

Hyperplane Priors

V. Dose

Centre for Interdisciplinary Plasma Science, Max-Planck-Institut für Plasmaphysik, EURATOM Association, Boltzmannstrasse 2, D-85748 Garching bei München, Germany

Abstract. The requirement of transformation invariance of a probability distribution is employed to derive prior probabilities for the coefficients of the equation describing a hyperplane. In two dimensions, this is a straight line, in three dimensions an ordinary plane etc. We treat the general case of n dimensions and propose a procedure to normalize the resulting distributions in order to make them proper and appropriate for model comparison problems.

INTRODUCTION

Fitting a straight line in two dimensions to a set of experimental data is one of the most important traditional regression problems in physics [1, 2]. The usual procedure employs the least squares method meaning that in a Bayesian sense the prior probability for the coefficients of the straight line is taken flat, improper. This choice is often a reasonable approximation but in a strict sense unacceptable since it violates the requirement of transformation invariance because a flat prior in a given system of coordinates does in general not remain flat after translations and rotations of that system. If the prior probability is uninformative, then we must insist that it remain invariant under such basic transformations. As a first introductory example we shall derive the prior distribution for a straight line in two dimensions parallel to one of the axes of the coordinate system. The only sensible transformation in this case is a translation. Our second introductory example will be the prior probability for the coefficient of a straight line in two dimensions passing through the origin. The only sensible coordinate transform for this case is a rotation of the coordinate system. We show that the resulting prior distribution satisfies intuition. These two special elementary cases constitute the elements of the general case of an $(n - 1)$ dimensional hyperplane in n -dimensional space.

THE TRANSFORMATION INVARIANCE EQUATION

Consider a probability density of a set of parameters $\vec{a}, p(\vec{a})$. $p(\vec{a})da_1 \cdots da_n$ is then an element of probability mass whose value must be independent from the system of coordinates which we use to evaluate its numerical value. Hence, for a different system of coordinates \vec{a}' we must require that

$$p(\vec{a})da_1 \cdots da_n = p(\vec{a}')da'_1 \cdots da'_n \quad (1)$$

yielding the functional equation

$$p(\vec{a}) = p(\vec{a}') \cdot \left| \frac{\partial(a'_1 \cdots a'_n)}{\partial(a_1 \cdots a_n)} \right| \quad (2)$$

where the determinant on the right-hand side is the Jacobian of the transformation $\vec{a} \rightarrow \vec{a}'$. Since any finite transformation $\vec{a} \rightarrow \vec{a}'$ can be constructed from an appropriate sequence of infinitesimal transformations it is sufficient to consider a single infinitesimal transformation. We denote the infinitesimal transformation which maps \vec{a} onto \vec{a}' as $T_\varepsilon(\vec{a})$. Equation (2) becomes then

$$p(\vec{a}) = p(T_\varepsilon(\vec{a})) \cdot \left| \frac{\partial T_\varepsilon(\vec{a})}{\partial \vec{a}} \right| \quad (3)$$

The right-hand side of (3) is a function of ε , the left hand side is independent of ε . Hence, by differentiating (3) by ε we obtain

$$\frac{\partial}{\partial \varepsilon} \left\{ p(T_\varepsilon(\vec{a})) \left| \frac{\partial T_\varepsilon(\vec{a})}{\partial \vec{a}} \right| \right\}_{\varepsilon=0} = 0 \quad (4)$$

the general functional equation which satisfies the requirement of transformation invariance. Well known special cases are obtained for \vec{a} scalar and subject to the similarity transformation $a' = a(1 + \varepsilon)$

$$\frac{\partial}{\partial \varepsilon} \{ p(a(1 + \varepsilon))(1 + \varepsilon) \}_{\varepsilon=0} = 0 \quad (5)$$

$$p'(a) \cdot a + p(a) = 0 \rightarrow p(a) = \frac{1}{a} \quad (6)$$

and, for a subject to a translation $a' = a + \varepsilon$,

$$p'(a) = 0 \rightarrow p(a) = \text{const.} \quad (7)$$

The former is the well known Jeffreys prior for a scale variable, the latter the improper prior for a location variable.

COORDINATE TRANSFORMS

The general problem which we are going to solve is to find a prior distribution for the coefficients \vec{a} of the linear form

$$a_1 x_1 + a_2 x_2 + \cdots + a_n x_n + 1 = 0 \quad (8)$$

In contrast to the examples of scale and location variable priors where the transform concerned directly the parameters, we shall now investigate the consequences of a transformation applied to \vec{x} . First we investigate one of the two simplest cases namely

$$a_1 x_1 + 1 = 0 \quad (9)$$

In two dimensional space this equation describes a straight line parallel to the x_2 axis. The only sensible transformation on x_1 is in this case a translation $x'_1 = x_1 + \varepsilon$. In the transformed coordinate system we have

$$(x_1 + \varepsilon)a'_1 + 1 = 0 \quad . \quad (10)$$

Ordering terms we get

$$x_1 \cdot \frac{a'_1}{1 + \varepsilon a'_1} + 1 = 0 \quad . \quad (11)$$

Comparison with (9) yields the transformation in a which is implied by the translation of x_1

$$a_1 = \frac{a'_1}{1 + \varepsilon a'_1} \rightarrow a' = \frac{a}{1 - a\varepsilon} \approx a + \varepsilon a^2 \quad (12)$$

where the latter equality results from a neglect of all terms of order ε^2 and higher. This is in line with the infinitesimal character of the transformation. The Jacobian of the transformation is readily calculated to be $(1 + 2a\varepsilon)$. Substituting these results into (4) we arrive at

$$\begin{aligned} \frac{d}{d\varepsilon} \{p(a + \varepsilon a^2)(1 + 2a\varepsilon)\}_{\varepsilon=0} &= 0 \quad , \\ p'a^2 + 2ap &= 0 \quad . \end{aligned} \quad (13)$$

Assuming $a \neq 0$, this differential equation may be simplified and easily solved to give

$$p(a) = \frac{1}{a^2} \quad . \quad (14)$$

This distribution is normalizable, because of the restriction $a \neq 0$. The question which cutoff a_{min} to choose will be addressed after treating the general case.

The second elementary example which we will treat, is a straight-line passing through the origin. The equation for this case is

$$ax_1 + x_2 = 0 \quad (15)$$

The only sensible transformation of the \vec{x} -coordinate system is in this case a rotation. Let φ be the angle of rotation between the \vec{x}' and the \vec{x} coordinate system. The transformation is then

$$\begin{aligned} x'_1 &= x_1 \cos \varphi - x_2 \sin \varphi \approx x_1 - \varepsilon x_2 \quad , \\ x'_2 &= x_1 \sin \varphi + x_2 \cos \varphi \approx \varepsilon x_1 + x_2 \quad . \end{aligned} \quad (16)$$

We substitute (16) into (15) and obtain

$$a'(x_1 - \varepsilon x_2) + \varepsilon x_1 + x_2 = 0 \quad . \quad (17)$$

Ordering in terms of x_1 and x_2 yields the implied transformation on a

$$a' = \frac{a + \varepsilon}{1 - a\varepsilon} \approx a + \varepsilon(1 + a^2) \quad . \quad (18)$$

The Jacobian of this transformation is $(1 + 2a\varepsilon)$ yielding the functional equation

$$\frac{\partial}{\partial \varepsilon} \{p(a + \varepsilon(1 + a^2))(1 + 2a\varepsilon)\}_{\varepsilon=0} = 0 \quad (19)$$

and the associated differential equation

$$p'(a^2 + 1) + 2ap = 0 \quad (20)$$

which is readily solved to yield

$$p(a) = \frac{1}{\pi} \frac{1}{1 + a^2} \quad . \quad (21)$$

The factor $1/\pi$ arises of course not from the solution of (20) but from subsequent normalization. This result can easily be shown to conform with intuition. Noting that the meaning of a is $a = -tg\psi$, where ψ is the angle of the straight line with respect to the positive x_1 -axis, the corresponding distribution of ψ is obtained from (2) as $p(\psi) = 1/\pi$. This means that the uninformative prior for a straight line passing through the origin assumes every orientation to be equally probable.

After these two prototypic examples we are well prepared to treat the general case of an $(n - 1)$ dimensional plane in n -dimensional space. The equations for the hyperplane in the primed and unprimed coordinate systems are

$$a_1x_1 + a_2x_2 + \cdots + a_nx_n + 1 = 0 \quad , \quad (22)$$

$$a'_1x'_1 + a'_2x'_2 + \cdots + a'_nx'_n + 1 = 0 \quad . \quad (23)$$

The general rotation in n -dimensional space may be expressed as a sequence of two-axis (j, k) rotations [5]. There exist $\binom{n}{2}$ such rotations e.g., one for $n = 2$, three for $n = 3$, six for $n = 4$ etc. Now we perform one such rotation: $x'_i = x_i$ for all i not equal to either j or k and

$$\begin{aligned} x'_j &= x_j - \varepsilon x_k \quad , \\ x'_k &= \varepsilon x_j + x_k \quad . \end{aligned} \quad (24)$$

Substituting the primed coordinates into (23) and collecting coefficients of x_i yields the implied transformation $a'_i = a_i$ for all i not equal to either j or k and

$$\begin{aligned} a_j = a'_j + \varepsilon a'_k & \longrightarrow a'_j = a_j - \varepsilon a_k \\ a_k = a'_k - \varepsilon a'_j & \longrightarrow a'_k = a_k + \varepsilon a_j \end{aligned} \quad (25)$$

Equation (33) may be simplified if we require $a_k \neq 0$ to

$$\sum_i a_i \frac{\partial p}{\partial a_i} + (n+1)p = 0 \quad . \quad (34)$$

Translation of other coordinates, say x_j , leads to exactly the same differential equation (34) however with side condition $a_j \neq 0$. We now use the previous result (28) that p can only be a function of $\rho = a_1^2 + a_2^2 + \dots + a_n^2$. Then $\partial p / \partial a_i = 2a_i dp / d\rho$ and we arrive at

$$2 \sum_i a_i^2 \frac{dp}{d\rho} + (n+1)p = 0 \quad (35)$$

which is readily integrated to yield

$$p(a_1 \dots a_n) = \frac{1}{\{a_1^2 + \dots + a_n^2\}^{\frac{n+1}{2}}} \quad , |a_i| > 0 \quad . \quad (36)$$

Note that this general result contains our previous simple case (14) for $n = 1$ as well as the priors for the coefficients of a straight-line in two dimensions and an ordinary plane in three dimensions derived by Kendall and Moran [3].

NORMALIZATION

The condition $|a_i| > 0$ in (36) means of course that the distribution is normalizable. Denote by $r^2 = a_1^2 + \dots + a_n^2$ and by r_0 the minimum value of r which we allow. The normalization integral Z is then

$$Z = \int d\Omega_n \int_{r_0}^{\infty} \frac{dr}{r^2} = \frac{2\pi^{n/2}}{\Gamma(n/2)} \cdot \frac{1}{r_0} \quad (37)$$

and the normalized distribution becomes for all $n \geq 1$

$$p(a_1 \dots a_n) = r_0 \frac{\Gamma(n/2)}{2\pi^{n/2}} \cdot \frac{1}{\{a_1^2 + \dots + a_n^2\}^{\frac{n+1}{2}}} \quad . \quad (38)$$

The remaining question concerns of course the choice of r_0 . Let x_1^*, \dots, x_n^* be a point through which the hyperplane is expected to pass, then for this point (22) holds. We now answer the question: "*What is the minimum value of r^2 given this point?*". The Lagrangian optimization problem is to find the minimum of

$$\Phi = a_1^2 + \dots + a_n^2 + 2\lambda(a_1 x_1^* + \dots + a_n x_n^* + 1) \quad . \quad (39)$$

The derivative with respect to a_k yields

$$\frac{\partial \Phi}{\partial a_k} = 2\tilde{a}_k + 2\lambda x_k^* = 0 \quad \rightarrow \quad \tilde{a}_k = -\lambda x_k^* \quad . \quad (40)$$

The derivative with respect to λ yields the hyperplane equation into which we substitute the solution (40) to obtain

$$-\lambda x_1^{*2} \cdots - \lambda x_n^{*2} + 1 = 0 \quad \rightarrow \quad \lambda = \left(\sum x_i^{*2} \right)^{-1} . \quad (41)$$

So the coefficients \tilde{a}_i become explicit and r_0^2 turns out to be

$$r_0^2 = \sum \tilde{a}_i^2 = \left(\sum x_i^{*2} \right)^{-1} . \quad (42)$$

A suitable choice of r_0 can therefore be obtained from that point \vec{x}^* which lies farthest from the origin of the data cloud. Clearly distances are invariant under rotations but not under translations. This suggests to dispose of the origin of the \vec{x} coordinate system. The highest lower limit of r_0 is obviously obtained if the origin is chosen to lie in the centre of gravity of the data vectors \vec{x}_j .

LOWER DIMENSIONAL HYPERPLANES

In the last section we have derived the prior probability for an $(n - 1)$ dimensional hyperplane in n -dimensional space. Though we believe that this is, from a practical point of view, the analysis of experimental data, the by far most important case it remains a matter of principle to show how to deal with lower, $(n - k)$ dimensional hyperplanes. The key to this problem is the parameter representation of a $(k \leq n - 1)$ dimensional vector space. The vector \vec{x} representing a point of this space is written as

$$\vec{x} = \vec{r} + \lambda_1 \vec{a}_1 + \cdots + \lambda_k \vec{a}_k \quad (43)$$

where \vec{r} is the vector which fixes the origin of the vector space $\vec{a}_1, \dots, \vec{a}_k$ which in turn spans the $(n - k)$ dimensional hyperplane. The vector equation (43) summarizes n scalar equations from which we may eliminate the arbitrary coefficients of the linear combination of $\{\vec{a}_i\}, \vec{\lambda}$. This leaves us with $(n - k)$ equations. A straight-line ($k = 1$) in three dimensions is therefore given by two equations. There are obtained by choosing one of the three equations (43) to express λ_1 in terms of $\vec{x}, \vec{r}, \vec{a}_1$ and then eliminate λ_1 from the remaining two equations. Choosing $\lambda_1 = (x_1 - r_1)/a_1$ the set of equations representing this straight-line is

$$\begin{aligned} \left(-\frac{a_2}{a_1} \right) x_1 + x_2 &= \left(r_2 - \frac{a_2}{a_1} r_1 \right) , \\ \left(-\frac{a_3}{a_1} \right) x_1 + x_3 &= \left(r_3 - \frac{a_3}{a_1} r_1 \right) . \end{aligned} \quad (44)$$

A straight-line in three dimensions is therefore characterized by four parameters, the quantities which we have put in brackets in (44). Introducing the definitions $a_2/a_1 = \alpha, a_3/a_1 = \beta, p = r_2 - a_2 r_1/a_1$ and $q = r_3 - a_3 r_1/a_1$ (44) becomes

$$\begin{aligned} x_2 &= \alpha x_1 + p , \\ x_3 &= \beta x_1 + q , \end{aligned} \quad (45)$$

and the transformation invariance prior for α, β, p, q is determined from the transformations implied on α, β, p, q by translating (3x) and rotating (3x) the \vec{x} -coordinate system. We realize from inspection of (45) that there is some equivalence between x_2, x_3 , but operations on x_1 (either translation of x_1 or rotation about x_1) will yield different implied transformations. It remains to be said that the necessary calculations are simple but lengthy. They lead to a prior probability independent of p, q

$$p(\alpha, \beta) = \frac{1}{\pi} \frac{1}{(1 + \alpha^2 + \beta^2)^2} \quad . \quad (46)$$

This result has previously been obtained also by Kendall and Moran [3]. The discussion above leads us to suspect that it will be difficult if not impossible to derive a general result for an $(n - k)$ dimensional hyperplane due to the in-equivalence of the coordinates which results from the process of parameter elimination from the general equation (44).

SUMMARY AND CONCLUSIONS

Prior probabilities for the coefficients of an $(n - 1)$ dimensional vector space have been derived for arbitrary n . A method has been devised to turn these distributions into proper priors which are required for a Bayesian model comparison. Important applications of these prior distributions arise in multivariate linear regression and in neural networks modeling of experimental data. A paper with an application in the latter class of problems is in preparation [4].

REFERENCES

1. Montgomery, D. C., and Peck, E. A., *Introduction to Linear Regression Analysis*, John Wiley & Sons, New York, 1982.
2. Gull, S. F., "Bayesian data analysis: straight line fitting," in *Maximum Entropy and Bayesian Methods*, edited by J. Skilling, Kluwer Academic, Dordrecht, 1989, p.511.
3. Kendall, M. G., and Moran, P. A. P., *Geometrical Probability*, Griffin, London, 1963.
4. von Toussaint, U., Dose, V., and Gori, S., in preperation.
5. Remember the introduction of Euler angles in classical mechanics. See e.g. Landau, L. D. und Lifschitz E. M., *Lehrbuch der theoretischen Physik I*, Akademie Verlag, Berlin, 1962.