

# Ignorance priors and transformation groups

← Back to Chapters

## Comments on 12.4

I find the discussion in 12.4 somewhat off the mark. “Statistical decision theory and bayesian analysis” by James O. Berger is a better reference.

To put the main issue up first: The prior 12.27 is in fact “best” in this problem, but not because 12.18 (which has multiple issues) is better than 12.30, but because the deduction of 12.36 from 12.30 is flawed, and a better analysis of 12.30 (using the framework of invariant decision rules) does indeed lead to 12.27!

I flesh some of this out in the writeup below, which is largely based Berger’s book, particularly sections 3.3.2 and 6.6.

---

The word “invariance” presupposes a group action. The most natural setting is that in which a group  $G$  acts on the space  $X$  in which we get data.

---

Example

1:  
(location-

scale

in

1D)

Take

$X =$

$\mathbb{A}^1$

the

affine

line,

and

$G$  the

group

of

(orientation-

preserving)

affine

trans-

for-

ma-

tions

of  $X$ .

As is

com-

mon

(for

ex-

am-

ple in

com-

puter

graph-

ics),

we

repre-

sent

ele-

ments

of  $X$

as

vec-

tors

$\vec{x} =$

$\begin{pmatrix} x \\ 1 \end{pmatrix}$

(the

line

$y = 1$

in-

side

$\mathbb{R}^2$ )

and

then

---

is in-  
deed  
affine-  
linear.  
The  
prod-  
uct  
in  $G$   
is  
then:

$$\begin{pmatrix} a & b \\ 0 & 1 \end{pmatrix} \begin{pmatrix} a' & b' \\ 0 & 1 \end{pmatrix} = \begin{pmatrix} aa' & ab' + b \\ 0 & 1 \end{pmatrix}$$

This  
is a  
non-  
commutative  
group.  
We  
have  
the  
fol-  
low-  
ing  
pair  
of  
group  
ho-  
mo-  
mor-  
phisms  
(known  
as a  
short  
exact  
sequence):

$$\mathbb{R} \hookrightarrow G \twoheadrightarrow \mathbb{R}_+$$

---

Here  
 $\mathbb{R}$  is  
 the  
 real  
 num-  
 bers  
 un-  
 der  
 addi-  
 tion,  
 and  
 $\hookrightarrow$   
 sends  
 a  
 num-  
 ber  
 to a  
 trans-  
 la-  
 tion  
 by  
 that  
 amount;  
 on  
 the  
 other  
 hand  
 $\mathbb{R}_+$  is  
 posi-  
 tive  
 reals  
 un-  
 der  
 mul-  
 tipli-  
 ca-  
 tion,  
 and  
 $\rightarrow$   
 sends  
 an  
 affine  
 map  
 to its  
 stretch-  
 ing  
 factor.

---

These  
 maps  
 do  
 not  
 de-  
 pend  
 on  
 any  
 par-  
 ticu-  
 lar  
 way  
 of  
 repre-  
 sent-  
 ing  
 $X$   
 and  
 $G$ .  
 On  
 the  
 other  
 hand,  
 there  
 are  
 “back-  
 ward  
 maps”  
 $\mathbb{R}_+ \rightarrow$   
 $G$   
 and  
 $G \rightarrow$   
 $\mathbb{R}$ ,  
 but  
 those  
 are  
 de-  
 pend  
 on  
 addi-  
 tional  
 choices,  
 like  
 the  
 choice  
 to  
 repre-  
 sent  
 $X$   
 and  
 $G$  in  
 terms  
 of  
 vec-  
 tors

---

Example  
 2:  
 (Location-  
 scale  
 in  
 higher  
 di-  
 men-  
 sions)  
 When  
 $X =$   
 $\mathbb{A}^n$   
 the  
 location-  
 scale  
 group  
 $G$  is  
 the  
 sub-  
 group  
 of  
 those  
 affine  
 trans-  
 forms  
 of  $X$   
 whose  
 lin-  
 ear  
 part  
 is a  
 pure  
 rescal-  
 ing  
 by a  
 (posi-  
 tive)  
 fac-  
 tor.  
 We  
 then  
 have  
 $\mathbb{R}^n \hookrightarrow$   
 $G \rightarrow$   
 $\mathbb{R}_+$ .

---

---

We are interested in distribution of the data, i.e. in probability distributions over  $X$ . Thus we consider a collection of  $\mathcal{P}$  of distributions over  $X$ .

**Defintion:** The collection  $\mathcal{P}$  is said to be invariant under the action of  $G$  if for any  $p \in \mathcal{P}$  and any  $g \in G$  the pushforward distribution  $g_*p$  is also in  $\mathcal{P}$ .

Since  $(g \cdot h)_*p = g_*(h_*p)$ , this means that  $G$  acts on  $\mathcal{P}$  as well. If  $\mathcal{P}$  is a parametric family, parametrized by space  $\Theta$  then we conclude that  $G$  acts on  $\Theta$ .

---

Example 1 continued:

- a) Let  $\mathcal{P}$  be the collection of all Gaussian distributions on  $\mathbb{A}^1$ . Note that in order to talk about the collection  $\mathcal{P}$  specifying the origin is not necessary. This collection is a invariant under the action of  $G$ . If we pick an origin, then we can use mean  $\mu$  and standard deviation  $\sigma$  as parameters for  $\mathcal{P}$ . We can also use representation of  $G$  by matrices that we have discussed above. Then  $\Theta = \mathbb{R} \times \mathbb{R}_+$  is the parameter space and  $G$  acts on by sending  $(\mu, \sigma)$  to  $(a\mu + b, a\sigma)$ .

(Together with  $x$  being sent to  $ax + b$  this appears as 12.30 in Jaynes.)

- b) Let  $\mathcal{P}$  be the collection of all Cauchy distributions on  $\mathbb{A}^1$ ; after choice of the origin this is the family with pdfs  $\frac{1}{\pi\gamma \left[1 + \left(\frac{x-x_0}{\gamma}\right)^2\right]}$ . We no longer have mean or standard deviation available as a parameters, but we do have location  $x_0$  and scale  $\gamma$ . They transform under  $G$  by the same formulas as  $(\mu, \sigma)$  did before.
- c) Let  $\mathcal{P}$  be the collection of all mixures of normal distributions on  $\mathbb{A}^1$ . After picking the origin, this is the collection of distributions which can be written as  $p = \sum_i^m w_i \mathcal{N}(\mu_i, \sigma_i^2)$  for some  $m \in \mathbb{N}$  and  $w_i > 0$  with  $\sum w_i = 1$ . This collection is invariant under  $G$ . When  $m$  is fixed the subfamily  $\mathcal{P}_m$  it is parametric, with parameters being  $3m$  dimensional vectors  $(\mu_i, \sigma_i, w_i)$ . If  $m$  is not fixed, however, the collection  $\mathcal{P}$  is not parametric in the usual sense of the word.

Example 3: Consider the family  $\mathcal{P}$  of all Gamma distributions on  $X = \mathbb{R}_+$ , with pdfs  $\Gamma_{(k,\theta)}(x) = \frac{1}{\Gamma(k)\theta} \left(\frac{x}{\theta}\right)^{k-1} \exp(-\frac{x}{\theta})$ . Here  $G = \mathbb{R}_+$  acts on  $X$  by multiplication,  $\mathcal{P}$  is invariant, and the action of  $a \in G$  sends  $(k, \theta)$  to  $(k, a\theta)$ .

---

Now  
a  
**prior**  
is a  
prob-  
abil-  
ity  
dis-  
tribu-  
tion  
 $\pi$   
over  
 $\mathcal{P}$ .  
This  
is  
easi-  
est to  
un-  
der-  
stand  
when  
the  
col-  
lec-  
tion  
 $\mathcal{P}$  is  
para-  
met-  
ric,  
so  
that  
we  
have  
 $\Theta \subset \mathbb{R}^d$ .  
When  
 $\Theta$  is  
open,  
we  
may,  
as  
usual,  
de-  
scribe  
the  
prior  
by its  
pdf,  
which  
we by  
abuse  
of  
nota-  
tion  
will



---


$$[g_*\pi](\theta) = \pi(g^{-1}(\theta))|J_{g^{-1}}(\theta)|$$


---

Example 1 continued: For the location-scale  $g = \begin{pmatrix} a & b \\ 0 & 1 \end{pmatrix}$  sends  $\theta = (m, s)$  to  $g(\theta) = (am + b, as)$ .

Then  $J_g(\theta) = \begin{pmatrix} a & b \\ 0 & a \end{pmatrix}$ , and so  $|J_g| = a^2$ ,  $|J_{g^{-1}}| = \frac{1}{a^2}$

$$[g_*\pi](\theta) = \frac{1}{a^2}\pi(g^{-1}(\theta)).$$


---

Now, one may argue that a transformation induced by  $g$  is simply a “change of coordinates” and the problems of forming a prior about  $\theta$  and about  $g(\theta)$  are equivalent, and so we should posit

$$[g_*\pi](\theta) = \pi(\theta)$$

A prior satisfying this is called **left invariant**. For such a prior we have in the differentiable case

$$\pi(\theta) = \pi(g^{-1}(\theta))|J_{g^{-1}}(\theta)|$$

In the cases where  $G$  acts transitively on  $\Theta$  this specifies  $\pi$  uniquely, up to a constant since if we set  $\pi(\theta_0) = C$  then for any  $\theta$  there exists  $g$  such that  $\theta = g(\theta_0)$  and then

$$\pi(\theta) = \pi(g^{-1}(\theta))|J_{g^{-1}}(\theta)| = C|J_{g^{-1}}(\theta)|$$

Example 1 continued: If  $\pi$  is left-invariant then taking  $\theta_0 = (m_0 = 0, s_0 = 1)$  and given  $\theta = (m, s)$  we have  $g = \begin{pmatrix} m & s \\ 0 & 1 \end{pmatrix}$  and

$$\pi(m, s) = \frac{C}{s^2}.$$

(This agrees with 12.36 in Jaynes.)

---

However, as pointed out by Berger (p.86) there is a logical flaw in requiring  $[g_*\pi](\theta) = \pi(\theta)$ . In fact, often the prior  $\pi$  in question is improper, and so is only determined up to a constant. Thus we can only require a weaker equality

$$[g_*\pi](\theta) = K(g)\pi(\theta)$$

for some  $g$ -dependent scaling function  $K(g)$ .

We will call priors satisfying this “weakly” invariant, and the ones with  $K(g) = 1$  “strictly” invariant.

In the differentiable case this leads to

$$K(g)\pi(\theta) = \pi(g^{-1}(\theta))|J_{g^{-1}}(\theta)|,$$

which can have many solutions.

---

Example 1 continued: Take  $\pi(\theta) = Cs^\alpha$ ,  $K\left(\begin{pmatrix} a & b \\ 0 & 1 \end{pmatrix}\right) = a^{-(\alpha+2)}$ . All of these solve the above equation.

---

The question of how to choose “invariant prior” (even when the transformation group is known) is hence not settled. One method is to use the framework of “invariant decision problems”, which suggests, that when the action of  $G$  on  $\Theta$  is “simply transitive”, one should use the **right invariant** prior. The story goes s follows:

Suppose that the group  $G$  acts on  $\Theta$  ins such a way that picking some starting  $\theta_0$  we have for any  $\theta \in \Theta$  a unique  $g \in G$  with  $g\theta_0 = \theta$ . Thus after making this choice of  $\theta_0$  we can identify this  $\theta$  with that unique  $g = f(\theta)$ . (Note that this  $f$  is a "map of  $G$ -sets:  $f(g\theta) = g \cdot f(\theta)$ , where  $\cdot$  is the multiplication in  $G$ .)

Now a distributions/measure over  $\Theta$  is a measure over  $G$ . “Strictly” invariant measures on  $\Theta$  correspond to what’s known as “left-invariant” measures on  $G$  (they are invariant under all left multiplication maps  $L_g : G \rightarrow G$  sending  $g' \in G$  to  $L_g(g') = gg'$ ). It is a theorem that such a measure is uniques up to scaling. There are, however, more “weakly” invariant measures (as we have seen), those that are “invarinat up to scaling function  $K(g)$ ”. Among those, there is unique up to scaling measure which is what is called **right-invariant**: it is invariant under all maps  $R_g : G \rightarrow G$  sending  $g'$  to  $R_g(g') = g'g$  (the uniqueness up to scale and the fact that right invariant measures are “weakly” left invariant are basic facts of the theory of Haar measures on groups).

We illustrate the result on the location-scale example.

(We postpone more the discussion of **why** the right-invariant prior is the best; see section 6.6 in Berger).

---

Example 1 continued: Using  $\theta_0 = (m_0 = 0, s_0 = 1)$  the map  $f : \Theta \rightarrow G$  sends  $(m, s)$  to  $\begin{pmatrix} s & m \\ 0 & 1 \end{pmatrix}$ .

The map  $R_g = R_{(a,b)}$  sends  $\begin{pmatrix} s & m \\ 0 & 1 \end{pmatrix}$  to

$\begin{pmatrix} s & m \\ 0 & 1 \end{pmatrix} \begin{pmatrix} a & b \\ 0 & 1 \end{pmatrix} = \begin{pmatrix} as & bs + m \\ 0 & 1 \end{pmatrix}$ , and thus has Jacobian matrix in  $(s, m)$  coordinates equal to  $\begin{pmatrix} a & 0 \\ b & 1 \end{pmatrix}$ , and determinant  $|J_{R_g}| = a$ , and  $|J_{R_{g^{-1}}}| = \frac{1}{a}$ .

The invariance condition  $R_g * \pi(g') = \pi(g')$  is then

$$\pi(g') = \pi(g^{-1}g') \frac{1}{a}$$

and, setting  $\pi(e) = C$  and  $g' = g$  we get

$$\pi(g) = \frac{C}{a}.$$

This is in fact the “invariant prior” 12.27.

So, to repeat what we said in the beginning, 12.27 is in fact “best” in this problem, but not because 12.18 (which has multiple issues) is better than 12.30, but because the deduction of 12.36 from 12.30 is flawed, and a better analysis of 12.30 (using the framework of invariant decision rules) does indeed lead to 12.27.