

WHAT IS BAYESIAN INFERENCE, WHEN DO I USE IT, AND WHY?

WHAT?

A PRINC. LEARN. FRAMEWORK (MLE, ERM, ...)
BASED ON HAVING A PROB OPINION & UPDATING IT

$$P(\theta|y) = \frac{P(y|\theta)P(\theta)}{P(y)}$$

WHEN?

ALWAYS! (COMP. INTENSIBLE, LAZY)

WHY?

EASIER:

⊖ COMPUTATION

- TO INTERPRET
- TO INCLUDE PRIOR KNOWLEDGE
- DECISION THEORY (!)

⊖ PRIORS

↑
LESS OF A CONCERN

- TO NOT OVERFIT (NO FITTING \Rightarrow NO OVERFITTING)

EXAMPLE: COVID-19 NATIONAL SURVEY

$$n = 1318, y_i \in \{0, 1\}, \sum_{i=1}^n y_i = 41$$

$$\hat{\theta} = \frac{\sum y_i}{n} = \frac{41}{1318} \approx 0.031 = 3.1\% \quad \times \text{ QUANTIFY UNCERTAINTY!}$$

- THIS IS A STAT. TASK
(QUANT. UNC. & MAKING DECISIONS)
- WE NEED A STAT. MODEL (A PROBABILISTIC MODEL)
(WE'LL OPT FOR A PARAMETRIC MODEL)

A PARAMETRIC MODEL FOR SUCCESS RATE

$$\theta \in [0, 1] \quad y_i | \theta = \begin{cases} 1, & \text{WITH } \theta \text{ PROB.} \\ 0, & \text{WITH } 1-\theta \text{ PROB.} \end{cases}$$

OR

$$P(y_i = k | \theta) = \theta^k (1-\theta)^{1-k}$$

OR

$$y_i | \theta \sim \text{Bernoulli}(\theta)$$

(LIKELIHOOD)

MAX. LIK. VIEW

- y_i ARE RANDOM SAMPLES FROM DGP (WE COULD SAMPLE MORE)
 \nwarrow R.V.S
- θ IS AN UNKNOWN (BUT CONSTANT)
- FIND $\hat{\theta}$ THAT MAX. THE LIKELIHOOD

$$\rightarrow \hat{\theta} = \underset{\theta}{\operatorname{argmax}} \prod_{i=1}^n p(y_i | \theta) = 3.1\%$$

R.V. OPTIMIZATION (FITTING)

HOW DO I QUANTIFY UNC.?

$$P(\theta < 0.02 | y) = ?$$

THIS QUESTION IS INAPPROPRIATE

$$\hat{\theta} - \theta \xrightarrow{D} N(0, \sigma_{\theta}^2)$$

THE CONF. INTERVAL

FOR EXAMPLE,

$$CI_{95\%} \approx \left[\underset{\text{R.V.}}{\hat{\theta}} - 2 \cdot \sigma_{\theta}, \underset{\text{R.V.}}{\hat{\theta}} + 2 \sigma_{\theta} \right]$$

(95% OF SUCH INTERVALS WILL
CONTAIN TRUE VALUE θ)

"95% OF THE CI WILL CONTAIN θ " (X)

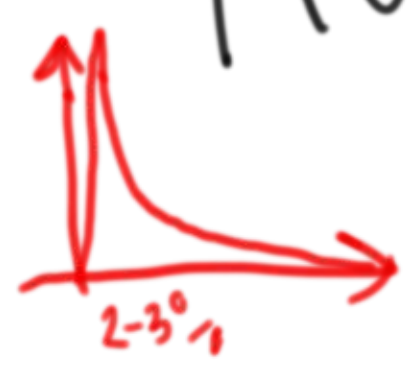
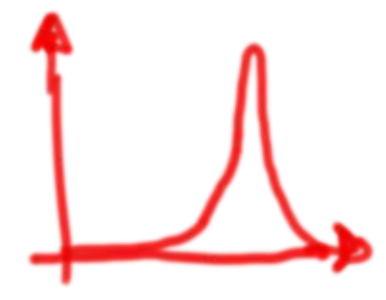
" θ IS IN CI WITH 95% PROB." (✓)

95% OF STUDIES (AT RANDOM) $CI_{95\%} = [0, 1]$
5% " " " " $CI_{95\%} = [-5, -4]$

95% CI

BAYESIAN VIEW

- y_i ARE SAMPLES FROM OUR DGP
- θ IS A CONSTANT (UNKNOWN)
 ↑
 R.V. ~~BE~~ TO REPRESENT MY UNCERTAINTY
- I'M GOING TO UPDATE MY OPINION COND. ON DATA



$$p(\theta|y) = \underbrace{p(y|\theta)}_{\text{LIKELIHOOD}} \underbrace{p(\theta)}_{\text{PRIOR}} / \underbrace{p(y)}_{\text{NORMALIZATION}}$$

POSTERIOR

OUR EXAMPLE

$$y_i | \theta \sim_{i.i.d} \text{Bernoulli}(\theta)$$

$$\theta \sim \text{Beta}(\alpha_0, \beta_0)$$

⊖ COMPUTATION

$$\theta | y \sim \text{Beta}(\underbrace{\alpha_0 + \sum y_i}_{\# \text{ POS.}}, \underbrace{\beta_0 + m - \sum y_i}_{\# \text{ NEG.}})$$

$$\int p(\theta) p(\theta) d\theta$$

ALWAYS DERIVE POST?
NO (RARELY)

A "NICE"



$$P(\theta > 2\% | y) = \int_{2\%}^1 p(\theta | y) d\theta$$

$$P(\theta > 2\%) = \int_0^{2\%} p(\theta) d\theta$$



⊕ EASIER TO INTERPRET

⊕ EASY TO INCL. PRIOR KNOWLEDGE

⊕ EASIER DECISION THEORY

⊖ COMPUTATION (= TIME & EFFORT)

LINEAR REGRESSION

$$y_i = \beta^T x_i + \varepsilon_i, \quad \varepsilon_i \sim N(0, \sigma^2)$$

$$y_i | \beta, \sigma^2, x_i \sim N(\beta^T x_i, \sigma^2)$$

$$P(y | \dots) = \left(\frac{1}{2\pi\sigma^2} \right)^{\frac{n}{2}} \exp \left\{ -\frac{1}{2\sigma^2} (y - X\beta)^T (y - X\beta) \right\}$$

SUM OF SQ. RESIDUALS

PRIOR:

$$P(\beta, \sigma^2) = P(\beta | \sigma^2) P(\sigma^2)$$

$$\beta \in \mathbb{R}^k$$
$$\sigma^2 \in \mathbb{R}$$

~~NS~~

$$\underbrace{N_k(\mu_\beta, \sigma^2 V_\beta) \times IG(\alpha, \beta)}_{\text{NIG}}$$

$$P(\beta, \sigma^2 | \dots)$$

REGULARIZATION

$$\beta^* = (\lambda I + X^T X)^{-1} X^T y$$

↑
RIDGE REGRESSION

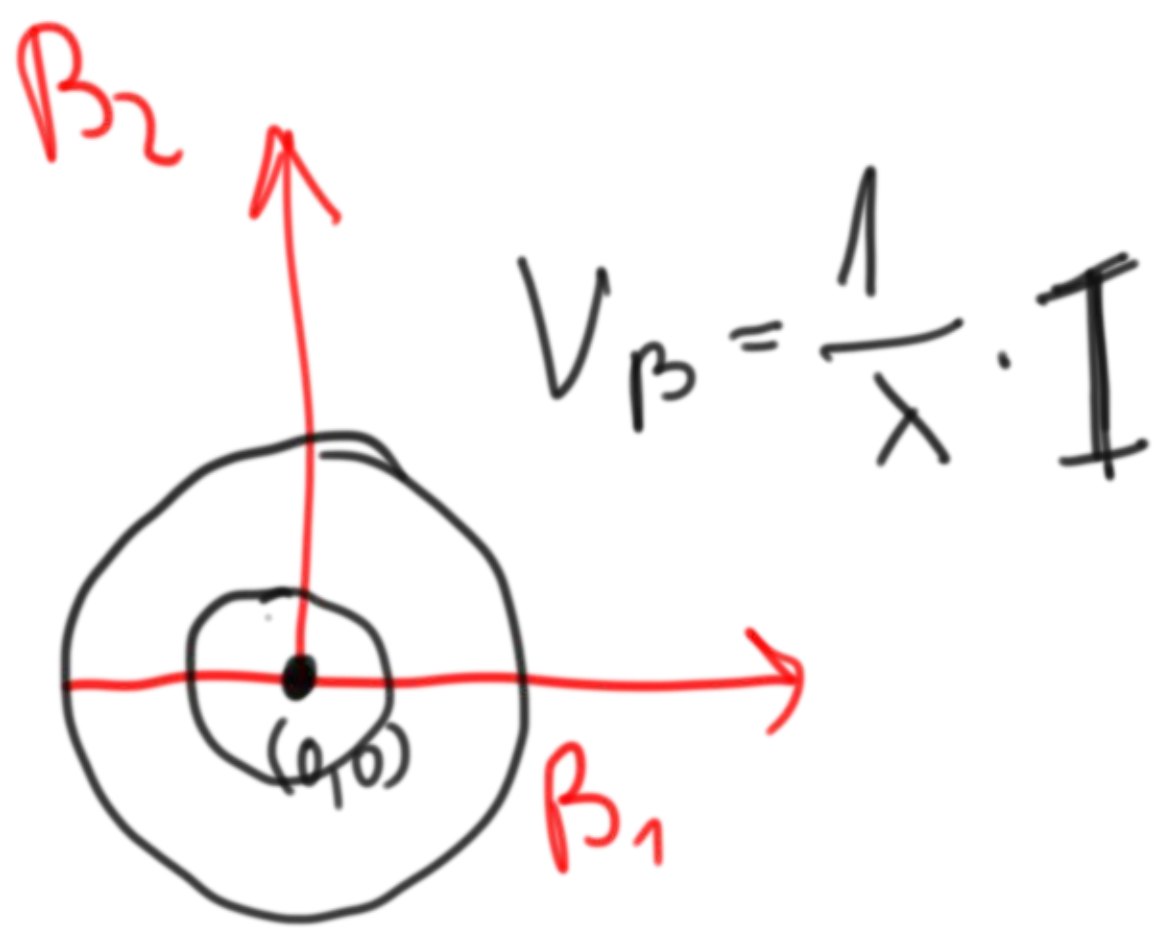
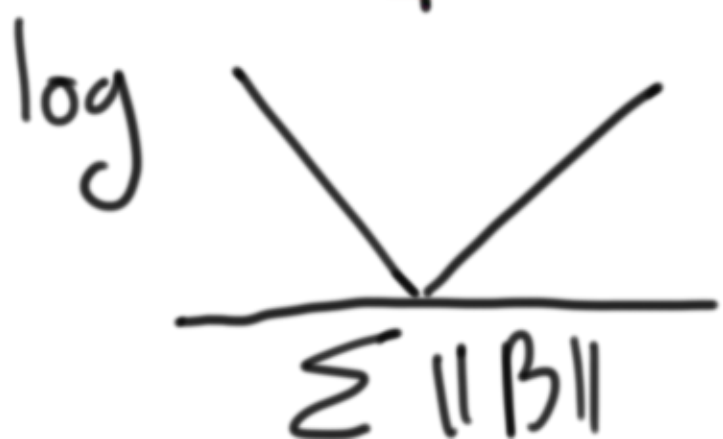
$$\beta_{OLS} = (X^T X)^{-1} X^T y$$

$$\beta | \sigma^2 \sim N_k(\mu_\beta, V_\beta)$$

LASSO (L1)

PRIOR ON β

$$\beta \sim \text{Laplace}(\cdot | \frac{1}{\lambda})$$



⊕ EASY TO INCLUDE PRIOR KNOWLEDGE
(REG.) IS AN EXAMPLE)

⊕ NO FITTING \Rightarrow NO OVERFITTING

⊖ COMPUTATION

COMPUTATION FOR BAYESIAN INF.

① 'NICE' POSTERIOR (CONJUGATE PRIOR)

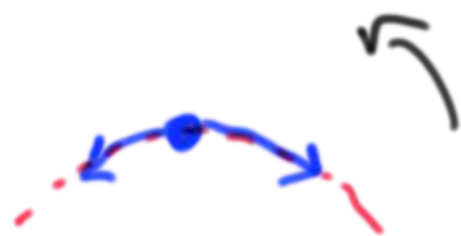
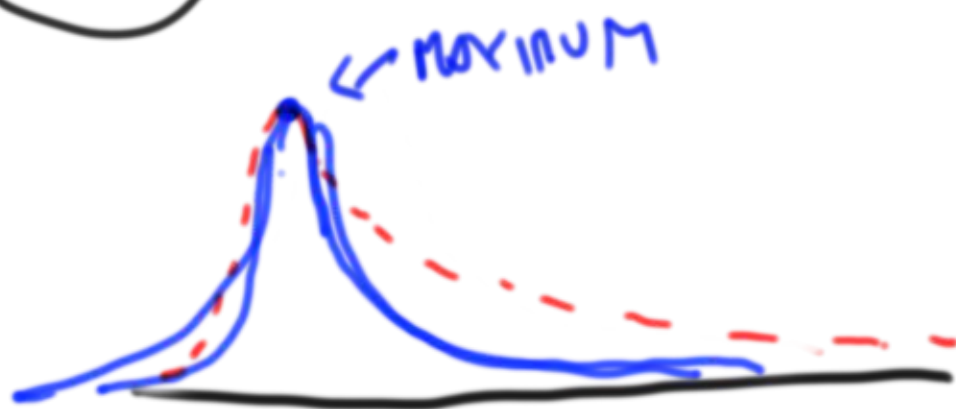
(RARE)

$$\int \theta p(\theta|y) d\theta$$

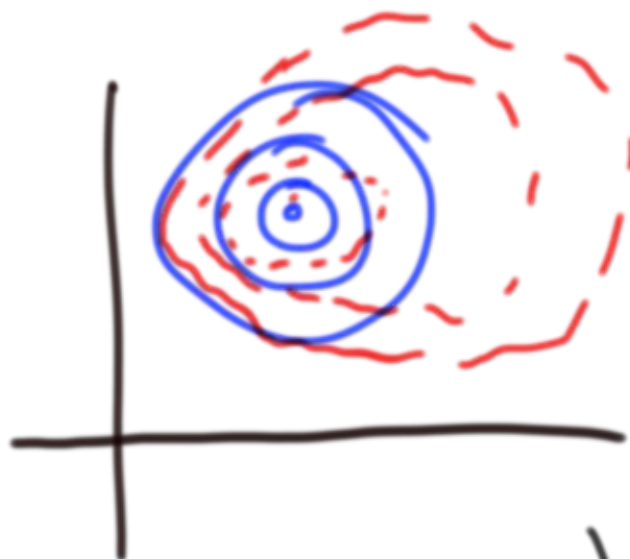
(NUMERICAL) $P(\theta > 2\%) = \int p(\theta|y) d\theta$

② MCMC (UNBIASED, COMP. INTENSIVE)

③ STRUCTURAL APPROX. (BIAS)



LAPLACE APPROX



$$\theta_{\text{MAP}} = \underset{\theta}{\text{argmax}} p(\theta|y)$$
$$p(y|\theta)$$

VARIENTIONAL INFERENCE