameter and deceleration parameter as a function of redshift. Extraction of these relations is tested successfully against models with features on various coherence scales, subject to certain statistical cautions.

#### I. INTRODUCTION

Cosmic acceleration is a fundamental mystery of great interest and importance to understanding cosmology, gravitation, and high energy physics. The cosmic expansion rate is slowed down by gravitationally attractive matter and sped up by some other, unknown contribution to the dynamical equations. While great effort is being put into identifying the source of this extra dark energy contribution, the overall expansion behavior also holds important clues to origin, evolution, and present

ing procedures have been suggested, e.g. [6], but tend to induce bias in the function reconstruction due to parametric restriction of the behavior or to have poor error control. Using a general orthonormal basis or principal component analysis is another approach, to describe the distance-redshift relation (e.g. [7]) or the deceleration parameter [8], or using a correlated prior for smoothness on the dark energy equation of state [9], but in practice a finite (and small) number of modes is significant beyond the prior, essentially reducing to a parametric approach, Gaussian processes [10] offer an interesting possibility for

ro-ph.CO] 10 Jul 2012

### Natural language processing

# Using Gaussian Processes for Rumour Stance Classification in Social Media

MICHAL LUKASIK, University of Sheffield KALINA BONTCHEVA, University of Sheffield TREVOR COHN, University of Melbourne ARKAITZ ZUBIAGA, University of Warwick MARIA LIAKATA, University of Warwick ROB PROCTER, University of Warwick

Social media tend to be rife with rumours while new reports are released piecemeal during breaking news. Interestingly, once an mine multiple reactions expressed by social media users in those situations, exploring their stance towards turns, ultimately enabling the flagging of highly disputed rumours as being potentially false. In this work, we set out to develop an automated, supervised classifier that uses multi-task learning to classify the stance expressed in each individual tweet in a rumourous conversation as either supporting, denying or questioning the rumour. Using a classifier based on Gaussian Processes, and exploring its effectiveness on two datasets with very different characteristics and varying distributions of stances, we show that our approach consistently outperforms competitive baseline classifiers. Our classifier is sepecially effective in estimating the distribution of different types of stance associated with a given rumour, which we set forth as a desired characteristic for a rumour-tasking system that will warn both ordinary users of Twitter and professional news practitioners when a rumour is being rebutted.

#### 1. INTRODUCTION

There is an increasing need to interpret and act upon rumours spreading quickly through social media during breaking news, where new reports are released piecemeal and often have an unverified

s.CL] 7 Sep 2016

## Bayesian optimization for hyperparameter tuning

#### **Algorithms for Hyper-Parameter Optimization**

#### James Bergstra The Rowland Institute Harvard University

bergstra@rowland.harvard.edu

#### Yoshua Bengio

Dépt. d'Informatique et Recherche Opérationelle Université de Montréal yoshua.bengio@umontreal.ca

#### Rémi Bardenet

Laboratoire de Recherche en Informatique Université Paris-Sud bardenet@lri.fr

#### Balázs Kégl

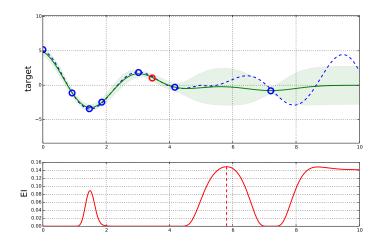
Linear Accelerator Laboratory Université Paris-Sud. CNRS balazs.kegl@gmail.com

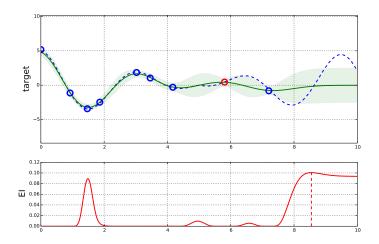
#### Abstract

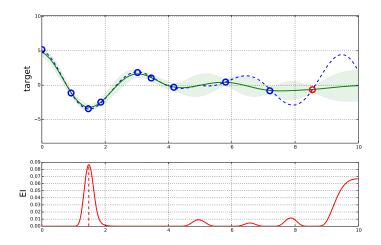
Several recent advances to the state of the art in image classification benchmarks have come from better configurations of existing techniques rather than novel approaches to feature learning. Traditionally, hyper-parameter optimization has been the job of humans because they can be very efficient in regimes where only a few

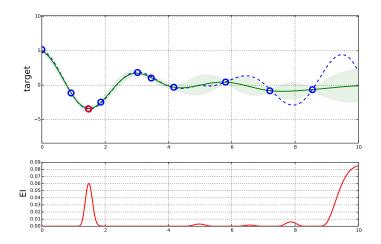
## **Bayesian optimization**

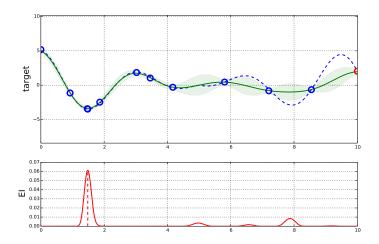
- $\triangleright$  Goal is to minimize a noisy f with N iterations, N small.
- ► Key application: find the hyperparameters of your ML algorithm that minimize the validation error.
- ► Idea is to sequentially
  - update your model on f,
  - optimize an aquisition criterion.
- Check out notebook 03 on https://github.com/rbardenet/bnp-course.

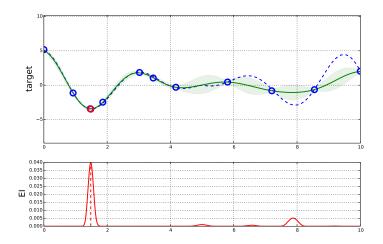


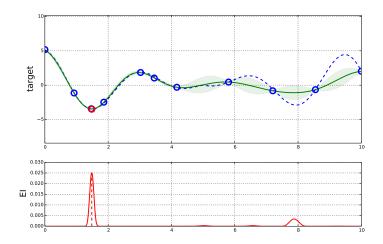


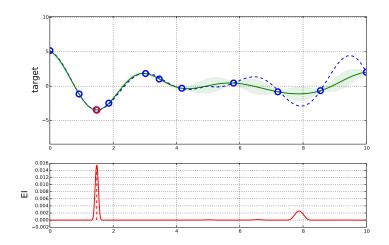


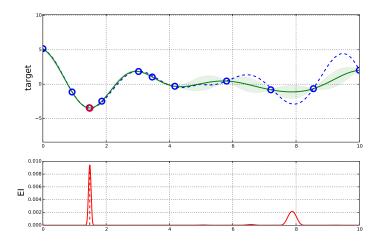


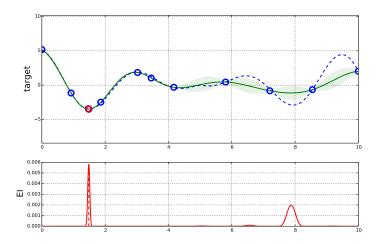












### Popular aquisition criteria

# Expected improvement<sup>10</sup>

$$\mathsf{El}(z) = \mathbb{E}\big(\max(m_N - f(z), 0) | (x, y)_{1:n}\big),$$

where  $m_N = \min_{1 \leq i \leq N} f(x_i)$ .

An easy computation yields

$$\mathsf{EI}(x) = \widetilde{\sigma}(x) \big( u \Phi(u) + \varphi(u) \big), \tag{3}$$

where

$$u = (m_n - \widetilde{m}(x))/\widetilde{\sigma}(x),$$

and  $\Phi$  and  $\varphi$  denote the cdf and pdf of the  $\mathcal{N}(0,1)$  distribution.

<sup>&</sup>lt;sup>10</sup>D. R. Jones. "A Taxonomy of Global Optimization Methods Based on Response Surfaces". In: *Journal of Global Optimization* 21 (2001), pp. 345–383.

## GP-UCB (Srinivas et al., 2010)

#### **GP-UCB**

$$\mathsf{GP\text{-}UCB}(z) = \widetilde{m}(z) + \beta \widetilde{\sigma}(z).$$

- $\triangleright$  If  $\beta$  properly tuned, can use bandit results.
- First criterion to give application theoretical results.

### Bayesian optimization for hyperparameter tuning

#### **Algorithms for Hyper-Parameter Optimization**

# James Bergstra The Rowland Institute Harvard University bergstra@rowland.harvard.edu

Rémi Bardenet

Laboratoire de Recherche en Informatique
Université Paris-Sud
bardenet@lri.fr

#### Yoshua Bengio Dépt. d'Informatique et Recherche Opérationelle Université de Montréal yoshua.bengio@umontreal.ca

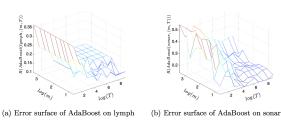
Balázs Kégl Linear Accelerator Laboratory Université Paris-Sud, CNRS balazs.kegl@gmail.com

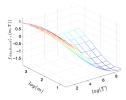
#### Abstract

Several recent advances to the state of the art in image classification benchmarks have come from better configurations of existing techniques rather than novel approaches to feature learning. Traditionally, hyper-parameter optimization has been the iob of humans because they can be very efficient in reeimes where only a few

Checkout hyperopt and spearmint.

# Going further: Hyperopt across datasets<sup>10</sup>





<sup>(</sup>c) The common latent ranker

<sup>&</sup>lt;sup>10</sup>R. Bardenet et al. "Collaborative hyperparameter tuning". In: *International Conference on Machine Learning (ICML)*. Atlanta, Georgia, 2013.

### Some useful hyperlinks

- Textbook by C. E. Rasmussen and C. K. I. Williams. Gaussian Processes for Machine Learning. MIT Press, 2006,
  - preat for understanding, methods, pointers to ML and stats.
- ▶ Videolecture by C. Rasmussen.
- ▶ lecture notes by P. Orbanz. *Lecture notes on Bayesian nonparametrics*. 2014.
  - mathematically clean, without losing the focus on ML.

#### Some open issues

- Fully Bayesian scalable approaches!
- Natural approaches to constrained GPs.
- Links with other models based on Gaussians and geometry.

#### Back to the roots

- Formulate HT across datasets and algorithms as a posterior expected loss problem, including computational constraints.
- Solve resulting dynamic programming problem.

#### References I

- Bardenet, R. et al. "Collaborative hyperparameter tuning". In: International Conference on Machine Learning (ICML). Atlanta, Georgia, 2013.
- Blei, D. M., A. Kucukelbir, and J. D. McAuliffe. "Variational inference: A review for statisticians". In: *Journal of the American Statistical Association* 112.518 (2017), pp. 859–877.
- Chu, W. and Z. Ghahramani. "Preference learning with Gaussian processes". In: *Proceedings of the 22nd International Conference on Machine Learning*. 2005, pp. 137–144.
- Douc, R., É. Moulines, and D. Stoffer. *Nonlinear time series*. Chapman-Hall, 2014.
- Jones, D. R. "A Taxonomy of Global Optimization Methods Based on Response Surfaces". In: *Journal of Global Optimization* 21 (2001), pp. 345–383.
- Orbanz, P. Lecture notes on Bayesian nonparametrics. 2014.

#### References II

- Parmigiani, G. and L. Inoue. *Decision theory: principles and approaches.* Vol. 812. John Wiley & Sons, 2009.
- Rasmussen, C. E. and C. K. I. Williams. *Gaussian Processes for Machine Learning*. MIT Press, 2006.
- Robert, C. P. The Bayesian choice: from decision-theoretic foundations to computational implementation. Springer Science & Business Media, 2007.
- Robert, C. P. and G. Casella. *Monte Carlo Statistical Methods*. New York: Springer-Verlag, 2004.
- Stigler, S. M. "The epic story of maximum likelihood". In: *Statistical Science* (2007), pp. 598–620.
- Wald, A. Statistical decision functions. Wiley, 1950.