

# Some facts about Gaussian linear models

Michael Wilensky

February 2025

## 1 Introduction

The purpose of this document is to expose some basic facts about Gaussian linear models.

## 2 Definitions

Suppose we have some data, denoted with the vector  $d$ , and a linear model for the data with Gaussian noise,

$$d = Ax + n, \tag{1}$$

where  $n$  is a zero-mean Gaussian noise vector,  $x$  is a vector holding the parameters of the linear model (the coefficients), and  $A$  is a matrix holding the predictors (basis functions). Given  $x$ ,  $A$ , and the noise covariance,  $N$ , and some prior information  $\mathcal{I}$ , we have

$$d|x, A, N, \mathcal{I} \sim \mathcal{N}(Ax, N) \tag{2}$$

i.e. the data would be normally distributed about a mean  $Ax$  with covariance  $N$ .

Generally we are interested in the case where  $x$  is unknown, which is related to  $d|a, A, N, \mathcal{I}$  via Bayes' theorem:

$$P(x|d, A, N, \mathcal{I}) = \frac{P(d|x, A, N, \mathcal{I})P(x|A, N, \mathcal{I})}{P(d|A, N, \mathcal{I})}. \tag{3}$$

For the remainder of this document, we will take  $A$ ,  $N$ , and  $\mathcal{I}$  as fixed prior information, and omit conditioning on them in the notation.

Equation 3 implies a need for a prior distribution,  $P(x)$ . In theory, this prior can be anything, but in practice it encodes a particular state of knowledge and therefore the choice of prior has an effect on the outcome of the inference. We will focus on Gaussian priors, eschewing discussions of appropriateness to other works. That is to say, *a priori*,

$$x \sim \mathcal{N}(\mu, C) \tag{4}$$

### 3 Properties of the Posterior

Inferences about  $x$  given  $d$  are captured by the probability distribution of  $x|d$ . First, what is its form? By multiplying  $P(d|x)$  and  $P(x)$ , which are two Gaussian densities, we can see by examining the exponent that the resulting density must be Gaussian. What are its parameters? The terms in the exponent are, for real vectors,<sup>1</sup>

$$-\frac{1}{2} \left( (d - Ax)^T N^{-1} (d - Ax) + (x - \mu)^T C^{-1} (x - \mu) \right) = -\frac{1}{2} \left( x^T (A^T N^{-1} A + C^{-1}) x + a(x) + b \right) \quad (5)$$

where  $a(x)$  is a linear function of  $x$ , and  $b$  is a constant. Their exact forms are not important because all we need to do is figure out what the covariance and mean of our new Gaussian are, which can be determined through a variety of methods, some more painful than others. Importantly, we can read off the posterior covariance,  $C'$ , from Equation 5:

$$C'^{-1} = A^T N^{-1} A + C^{-1}. \quad (6)$$

To get the mean, we take the gradient with respect to  $x$  of the left side of Equation 5 to obtain linear equation that defines the location of the extremum (equivalent to the mean). I'm going to skip over formalities, but in general, one should break this down component-wise and verify the following relationship:

$$C'^{-1} \mu' = A^T N^{-1} d + C^{-1} \mu \quad (7)$$

where  $\mu'$  is the posterior mean. This gives us everything we need to know analytically about the posterior. If this is the extent of the model, then the problem has a fully analytic solution and this is the end of the journey up to explorations of particular instances of the quantities involved. However, often Gaussian linear models appear in larger hierarchical models as conditional distributions of joint posteriors that are not necessarily easily understood analytically. In this case it is useful to know how to generate samples from Gaussian distributions, as part of e.g. a Gibbs sampler.

### 4 Generating Gaussian samples

In general, if  $\mu'$  and  $C'$  are already calculated, one may draw a standard normal sample,  $\omega$ , and form a sample,  $s$ , via

$$s = \mu' + C'^{1/2} \omega \quad (8)$$

where  $C'^{1/2}$  is any matrix,  $M$ , satisfying  $MM^T = C'$ . A common choice is the Cholesky decomposition for its numerical properties but theoretically this is not strictly necessary.

---

<sup>1</sup>which is a space we can always choose to work in when faced with complex vectors – not sure what happens with e.g. quaternions

In typical sampling applications,  $C$  and  $N$  are easily invertible, but  $C'^{-1}$  is not easily invertible to obtain  $C$ . Furthermore, it is also typical in Gibbs sampling applications for  $A$  or  $C$  to be a function of other model parameters not included in  $x$ , meaning an explicit inversion of  $C'^{-1}$  may have to happen for every new sample to use Equation 8 naively. This is often slow and numerically unpredictable compared to various alternatives. More often, we take advantage of the fact that  $C'^{-1}$  is easily calculated and instead use a fast linear solver<sup>2</sup> to instead solve a linear equation for  $s$ . For example, one could solve

$$C'^{-1}s = C'^{-1}\mu' + C'^{-1/2}\omega, \quad (9)$$

which can be written

$$(A^T N^{-1}A + C^{-1})s = A^T N^{-1}d + C^{-1}\mu + C'^{-1/2}\omega. \quad (10)$$

Now, we can play one more trick with the fluctuation term (the one involving  $\omega$ ), to get an equivalent set of samples. In particular, we can instead solve

$$(A^T N^{-1}A + C^{-1})s = A^T N^{-1}d + C^{-1}\mu + A^T N^{-1/2}\omega_0 + C^{-1/2}\omega_1 \quad (11)$$

where  $\omega_0$  and  $\omega_1$  are standard normal random vectors. This is usually preferable because, for example, recalculating  $A^T N^{-1/2}$  for each sample where  $A$  changes is generally easier than recalculating  $(A^T N^{-1}A + C^{-1})^{1/2}$ . To see why this works, first see that

$$C'^{-1}\langle s \rangle = C'^{-1}\mu' \quad (12)$$

i.e.

$$\langle s \rangle = \mu'. \quad (13)$$

where this expectation is over *realizations of  $\omega_0$  and  $\omega_1$* . Then see that the fluctuation terms obey

$$\text{Cov}\left(A^T N^{-1/2}\omega_0 + C^{-1/2}\omega_1, A^T N^{-1/2}\omega_0 + C^{-1/2}\omega_1\right) = C'^{-1} \quad (14)$$

implying

$$\text{Cov}(C'^{-1}s, C'^{-1}s) = C'^{-1}\text{Cov}(s, s)C'^{-1} = C'^{-1} \quad (15)$$

further implying

$$\text{Cov}(s, s) = C'. \quad (16)$$

In some instances, further speedups can be made if  $C'-1$  is sparse, or can be made to be sparse with preconditioning (rumor is that this is usually done by finding an approximate inverse that is easy to compute or does not need to be recomputed ever).

---

<sup>2</sup>for example, NUMPY's linear algebra package wraps LAPACK, which is a well-understood high-performing FORTRAN-based linear algebra package

## 5 Model selection

Consider two linear models,  $A_1$  and  $A_2$ . So long as we have proper priors for the parameters of these models, we can compare the marginal likelihoods of the models fully analytically. If we are willing to assign prior probabilities to the models, we can perform Bayesian model selection. In particular, omitting a bunch of things for brevity, and defining for the  $i$ th model

$$C'_i \equiv N + A_i C_i A_i^T \quad (17)$$

where  $C_i$  is the prior covariance for the  $i$ th model. Let  $\mu_i$  be the prior mean for the  $i$ th model. Then define

$$\chi_i^2 \equiv (d - A_i \mu_i)^T C_i'^{-1} (d - A_i \mu_i). \quad (18)$$

This lets us see that

$$\log \frac{P(A_1|d)}{P(A_2|d)} = -\frac{1}{2} \left( \log \frac{|C'_1|}{|C'_2|} + \chi_1^2 - \chi_2^2 \right). \quad (19)$$

Roughly speaking, this will favor model 1 if it has smaller volume and smaller squared deviations with respect to the prior parameters  $\mu_1$  and  $C_1$  compared to the equivalent quantities involving model 2.