

Bayesian Inference

It's not about true param., it's about true parameter Distribution!

Bayesian View: Data are fixed in a certain way. We know 100% for sure that the data is there, it's all the information we have right now!

The uncertainty lies in the parameter which produced the data. Parameter viewed as a (vector-valued) random variable.

Underlying subjective interpretation of probability which can be more intuitive than frequentist (repetitions of experiment impossible in reality)

→ We probabilistically express our knowledge/belief about parameter(s) Θ given limited amount of obs. data (sample) data are there, don't change & we don't get more (or One case is useless Clinical Test where $P(B|A) = P(B|\bar{A}) = 0.5$)

Prob. of B given A , prior probability of A Prior $P(A) = 0/1 \Rightarrow$ Posterior $P(B|A)$ also $0/1$

Inverse Probability, Bayes Theorem] $P(A|B) = \frac{P(B|A) \cdot P(A)}{P(B|A) \cdot P(A) + P(B|\bar{A}) \cdot P(\bar{A})} = \frac{P(B|A) \cdot P(A)}{P(B)}$ (more general) $P(A_i|B) = \frac{P(B|A_i) \cdot P(A_i)}{\sum_{j=1}^n P(B|A_j) \cdot P(A_j)}$

Bayes Theorem $\hat{=}$ Update Rule for subjective probability This here describes the discrete case in simplified form for two events, B as new evidence

data $x = (x_1, \dots, x_n)$ (continuous parameter: π) $p(\pi|x) = \frac{f(x|\pi) \cdot p(\pi)}{\int f(x|\pi) \cdot p(\pi) d\pi}$ Likelihood of data x prior distr. density $f(x|\pi) \hat{=}$ Likelihood if we "fixate" the observed data x to interval $[0,1]$ Likelihood $\hat{=}$ probability cause it has not been weighted yet by prior and not normalized (in case of discrete π) Posterior up to multiplicative constant $p(\pi|x) \propto f(x|\pi) \cdot p(\pi) \propto f(x|\pi)$ for exact prior

• We could also use notation $f(\pi|x)$, $f(\pi)$ instead. Notation above highlight prior & posterior!

• We should actually write $f(x|\pi)$ for discrete RV π below, in general $f(x|\pi)$ (we can define integral as $\int f(x|\pi) p(\pi) d\pi$ Integration over π inherently Bayesian)

(Conditional) Posterior Distr. $\pi|x=x$ \hookrightarrow Then: $f(x) = \sum_{\pi} f(x|\pi) P(\pi=\pi)$ continuous π : $f(x) = \int f(x|\pi) p(\pi) d\pi$ We "integrate out parameter"

$f(x)$ is marginal prob. of observing the data, across whole parameter space

- Form posterior by combining prior knowledge/assumptions/beliefs about parameter with new observed data \Rightarrow using Bayes reduces our uncertainty about Θ

[Random distribution of estimates for Θ ; posterior denotes our new current belief about prob. distr. of param.]

Subjectivity is implied. There is never complete certainty / perfect knowledge about anything or else there would be no need for statistics in the first place!!

- Bayesian Learning can be sequential. Posterior becomes new prior \sim new data \dots new posterior

Bayes-Inference formal 1. Set probabilistic model for (i.i.d.) data y $F = \{f(\cdot, \theta) | \theta \in \Theta\}$ We only "know" this conditional data density, $f(y)$?? From there we derive Likelihood $f(y|\theta)$

2. For Θ impose a prior distribution that captures our knowledge/beliefs about param. $\theta \sim p(\cdot, w)$ often just $\theta \sim p(\theta)$ for simplicity which may itself depend on additional hyperparameters w $L(\theta) = \exp(\eta(\theta))$ \Rightarrow then posterior $p(\theta|y)$

3. Calculate Θ 's posterior distribution $p_{\theta|y}(\theta|y) = \frac{\prod_{i=1}^n f(y_i, \theta) p(\theta, w)}{\int \prod_{i=1}^n f(y_i, \pi) p(\pi, w) d\pi}$ with $f(y, w) = \int \prod_{i=1}^n f(y_i, \theta) p(\theta, w) d\theta$

Conjugate Prior Distribution Leads to posterior from same distr. family as prior i.e. $p_{\theta|y} \in \mathcal{P} = \{p(\cdot, \theta) | w \in \mathcal{W}\}$ $\mathcal{P} \hat{=}$ prior structure model class for given probability model F . Not all prob. models in F have such conjugate prior!

List of Conjugate Priors and respective Posteriors for Common Exponential Family Distributions

	Prob. Model / Likelihood	Parameter	Prior Conjugate	Posterior with updated hyperparameters	Jeffreys Prior
discrete f	Bern(π), Bin(n, π)	$\pi \in (0,1)$	Beta(α, β)	Beta($\alpha+y, \beta+n-y$)	$p^*(\pi) \propto \pi^{-1/2} (1-\pi)^{-1/2} \hat{=}$ Beta($\frac{1}{2}, \frac{1}{2}$) core
	Geom(π)	$\pi \in (0,1)$	Beta(α, β)	Beta($\alpha+n, \beta+\sum y_i$)	
	Neg. Bin(r, π)	$r > 0, \pi \in (0,1)$	Beta(α, β)	Beta($\alpha+r, \beta+\sum y_i$)	
	Poisson(λ)	$\lambda > 0$	Gamma(α, β)	Gamma($\alpha+\sum y_i, \beta+n$)	$p^*(\lambda) \propto \lambda^{-1/2} \hat{=}$ (improper conjugate prior) special case of Gal(α, β) $\alpha \rightarrow 1/2, \beta \rightarrow 0$
continuous f	Exp(λ)	$\lambda > 0$	Gamma(α, β)	Gamma($\alpha+n, \beta+\sum y_i$) $\sum y_i = n\bar{y}$	$p^*(\lambda) \propto \lambda^{-1} \hat{=}$ special case of Gal(α, β) $\alpha \rightarrow 0, \beta \rightarrow 0$
	Gamma(known α, β)	$\beta > 0$	Gamma(α, β)	Gamma($\alpha+n, \beta+\sum y_i$)	
	$N(\mu, \text{known } \sigma^2)$	$\mu \in \mathbb{R}$	$N(\mu_0, \tau^2)$	$N(\frac{\frac{M_0}{\tau^2} + \frac{n\bar{y}}{\sigma^2}}{\frac{1}{\tau^2} + \frac{n}{\sigma^2}}, \frac{1}{\frac{1}{\tau^2} + \frac{n}{\sigma^2}})$	$p^*(\mu) \propto \text{const.}$
	$N(\text{known } \mu, \sigma^2)$	$\sigma^2 > 0$	Inv-Gamma(α, β)	Inv-Gamma($\alpha+n/2, \beta+\frac{1}{2} \sum (y_i - \mu)^2$)	$p^*(\sigma^2) \propto 1/\sigma^2$
	$N(\mu, \sigma^2)$	(μ, σ^2)	Normal Inv-Gamma	Normal Inv-Gamma updated	

(connection betw. Bayes and MLE: For fixated i.i.d. data get posterior $p(\theta|y) = \frac{L(\theta) p(\theta)}{\int L(\pi) p(\pi) d\pi}$, Let's define posterior mode $\hat{\theta}_{\text{post, mode}} = \arg \max_{\theta \in \Theta} p(\theta|y)$ For constant FLAT prior $p(\theta) \Rightarrow p(\theta|y) \propto L(\theta)$ and posterior mode equals Max. Likelihood estimate $\arg \max_{\theta \in \Theta} L(\theta)$ [Laplace Prior makes Bayesian frequentistic]

Proof uses penalized Likelihood $f(\theta) + \log p(\theta)$

Bernstein-von-Mises Theorem: For massive sample size n and appropriate prior it holds $\theta \sim N(\hat{\theta}, \frac{1}{n} I(\hat{\theta})) \Rightarrow$ Bayesian Inference corresponds to ML-inf.

This however says nothing about finite samples in reality! Big upside of Bayesian is that it still works when CLT and thus MLE fails, e.g. extreme data or just 1 observ. Posterior always calculable

everything proper to this (e.g. scales) $\hat{=}$ Posterior Core $\rightarrow p(\theta|y)$ higher and more dense \wedge as $n \uparrow \Rightarrow$ Prior has less influence, "weak" uniform prior requires larger n to converge precise poster.

Posterior $p(\theta|y) = \frac{f(x|\theta) p(\theta)}{\int f(x|\pi) p(\pi) d\pi}$ contains ALL information about θ after observing data

Bayesian Denominator / Normalization Constant! Posterior dens. concentrate around more likely θ values

Posterior shape reflects shape of prior & Likelihood. More informative (less flat) prior \rightarrow prior influence \uparrow

Bayes-Principle: All conclusions drawn from posterior distr.

- for multidim. Θ we look at marginal distr. $\theta_j | x$
- point estimators (Posterior-EV / Median, Maximum MAP etc.)
- Interval estimators (Credibility Interval) and Tests
- Model comparison
- Prediction

How to choose Prior

On the one hand we want to consider prior knowledge by using an (informative) prior, learn sequentially etc.

then data evidence has only tiny influence

On the other hand the prior choice can heavily skew our results. Too informative prior basically determines posterior already!

Obviously it matters how good our prior knowledge is, how confident are in our assumptions (keep in mind that there is always subjectivity involved!)

Is not a good idea to artificially make the prior less informative than an honest assessment calls for!

Subjective Prior sometimes useful. Ideally based on reliable expert knowledge. More on Page Two...

Poisson $\lambda \sim \text{Ga}(a, b)$: Choose a, b value s.t. Prior $E(\lambda) = \frac{a}{b}$ matches beliefs

vs. for achieve both

Within a specific family of prior distr. we can "measure" informativeness through the variance of param. $\text{Var}(\theta)$ given $\theta \sim p(\cdot, w)$

controlled by priors hyperparameters

• Larger Variance \iff less condensed prior density, meaning less prior information

• Tiny Variance \iff highly concentrated prior, very informative

Example: Often choose $\text{Ga}(a=0.001, b=0.001)$ as prior for Poisson λ

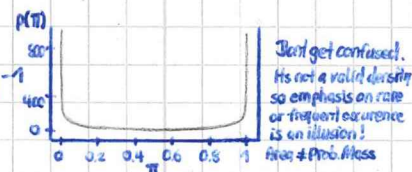
That leads to high $\text{Var}(\lambda) = \frac{a}{b^2} = 1000$ for Prior λ , low informativeness

But what do we do if there is no Prior Information at all OR (as often an issue in reality) its not really quantifiable? Then we aim for Uninformative Prior

Going back to Billard Ball / Lottery Example: Prior $\text{Beta}(a=1, b=1)$ where we interpreted a as "prior Nr. of successes" b as "prior Nr. of failures"

But $a=1, b=1$ seems to have given some info. Intuitively no information implies $a=0, b=0$

\implies inputting $a=0, b=0$ in core of Beta distribution produces the Haldane Prior $p(\pi) \propto \pi^{-1}(1-\pi)^{-1}$



Careful! Haldane is an improper prior $\int_0^1 \pi^{-1}(1-\pi)^{-1} d\pi \neq 1$ divergent Normalising $\frac{1}{B(a,b)} = \frac{\Gamma(a)\Gamma(b)}{\Gamma(a+b)} = 1$ not working here

$$\int_0^1 \left(\frac{1}{\pi} - \frac{1}{1-\pi}\right) d\pi = \lim_{\epsilon \rightarrow 0} [\log \pi + \log(1-\pi)]_0^\epsilon = \infty$$

Its a borderline case derived from $\text{Beta}(a,b)$ with both $a, b \rightarrow 0$

In general definition of degenerate/improper prior $\int_{\Theta} p(\theta; w) \neq 1$ divergent Many flat priors are improper

if $x > 0, n-x > 0 \triangleq \text{Beta}(x, n-x)$

Using Haldane prior $\pi \sim \text{Beta}(0,0)$ which is not a real (Beta-) Distr. still functions a conjugate Prior $p(\pi|x) \propto \pi^{-1}(1-\pi)^{-1} \cdot \pi^x(1-\pi)^{n-x} = \pi^{x-1}(1-\pi)^{n-x-1}$

Even with improper Prior, the posterior is usually well defined ✓ (Convergence to proper Posterior)

For $n \rightarrow \infty$ and finite dimensional param. space always the case that Posterior integrates to one

Normalizing by data: Sample n large enough, Likelihood (usually sharp locally concentrated) "strong enough" to make $\int f(x|\theta) p(\theta) d\theta < \infty$ finite Bayesian Denominator

Analogy: Flashlight illuminates finite spot in infinite fog, data cut out finite volume from ∞ area

Interesting: Non-informative Prior might be edge case of conjugate Prior

Laplace Prior 1 Laplace-Principle: Without knowledge should initially assume equal probability for all elementary events $P(w_i) = 1/L$ discrete uniform distr.

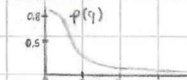
In general defined as $p(\theta) \propto \text{constant}$ \implies flat prior, but not necessarily totally uninformative as explained above with $\text{Beta}(1,1) \iff U(0,1)$ prior

continuous uniform distr. for constrained θ smt. refer to $p(\theta) \propto \text{const.}$ as "uniform distr. on \mathbb{R}^d "

If parameter space unconstrained then Laplace Prior is improper

Problem: Laplace-Prior on transformed Parameter can lead to different Prior & Post. Family (unlike Haldane-Prior which is trans. invariant)

Example: Instead of π use odds $q = \frac{\pi}{1-\pi}$ density transform with $\pi \sim U(0,1)$ $\pi(q) = \frac{q}{q+1}$ $\frac{\partial \pi(q)}{\partial q} = \frac{1}{(q+1)^2}$ $f_q(q) = f_\pi(\pi(q)) \cdot \left| \frac{\partial \pi(q)}{\partial q} \right| = \frac{1}{(1+q)^2} \cdot 1$



Transformed Prior not uniform, not constant flat. Thus Posterior looks different, no longer just influenced by Likelihood as for original param

This also showcases that this Laplace-Prior does contain a tiny bit of information, transformed Prior is pretty informative

Jeffreys Prior non-informative Prior, invariant to Parameter - Transf. (bijective Φ) $p^*(\Phi(\theta))$ is again Jeffreys Prior for $\Phi(\theta)$ and therefore uninformative

$p^*(\theta) \propto \pm^{1/2}(\theta)$ Jeffreys prior defined as sq. root of Fisher Information Matrix Jeffreys Prior can, but doesn't have to be improper

Under certain conditions: $\pm(\theta) = \text{Var}(s(x;\theta)) = E[(s(x;\theta))^2]$ as $E(s(x;\theta)) = 0$ score-fct. $\frac{\partial}{\partial \theta} \log f(x|\theta)$

$\pm(\theta) = -E\left[\frac{\partial^2 \log f(x|\theta)}{\partial \theta \partial \theta^T}\right] = E\left[\left(\frac{\partial}{\partial \theta} \log f(x|\theta)\right)^2\right]$ Interpretable as expected curvature of log-Likelihood

Frequentist: integration over data x !

Fisher Info indicates degree of param. information in a random sample of data, param. viewed as fixed

Jeffreys Prior for Binomial Distribution $\log(f(x|\pi)) = x \log \pi + (1-x) \log(1-\pi)$, $S_x(\pi) = \frac{x}{\pi} - \frac{1-x}{1-\pi}$, $\pm(\pi) = E[S_x(\pi)] = \frac{1}{\pi} - \frac{1-\pi}{1-\pi} = \frac{1}{\pi(1-\pi)}$ $p^*(\pi) \propto \pm(\pi) = \pi^{-1/2}(1-\pi)^{-1/2} = \frac{1}{\pi(1-\pi)} \triangleq \text{Beta}(\frac{1}{2}, \frac{1}{2})$ density core

Jeffreys Prior for Exp. Distribution $\log(f(x|\lambda)) = \log \lambda - \lambda x$, $S_x(\lambda) = \frac{1}{\lambda} - x$ (improper, else case conjugate) $\pm(\lambda) = \text{Var}(S_x) = \text{Var}(\frac{1}{\lambda} - x) = \text{Var}(x) = 1/\lambda^2 \implies p^*(\lambda) \propto \lambda^{-1/2}$ special case of $\text{Ga}(a,b)$ Poisson $p^*(\lambda) \propto \lambda^{-1/2}$ special case of $\text{Ga}(a,b)$

Expect $\pm(\theta)$ to have much shape max at thing! No π shape similar Haldane but its proper Prior!

Guidelines/Approaches 1. Prior-Hyperparam. choice should closely reflect Prior Knowledge about θ (subjective! maybe hard...)

2. Prior should be as uninformative / flat as possible (esp. if no knowledge) and/or be independent of parametrization

Sometimes we prefer rel. "flat" but proper Priors over improper Laplace or improper Jeffreys (choose high var Prior)

Inherently NOT Bayesian Empirical Bayes

3. Estimate Prior-Hyperparams from Data (related w still fixed)

$\tilde{f}(w) = \int_{\Theta} f(w, \theta) p(\theta|w) d\theta$ marginal distr. distr. given w

new Likelihood $\tilde{f}(w) = \int_{\Theta} f(w, \theta) p(\theta|w) d\theta$

Subjective Priors

One issue: Expert Knowledge Strong beliefs typically not given with a specific param. density



Solution: • For discrete Params should be pretty easy to formulate Probabilities that add up to one

- For continuous Params we could discretize it into intervals, find parametric density that "matches" knowledge or use sth. like Maximum Entropy Prior

Parametric subjective Prior i) Choose matching param distribution (symmetric gaussian, skewed gamma, beta etc.)

In the easiest case fits the conjugate Prior

ii) Set Prior Hyperparameters so that Prior Knowledge is captured with unique theoretical Hyperparam. combination

(mean, variance, or more robust statistics like median, quantiles & perhaps credibility interval)

If there are several Hyperparams we need more than just a single information to get unique Hyperparams

Lack of additional information? Make reasonable assumption or use information criterion under constraint of given knowledge

Example: • Prior Knowledge on Bernoulli π is mean = 0.3 and 95% credibility interval $[0.1, 0.5]$ i.e. $P(0.1 < \pi < 0.5) = 0.95$

Conjugate Prior $\pi \sim \text{Beta}(a, b)$ with EV $\frac{a}{a+b} \Rightarrow$ map it to mean: $\frac{a}{a+b} = 0.3 \Leftrightarrow b = \frac{7}{3}a$

Additional info about Cred. Interv.: Numerically solve Min. deviation of $(p_{\text{beta}}(0.5, a, \frac{7}{3}a) - p_{\text{beta}}(0.1, a, \frac{7}{3}a))$ from 0.95 regarding e.g. L2 Loss

\Rightarrow In R we can use `optimize(f, interval)` to search interval $c(0, 1000)$ which yields $a \approx 5.45 \leadsto b \approx 12.72$

*Same as above but only info about median = 0.36. This gets much more complex, also because Median(Beta) has no general form and needs to be approx.

We could numerically solve $\text{Med}(\text{Beta}(a, b)) = 0.36$ under $a+b = d$ with d measuring degree of confidence in knowledge (low $d=2$, medium $d=10$, high $d=50$)

(Another option is a very conservative Minimal Information prior that maximizes Entropy under Median-constraint? Applying MEP requires using robust Median as $\hat{E}(\theta)$)

Maximum Entropy Prior Basic idea is to choose minimally informative prior distr. that still uses Prior Information

Definition Shannon-Entropy: $H(\theta) = E[-\log f(\theta)] = - \int_{\theta} \log(p(\theta)) p(\theta) d\theta = - \sum_{\theta} \log(p(\theta)) p(\theta)$
regarding prior
for discrete θ

Prior Knowledge appears in form of certain of certain moments $E(g_k(\theta))$ of random variable, here parameter θ

↓
There exists a Max. Entropy Prior (MEP) $p_{\text{MEP}}(\theta) = \frac{1}{Z} \cdot \exp \left(\sum_{k=1}^K c_k g_k(\theta) \right)$ which preserves the moments
normalisation const. coeff. θ^k $E(g_1(\theta)) = m_1(\theta) = E(\theta)$ etc.