

BML lecture #1: Bayesics

<http://github.com/rbardenet/bml-course>

Rémi Bardenet

remi.bardenet@gmail.com

CNRS & CRISTAL, Univ. Lille, France



Instructors

Besides me, the lecturers are

- ▶ Gabriel Victorino Cardoso (Mines Paris, PSL Univ.),
- ▶ Julyan Arbel (Inria Grenoble).



Course outline

I. What is Bayesian ML?

- ① Warmup: Linear regression
- ② Maximizing expected utility

II. How do we perform Bayesian ML?

- ① Monte Carlo methods
- ② Variational Bayes

III. Why do you want to apply BML?

Foundations

IV. Nonparametrics

V. Generative models in BML.

VI. Student seminar.

Homework: exercise sheet, practicals

Final 4x3

← projects start

Final
Julian 3hrs
Gabriel 2x3
1x2 hrs

- 1 A warmup: Estimation in regression models**
- 2 ML as data-driven decision-making**
- 3 Subjective expected utility**
- 4 Specifying joint models**
- 5 50 shades of Bayes**

- ▶ [...] practical methods for making inferences from data, using probability models for quantities we observe **and for quantities about which we wish to learn.**
- ▶ The essential characteristic of Bayesian methods is their **explicit use of probability for quantifying uncertainty** in inferences based on statistical data analysis.
- ▶ Three steps:
 - 1 Setting up a full probability model,
 - 2 Conditioning on observed data, calculating and interpreting the appropriate “posterior distribution”,
 - 3 Evaluating the fit of the model and the implications of the resulting posterior distribution. In response, one can alter or expand the model and repeat the three steps.

Notation that I will try to stick to

- ▶ $y_{1:n} = (y_1, \dots, y_n) \in \mathcal{Y}^n$ denote observable data/labels.
- ▶ $x_{1:n} \in \mathcal{X}^n$ denote covariates/features/hidden states.
- ▶ $z_{1:n} \in \mathcal{Z}^n$ denote hidden variables.
- ▶ $\theta \in \Theta$ denote parameters.
- ▶ X denotes an \mathcal{X} -valued random variable. Lowercase x denotes either a point in \mathcal{X} or an \mathcal{X} -valued random variable.

- ▶ Whenever it can easily be made formal, we write densities for our random variables and let the context indicate what is meant. So if $X \sim \mathcal{N}(0, \sigma^2)$, we write

$$\mathbb{E}h(X) = \int h(x) \frac{e^{-x^2/2\sigma^2}}{\sigma\sqrt{2\pi}} dx = \int h(x)p(x)dx.$$

Similarly, for $X \sim \mathcal{P}(\lambda)$, we write

$$\mathbb{E}h(X) = \sum_{k=0}^{\infty} h(k) e^{-\lambda} \frac{\lambda^k}{k!} = \int h(x)p(x)dx$$

- ▶ All pdfs are denoted by p , so that, e. g.

$$\begin{aligned}\mathbb{E}h(Y, \theta) &= \int h(y, \theta)p(y, \theta) dy d\theta \\ &= \int h(y, \theta)p(y, x, \theta) dx dy d\theta \\ &= \int h(y, \theta)p(y, \theta|x)p(x) dx dy d\theta\end{aligned}$$

- 1 A warmup: Estimation in regression models**
- 2 ML as data-driven decision-making**
- 3 Subjective expected utility**
- 4 Specifying joint models**
- 5 50 shades of Bayes**

1 A warmup: Estimation in regression models

2 ML as data-driven decision-making

3 Subjective expected utility

4 Specifying joint models

5 50 shades of Bayes

Inference in regression models

$$\begin{array}{c} \cancel{y} \\ | \\ \vdots \\ n \end{array} \quad y_i = x_i^T \theta + \varepsilon_i \quad \varepsilon_i \sim \text{iid } N(0, \sigma^2) \\ \underline{y = X\theta + \varepsilon}$$



Thm: $\hat{\theta}_{\text{MLE}}$ Eargmin $\mathcal{L}(y | x\theta, \sigma^2)$
 $\hat{\theta}_{\text{MLE}} \sim N(\theta, \sigma^2 (X^T X)^{-1})$

Proof: Assume $X^T X$ is full-rank

then $\hat{\theta}_{\text{MLE}} = (X^T X)^{-1} X^T y$. So

$$\begin{aligned} \mathbb{E} \hat{\theta}_{\text{MLE}} &= (X^T X)^{-1} X^T \mathbb{E}[y | \theta, X] \\ &= (X^T X)^{-1} X^T X \theta = \theta. \end{aligned}$$

$$V \hat{\theta}_{\text{MLE}} = (X^T X)^{-1} X^T \mathbb{V}[y | \theta, X] X (X^T X)^{-1}$$

Rq: $P_{y|X\theta} (\theta \in A_{\alpha}(x,y)) > 1 - \alpha$

Fisher (1st half
of 20th century)

$$\frac{(y - \hat{y})^T (X^T X)(\theta - \hat{\theta})}{\sigma^2} = \chi^2$$

Inference in regression models

$$* L(\theta, \hat{\theta}) = \| \theta - \hat{\theta} \|^2$$

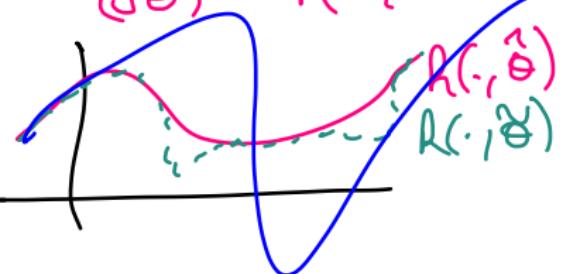
$$* R(\theta, \hat{\theta}) = \mathbb{E}_{y|x,\theta} L(\theta, \hat{\theta}(x,y)).$$

Def: $\hat{\theta}_{\text{minimax}} \in \arg\min_{\hat{\theta}} \sup_{\theta} R(\theta, \hat{\theta})$

$$\hat{\theta}_{\text{bayes}} \in \arg\min_{\theta} \mathbb{E}_{\text{emp}} R(\theta, \hat{\theta})$$

Def: $\hat{\theta}$ is inadmissible if $\exists \theta$

- (f) $R(\theta, \hat{\theta}) \leq R(\theta, \tilde{\theta})$
- (g) $R(\theta, \hat{\theta}) < R(\theta, \tilde{\theta})$



Wald + 1950

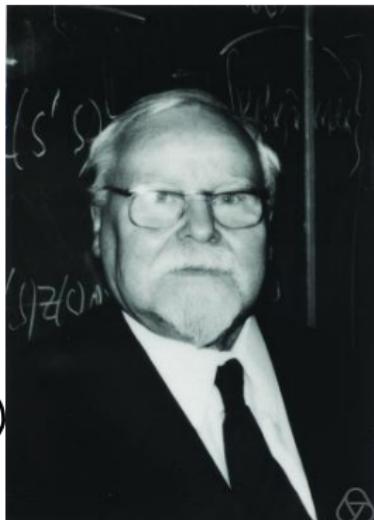
Theorem (ext. of James-Stein estimator)

$\hat{\theta}_{\text{MLE}}$ is inadmissible in linear regression.

Inference in regression models

$$\text{Let } \lambda > 0, \hat{\theta}_\lambda \in \underset{\theta}{\operatorname{argmin}} \|y - X\theta\|^2 + \lambda \|\theta\|^2 \\ = (X^T X + \lambda I)^{-1} X^T y. \quad (\#)$$

Proposition: $\hat{\theta}_\lambda$ s.t. ($\| \theta \| \leq B$)
 $R(\theta, \hat{\theta}_\lambda) < R(\theta, \hat{\theta}_{\text{MLE}})$.



Tibor
Hoerl & Kennard

$$\begin{aligned} \text{Proof: } R(\theta, \hat{\theta}_\lambda) &= \mathbb{E}_{y|X, \theta} \|\theta - \hat{\theta}_\lambda\|^2 \\ &= \|\theta - \mathbb{E}_{y|X, \theta} \hat{\theta}_\lambda\|^2 + \text{Tr } V \hat{\theta}_\lambda \quad (\#) \end{aligned}$$

* Assume for simplicity $X^T X = I$,

$$\mathbb{E} \hat{\theta}_\lambda = \frac{1}{1+\lambda} X^T \mathbb{E} y = \frac{1}{1+\lambda} \theta.$$

$$\text{Tr } V \hat{\theta}_\lambda = \frac{1}{(1+\lambda)^2} \text{Tr} [X^T V y X] = \frac{1}{(1+\lambda)^2} \sigma^2 d$$

$$\text{So } R(\theta, \hat{\theta}_\lambda) = \|\theta\|^2 \left(1 - \frac{1}{1+\lambda}\right)^2 + \frac{\sigma^2 d}{(1+\lambda)^2} \leq B^2 \frac{d^2}{(1+\lambda)^2} + \frac{\sigma^2 d}{(1+\lambda)^2}$$

Inference in regression models

$$\hat{\theta}_{\text{Bayes}} \in \arg \min_{\theta} \mathbb{E}_{\theta \sim p} \mathbb{E}_{y|x,\theta} L(\theta, \hat{\theta})$$

Proposition: Let $\theta \sim \mathcal{N}(0, \sigma^2 I)$ (*)

$$\text{Let } y|x, \theta \sim \mathcal{N}(x\theta, \sigma^2 I)$$

Then $\hat{\theta}_{\text{Bayes}} = \hat{\theta}_1 \text{ for } d = \dots$

$$\begin{aligned} \text{Proof: } & \mathbb{E}_{\theta} \mathbb{E}_{y|x,\theta} L(\theta, \hat{\theta}) = \\ & \int \| \theta - \hat{\theta} \|^2 p(y|x, \theta) p(\theta) dy d\theta = \\ & \int \left[\| \theta - \hat{\theta} \|^2 p(\theta|x, y) \right] p(y|x) dy \end{aligned}$$



Bayes (18th century)

Thus $\hat{\theta}_{\text{Bayes}} := \mathbb{E}(\theta | X, y) \in \arg \min_{\theta} \mathbb{E}_{\theta \sim p} R(\theta, \hat{\theta})$.

$$\text{Now, note that } p(\theta | X, y) = \frac{p(y, \theta | X)}{p(y | X)} = c(\theta) \left(\left(\frac{1}{\sigma^2} I + \frac{1}{n} X^T X \right)^{-1} X^T y \right) \dots$$

Inference in regression models

1 A warmup: Estimation in regression models

2 ML as data-driven decision-making

3 Subjective expected utility

4 Specifying joint models

5 50 shades of Bayes

Describing a decision problem under uncertainty

- ▶ A state space \mathcal{S} ,
Every quantity you need to consider to make your decision.
- ▶ Actions $\mathcal{A} \subset \mathcal{F}(\mathcal{S}, \mathcal{Z})$,
Making a decision means picking one of the available actions.
- ▶ A reward space \mathcal{Z} ,
Encodes how you feel about having picked a particular action.
- ▶ A loss function $L : \mathcal{A} \times \mathcal{S} \rightarrow \mathbb{R}_+$.
How much you would suffer from picking action a in state s .

Classification as a decision problem

- ▶ $\mathcal{S} = \mathcal{X}^n \times \mathcal{Y}^n \times \mathcal{X} \times \mathcal{Y}$, i.e. $s = (x_{1:n}, y_{1:n}, x, y)$.
- ▶ $\mathcal{Z} = \{0, 1\}$.
- ▶ $\mathcal{A} = \{a_g : s \mapsto 1_{y \neq g(x; \mathbf{x}_{1:n}, \mathbf{y}_{1:n})}, \quad g \in \mathcal{G}\}$.
- ▶ $L(a_g, s) = 1_{y \neq g(x; \mathbf{x}_{1:n}, \mathbf{y}_{1:n})}$.

PAC bounds; see e.g. (Shalev-Shwartz and Ben-David, 2014)

Let $(x_{1:n}, y_{1:n}) \sim \mathbb{P}^{\otimes n}$, and independently $(x, y) \sim \mathbb{P}$, we want an algorithm $g(\cdot; x_{1:n}, y_{1:n}) \in \mathcal{G}$ such that if $n \geq n(\delta, \varepsilon)$,

$$\mathbb{P}^{\otimes n} [\mathbb{E}_{(x,y) \sim \mathbb{P}} L(a_g, s) \leq \varepsilon] \geq 1 - \delta.$$

In particular, $\mathbb{E}_{x_{1:n}, y_{1:n}, x, y} L(a_g, s) \leq (1-\delta)\varepsilon + \delta$

- ▶ $S = \mathcal{X}^n \times \mathcal{Y}^n \times \mathcal{H} \times \mathcal{Y}$
- ▶ $Z = R$
- ▶ $A = \{a_g : s \mapsto y - g(x_1, x_2, \dots, x_n), g \in G\}$
- ▶ $L : a_g, s \mapsto \|y - g(x_1, x_2, \dots, x_n)\|^2$
 $\quad (= -u(a_g(s)))$

"Bayesian" Reason: for $n \gg n(\varepsilon, \delta)$,

$$\Pr_{\mathcal{P}}^{\otimes n} \left[\mathbb{E}_{\mathcal{P}} L(a_g, s) \leq \varepsilon \right] \geq 1 - \delta$$

Estimation as a decision problem (credible intervals).

- $\mathcal{S} = \mathbb{Y}^n \times \Theta$, $\Theta = \mathbb{R}$ for simplicity.
- $\mathcal{Z} = \{0,1\}^n \times \mathbb{R}_+$
- $\mathcal{A} = \{a_I : s \mapsto (\mathbb{1}_{\theta \notin I}, |I|)$
- $L : a_I s \mapsto \mathbb{1}_{\theta \notin I} + \delta |I|$ for some $\delta > 0$.

$$\mathbb{E}_{y_{j:n}|\theta} L(a_I, s) = \mathbb{E}_{y_{j:n}|\theta} \mathbb{1}_{\theta \notin I(y_{j:n})} + \delta \mathbb{E}_{y_{j:n}|\theta} |I(y_{j:n})|$$

Exo

- ▶ $\mathcal{S} =$
- ▶ $\mathcal{Z} =$
- ▶ $\mathcal{A} =$
- ▶

Density estimation as a decision problem

- ▶ $S = \mathbb{Y}^n \times \mathbb{Y}$
 - ▶ $Z =$
 - ▶ $A = \{a_{\pi, \ell(-1)}(y_{1:n}, y) = \int \pi(z) \ell(y|z) dz, \pi, \ell \in \mathcal{P}\}$
 - ▶ $L: a_s \mapsto -\log a_{\pi, \ell}(y_{1:n}, y)$
- $\#_{Y_{1:n}, Y} \left[-\log \int \pi(z|y_{1:n}) \ell(y|z, y_{1:n}) dz \right]$
is the risk.

1 A warmup: Estimation in regression models

2 ML as data-driven decision-making

3 Subjective expected utility

4 Specifying joint models

5 50 shades of Bayes

The subjective expected utility principle

- 1 Choose $\mathcal{S}, \mathcal{Z}, \mathcal{A}$ and a loss function $L(a, s)$,
- 2 Choose a distribution p over \mathcal{S} ,
- 3 Take the corresponding Bayes action

$$a^* \in \arg \min_{a \in \mathcal{A}} \mathbb{E}_{s \sim p} L(a, s). \quad (1)$$

Corollary: minimize the posterior expected loss

Now partition $s = (s_{\text{obs}}, s_u)$, then

$$a^* \in \arg \min_{a \in \mathcal{A}} \mathbb{E}_{s_{\text{obs}}} \mathbb{E}_{s_u | s_{\text{obs}}} L(a, s).$$

In ML, $\mathcal{A} = \{a_g\}$, with $g = g(s_{\text{obs}})$, so that (1) is equivalent to $a^* = a_{g^*}$, with

$$g^*(s_{\text{obs}}) \triangleq \arg \min_g \mathbb{E}_{s_u | s_{\text{obs}}} L(a, s).$$

1 A warmup: Estimation in regression models

2 ML as data-driven decision-making

3 Subjective expected utility

4 Specifying joint models

5 50 shades of Bayes

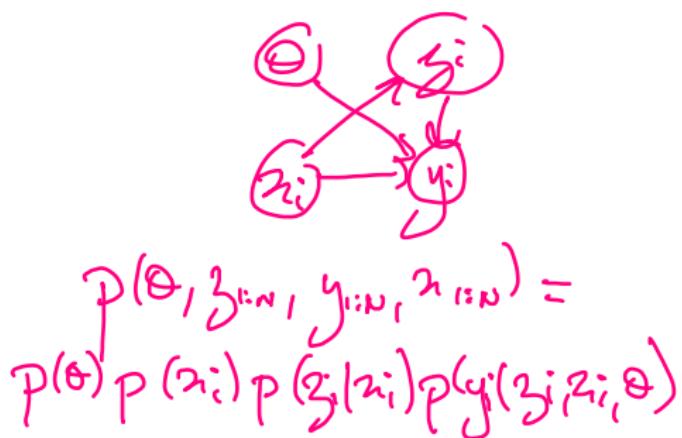
Pause: A recap on probabilistic graphical models

- ▶ Often we specify joint distributions through their conditionals.
- ▶ PGMs (aka “Bayesian” networks) represent the dependencies in a joint distribution $p(s)$ by a directed acyclic graph $G = (E, V)$,

$$p(s) = \prod_{v \in V} p(s_v | s_{\text{pa}(v)}).$$



$$p(\theta, y_{1:N}) = \prod_{i=1}^N p(y_i | \theta) p(\theta)$$



Pause: A recap on probabilistic graphical models

Check (Murphy, 2012, Section 10.5) for a refresher.

Theorem

Given a PGM $G = (V, E)$, and $A, B, C \subset V$, then $A \perp B | C$ iff every path from a node in A to a node in B is d -blocked by C .

d -blocking

An undirected path P in G is d -blocked by $E \subset V$ if at least one of the following conditions hold.

- ▶ P contains a “chain” $a \rightarrow b \rightarrow c$ and $b \in E$.
- ▶ P contains a “tent” $a \leftarrow b \rightarrow c$ and $b \in E$.
- ▶ P contains a “v-structure” $a \rightarrow b \leftarrow c$ and neither b nor any of its descendants are in E .



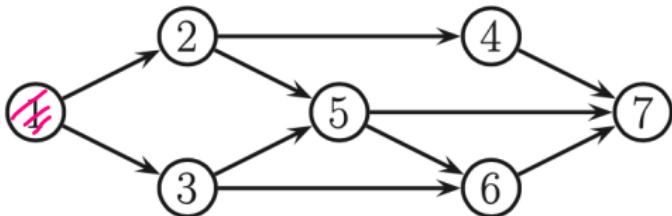


Figure 10.11 A DGM.

- ▶ Does $x_2 \perp x_6 | x_5, x_1$?
- ▶ Does $x_2 \perp x_6 | x_1$?
- ▶ Write the joint distribution as factorized over the graph.

$$p(x_1)p(x_2|x_1)p(x_3|x_1)p(x_4|x_2)p(x_5|x_2, x_3)p(x_6|x_5, x_3)p(x_7|x_4, x_6, x_5)$$

Estimation as a decision problem: point estimates

Estimation as a decision problem: credible intervals

Choosing priors (see Exercises)

Classification as a decision problem

Regression as a decision problem 1/2

Regression as a decision problem 2/2

Dimensionality reduction as a decision problem

Clustering as a decision problem

Topic modelling as a decision problem

- 1 A warmup: Estimation in regression models**
- 2 ML as data-driven decision-making**
- 3 Subjective expected utility**
- 4 Specifying joint models**
- 5 50 shades of Bayes**

An issue (or is it?)

Depending on how they interpret and how they implement SEU, you will meet many types of Bayesians (46656, according to Good).

A few divisive questions

- ▶ Using data or the likelihood to choose your prior; see Lecture 4.
- ▶ Using MAP estimators for their computational tractability, like in inverse problems

$$\hat{x}_\lambda \in \arg \min \|y - Ax\| + \lambda \Omega(x).$$

- ▶ When and how should you revise your model (likelihood or prior)?
- ▶ MCMC vs variational Bayes (more in Lectures #2 and #3)

References I

- [1] A. Gelman et al. *Bayesian data analysis*. 3rd. CRC press, 2013.
- [2] K. Murphy. *Machine learning: a probabilistic perspective*. MIT Press, 2012.
- [3] S. Shalev-Shwartz and S. Ben-David. *Understanding machine learning: From theory to algorithms*. Cambridge university press, 2014.