# Non-Informative Priors:

A non-informative Prior (an objective prior) is a prior which is somehow automatic, reflecting the lack of any initial knowledge about the parameter. It may have no probabilistic interpretation and so does not have to be a valid probability distribution. Non-informative priors can be used when little no reliable information is available.

## (1) Uniform Priors:

Definition: The uniform prior or flat prior is the prior $P(\vartheta) \propto 1$.

This is the obvious choice of lack of information; every value being equally likely. Under this prior,

$$P(\vartheta|x) = \frac{L(\vartheta, x)}{\int_{\Theta} L(\vartheta, x)\, d\vartheta}$$

which is well defined as long as

$$\int_{\Theta} L(\vartheta, x)\, d\vartheta < \infty .$$

Example: Let $X \sim Exp(\vartheta)$ and $P(\vartheta)=1$
The marginal likelihood is $\int_0^{\infty} e^{-\vartheta x}\, d\vartheta$
which is finite for all $x > 0$, so
the posterior is well-defined.

Properties: Let $n = \log \vartheta$. Then the

prior for $n$ is $\qquad \rightsquigarrow \vartheta = e^n$

$$\tilde{P}(n) = P(\vartheta(n)) \frac{d\vartheta}{dn} = \frac{d\vartheta}{dn} = 1 \cdot e^n \neq 1$$

After reparametrisation the prior is not

flat anymore. In fact, as a prior in $\eta$, $\tilde{P}$ is very informative (large values are more likely than small ones)

## (2) Jeffrey's Prior:

(W)e need a prior which does not depend on the parametrisation.

**Definition:** In one-dimensional case Jeffrey's prior is given by:

$$P(\theta) \propto \left( I(\theta) \right)^{1/2}$$

where $I(\theta) = E_\theta \left( -\frac{\partial^2}{\partial \theta^2} \ell(\theta, x) \right)$ is the Fisher Inf.

**Remark:** If $\theta = g(\psi)$ for some one-to-one differentiable function $g$ then the reparametrised prior is

$$\tilde{P}(\psi) \propto P(g(\psi)) \cdot |g'(\psi)|$$

Recall, that $I(\psi) = (g'(\psi))^2 I_\theta$

So $\sqrt{I_\psi} = \sqrt{I_\theta} |g'(\psi)|$.

Hence, $\tilde{P}(\psi) \propto \sqrt{I_\psi}$

⚠ Jeffrey's prior is invariant under reparametrisation.

<u>Definition</u>: The k-dimensional Jeffrey's prior is given by

$$P(\theta) \propto |I_\theta|^{1/2}$$

where $|I_\theta| = \det I_\theta$, $I_\theta$ is the Fisher Information matrix, so under the standard regularity assumptions

$$(I_\theta)_{ij} = -\mathbb{E}_\theta \left( \frac{\partial^2}{\partial\theta_i \partial\theta_j} \ell(\theta, x) \right)$$

**Example:** Suppose $X \sim \text{Poi}(\theta)$ so that $f(x, \theta) = \dfrac{e^{-\theta} \cdot \theta^x}{x!}$ for $x = 0, 1, 2, \dots$

The Jeffrey's prior is 
$$\left( \begin{array}{l} \log f(x; \theta) = \\ = x \log \theta - \theta - \log x! \end{array} \right)$$

$$P(\theta) \propto \left( I_X(\theta) \right)^{1/2}$$

$$= \left( E \left( \frac{X}{\theta} - 1 \right)^2 \right)^{1/2}$$

$$= \left( \sum_{x=0}^{\infty} f(x, \theta) \left( \frac{x - \theta}{\theta} \right)^2 \right)^{1/2}$$

$$= \left( e^{-\theta} \sum_{x=0}^{\infty} \frac{\theta^x}{x!} \left( \frac{x^2}{\theta^2} - \frac{2x}{\theta} + 1 \right) \right)$$

$$= \left( \frac{1}{\theta^2} E \left( E (X - \theta)^2 \right) \right)^{1/2}$$

$$= \theta^{-1/2} \, .$$

**Note:** This is an improper prior.

# (3) Maximum Entropy Prior:

This prior is inspired by information theory.

## Definition: The entropy of a pdf P is defined as

$$H(P) = -\int_{\Theta} P(\vartheta) \log P(\vartheta) \, d\vartheta.$$

A maximum entropy probability distribution has entropy that is at least as great as that of all other members of a specified class of probability distributions. This ensures the least biased or the most non-informative choice, assuming no other knowledge is available. This approach minimizes prior assumption or information

**Example:** Suppose we know the mean and variance of $\vartheta$. Then the maximum entropy distribution is the Gaussian distribution.

$$P(\vartheta) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(\vartheta-\mu)^2/(2\sigma^2)}$$

**Note:** This is because Gaussian distribution has the highest entropy among all distributions on $\mathbb{R}$ with a given mean $\mu$ and variance $\sigma^2$.

**Constraints**
- $\int P(\vartheta) d\vartheta = 1$

- $\int \vartheta P(\vartheta) d\vartheta = \mu$

- $\int (\vartheta-\mu)^2 P(\vartheta) d\vartheta = \sigma^2$

We want to maximize entropy:

$$H(P) = -\int_{-\infty}^{+\infty} P(\vartheta) \log P(\vartheta) d\vartheta$$

Subject to constraints.

Use lagrangian multipliers

$$\mathcal{L}(P) = - \int P(\vartheta) \log P(\vartheta) \, d\vartheta + \lambda_1 \int P(\vartheta) \, d\vartheta - 1$$

$$+ \lambda_2 \left( \int \vartheta \, P(\vartheta) \, d\vartheta - \mu \right)$$

$$+ \lambda_3 \left( \int (\vartheta - \mu)^2 \, P(\vartheta) \, d\vartheta - \sigma^2 \right)$$

$$\Longrightarrow \quad \frac{\partial \mathcal{L}}{\partial P(\vartheta)} = -\log P(\vartheta) - 1 + \lambda_1 + \lambda_2 \vartheta$$

$$+ \lambda_3 (\vartheta - \mu)^2 = 0$$

Solve for $P(\vartheta)$

$$\Longrightarrow P(\vartheta) = \exp \left( \lambda_2 \vartheta + \lambda_3 (\vartheta - \mu)^2 + \lambda_1 - 1 \right)$$

$$\Longrightarrow P(\vartheta) = C \exp \left( \lambda_2 \vartheta + \lambda_3 (\vartheta - t)^2 \right)$$

if $\lambda_3 < 0$ then the exponent is

concave $P(\vartheta)$ integrates to 1

choosing $\lambda_2 = 0$ $\quad \lambda_3 = -\frac{1}{2\sigma^2}$

$$P(\vartheta) = C \cdot \exp\left(-\frac{(\vartheta - \mu)^2}{2\sigma^2}\right)$$

we normalize

$$\int_{-\infty}^{+\infty} C \cdot \exp\left(-\frac{(\vartheta - \mu)^2}{2\sigma^2}\right) d\vartheta = 1$$

we solve for $C$

this gives $C = \frac{1}{\sqrt{2\pi\sigma^2}}$.

Theorem: Let

$$P(\vartheta) = \exp\left(\sum_{i=1}^{K} \lambda_i T_i(\vartheta) - B(\lambda)\right), \forall \vartheta \in \Theta.$$

be a probability density function

and suppose that

$$\int T_i(x) p(\theta) dx = t_i \qquad \text{for } i = 1, \ldots, k \quad (*)$$

Then $p$ uniquely maximizes $H(p)$ among all densities satisfying the constraint.

Proof: Let $\Pi$ be the class of distributions satisfying $(*)$.

Recall that for 2 distributions $\Pi_1 << \Pi_2$ the Kullback-Leibler divergence $KL(\Pi_1 \| \Pi_2)$ is defined through:

$$KL(\Pi_1 \| \Pi_2) = \int \Pi_1(dx) \log\left(\frac{d\Pi_1}{d\Pi_2}\right)$$

where $\frac{d\pi_1}{d\pi_2}$ is the radon Nikodym
derivative. If $\pi_1$ is not absolute
continuous w.r.t $\pi_2$ we set

$KL(\pi_1 \| \pi_2) = + \infty$. It is a simple

application of Jensen's inequality
to check that $KL(\pi_1 \| \pi_2) \geq 0$

Let $p'$ be an element on $\pi$

$$H(p') = -\int p'(x) \log p'(x) \, dx$$

$$= -\int p'(x) \log \left( \frac{p'(x)}{p(x)} \right) \, dx$$

$$- \int p'(x) \log p(x) \, dx$$

$$= -KL(p' \| p) - \int p'(x) \log p(x) \, dx$$

$$= -KL\left(P' \| P\right) - \int P'(x) \left(\sum_{i=1}^{P} \lambda_i T_i(x)\right) dx + B(\lambda)$$

and since $P' \in \Pi$, it satisfies
the same moment constraints as $P$

So,
$$\int P' \log P = \sum_{i=1}^{K} \lambda_i t_i - B(\lambda)$$

$$= \int \Pi \log \Pi$$

There fore
$$H(P') = -KL(P' \| P) + H(P)$$

$$\Rightarrow \quad H(P) - H(P') = KL(P' \| P) \geq 0$$

Equality holds only when
$P' = P$, so $P$ is uniquely
maximizes the entropy

under the given constraint (*) □

**Note:** Theorem confirms that the maximum entropy distribution under moment constraint is the exponential family distribution.

**Example:** In the Gaussian example:

$$\mathbb{E}\left(T_1(\vartheta)\right) = \mu \qquad \mathbb{E}\left(T_2(\vartheta)\right) = \sigma^2$$

where $T_1(\vartheta) = \vartheta \qquad T_2(\vartheta) = (\vartheta - \mu)^2$

By the previous theorem the maximum entropy prior is of the form:

$$p(\vartheta) \propto \exp\left(\lambda_1 \vartheta + \lambda_2 (\vartheta - \mu)^2\right)$$

The 2 constraints then imply that $\lambda_1 = 0$ $\lambda_2 = -\frac{1}{2\sigma^2}$ .