

Lecture 15: Linear estimation theory (Foundations of objective mapping)

Recap

We've been looking at strategies to use our knowledge of the relationship between two data fields in order to infer the parameters that describe this relationship. For example, we used estimates of sea level rise to fit a functional form to the data, and we looked at the correlation between different variables. Now we're going to flip this relationship around to ask how to use prior knowledge about the covariance of two different quantities in order to infer a quantity that we didn't measure. The quintessential example of this comes from the challenge of mapping irregularly spaced Argo data (our measured quantity) onto a regular grid (unmeasured values).

To get started, we'll return to the definitions that we wrote down when we first considered correlation and covariance.

Linear estimation theory

Sometimes we may expect on theoretical grounds that there is a linear relationship between observable variables. In this case, we may want to find the best linear model. Let

$$\hat{y} = \alpha x \quad (1)$$

be the prediction of y where the variables x and y have zero mean

$$\langle x \rangle = \langle y \rangle = 0 \quad (2)$$

Our goal is to choose a value for α that is optimum in some sense. Previously we did this by computing the covariance between x and y , but now we're going to assume that we don't know y . Nonetheless, we'll proceed in the same way that we did when we derived expressions for correlation or least squares fits. Define the error as the difference between the model and the observed y :

$$\epsilon = \hat{y} - y \quad (3)$$

A reasonable optimization criterion is to minimize the **mean square error** (MSE)

$$\langle \epsilon^2 \rangle = \langle (\hat{y} - y)^2 \rangle = \langle (\alpha x - y)^2 \rangle = \alpha^2 \langle x^2 \rangle - 2\alpha \langle xy \rangle + \langle y^2 \rangle. \quad (4)$$

To minimize this with respect to α , we differentiate by α and set the result to zero.

$$\frac{\partial \langle \epsilon^2 \rangle}{\partial \alpha} = 2\alpha \langle x^2 \rangle - 2\langle xy \rangle. \quad (5)$$

Thus α , sometimes called the gain, is

$$\alpha = \frac{\langle xy \rangle}{\langle x^2 \rangle}. \quad (6)$$

This should be familiar: we derived this previously when we looked at regression coefficients.

Recall the definition of the correlation

$$\rho = \frac{\langle xy \rangle}{\sqrt{\langle x^2 \rangle \langle y^2 \rangle}}. \quad (7)$$

With this value of α , the MSE is

$$\langle \epsilon^2 \rangle = \langle y^2 \rangle [1 - \rho^2] = \langle y^2 \rangle - \langle \hat{y}^2 \rangle. \quad (8)$$

So the MSE is simply the difference between the variance in y and the variance explained by the model \hat{y} . As we noted before, the skill is the fraction of variance explained by the model

$$\frac{\langle \hat{y}^2 \rangle}{\langle y^2 \rangle} = \frac{\langle xy \rangle^2}{\langle x^2 \rangle \langle y^2 \rangle} = \rho^2. \quad (9)$$

Note that the error ϵ is uncorrelated with x :

$$\langle \epsilon x \rangle = 0 \quad (10)$$

The statistical model expressed above has a smaller MSE than any other dynamical or statistical model could have given the same data x and variable to predict y . The statistics needed for the prediction (6) are the covariance $\langle xy \rangle$ between the data and the variable that we want to predict, and the variance of the data $\langle x^2 \rangle$. To evaluate the skill (9) of the prediction, we also need the variance of the predicted variable $\langle y^2 \rangle$. As long as these statistics are available, the prediction is statistically optimum. A linear statistical model can also be used to improve the prediction of a dynamical model. Using the output of the dynamical model as the data x , and statistics of the covariance of the model output with predicted variables, a linear statistical model is guaranteed to improve the prediction. This sort of approach is central to weather prediction.

Nonzero mean

So far, we have considered only variables with zero mean (2). If x and y have nonzero means, what would be the best model? Consider the model

$$\hat{y} = \alpha x + \beta, \quad (11)$$

where the constant β is added to take into account nonzero means. The mean square error is

$$\langle \epsilon^2 \rangle = \langle (\alpha x + \beta - y)^2 \rangle. \quad (12)$$

Differentiate (12) with respect to α and β , and set the results to zero

$$\frac{\partial \langle \epsilon^2 \rangle}{\partial \alpha} = 2 \langle x(\alpha x + \beta - y) \rangle = 2\alpha \langle x^2 \rangle + 2\beta \langle x \rangle - 2 \langle xy \rangle \quad (13)$$

$$\frac{\partial \langle \epsilon^2 \rangle}{\partial \beta} = 2 \langle \alpha x + \beta - y \rangle = 2\alpha \langle x \rangle + 2\beta - 2 \langle y \rangle \quad (14)$$

Equations (13-14) are solved for the unknowns α and β :

$$\alpha = \frac{\langle xy \rangle - \langle x \rangle \langle y \rangle}{\langle x^2 \rangle - \langle x \rangle^2} \quad (15)$$

$$\beta = \langle y \rangle - \alpha \langle x \rangle \quad (16)$$

Writing the variables in terms of their means and fluctuations

$$x = \langle x \rangle + x' \quad (17)$$

$$y = \langle y \rangle + y' \quad (18)$$

$$(19)$$

we find the gain to be

$$\alpha = \frac{\langle x' y' \rangle}{\langle x'^2 \rangle}. \quad (20)$$

Then the optimum model is, using (16) in (11)

$$\hat{y} = \alpha(x - \langle x \rangle) + \langle y \rangle \quad (21)$$

or equivalently

$$\hat{y}' = \alpha x'. \quad (22)$$

So the optimum model is just as before for variables with zero mean. Thus, the best approach is always to take the means out of data before proceeding with linear estimation.

Multiple variables

Consider the prediction of a scalar y from a vector of data \mathbf{x} . A linear model would be

$$\hat{y} = \mathbf{a}^T \mathbf{x} \quad (23)$$

where \mathbf{a} is the gain vector. The MSE is

$$\langle \epsilon^2 \rangle = \langle (\hat{y} - y)^2 \rangle = \langle (\mathbf{a}^T \mathbf{x} - y)^2 \rangle. \quad (24)$$

Minimizing

$$\frac{\partial \langle \epsilon^2 \rangle}{\partial \mathbf{a}} = 2\mathbf{a} \langle \mathbf{x} \mathbf{x}^T \rangle - 2 \langle \mathbf{x} y \rangle = 0 \quad (25)$$

$$\mathbf{a} = \langle \mathbf{x} \mathbf{x}^T \rangle^{-1} \langle \mathbf{x} y \rangle. \quad (26)$$

The MSE is then

$$\langle (\hat{y} - y)^2 \rangle = \langle y^2 \rangle - \langle y \mathbf{x}^T \rangle \langle \mathbf{x} \mathbf{x}^T \rangle^{-1} \langle \mathbf{x} y \rangle = \langle y^2 \rangle - \langle \hat{y}^2 \rangle, \quad (27)$$

and the skill is

$$\frac{\langle \hat{y}^2 \rangle}{\langle y^2 \rangle} = \frac{\langle y \mathbf{x}^T \rangle \langle \mathbf{x} \mathbf{x}^T \rangle^{-1} \langle \mathbf{x} y \rangle}{\langle y^2 \rangle}. \quad (28)$$

The solution above obviously requires the data covariance matrix $\langle \mathbf{x} \mathbf{x}^T \rangle$ to be invertible. This is understood as each datum providing at least some new information, even if the data are not uncorrelated with each other. A completely redundant datum would result in a rank deficient, therefore singular, matrix. As real data almost always have some noise, the noise shows up as an augmentation of the data–data covariance matrix along the diagonal. In essence, we add $\sigma^2 \mathbf{I}$ to the

theoretical covariance matrix. The large non-zero diagonal of the covariance matrix means that it is almost always invertible.

The equations derived here describe the core method used to map irregularly spaced data—what oceanographers call “objective mapping” and statisticians call “kriging”. The core principle is that we use a priori knowledge of the data–data and data–model covariances to make inferences about unknown quantities.

Testing with interpolation

Linear interpolation is a good place to start testing our objective mapping procedure. Suppose we measure at two data points: $t = +1$ and $t = -1$. Then our mapped values will depend on the covariances between our observations ($x(+1)$ and $x(-1)$) and between the observations and the mapped quantities.

For testing purposes, we’ll assume a Gaussian covariance of the form

$$\langle x(t_1)x(t_2) \rangle = A \exp \left[\frac{-(t_1 - t_2)^2}{T^2} \right]. \quad (29)$$

You’ll notice that in this case, the covariance only depends on the spatial separation $t_1 - t_2$. We can write this as a matrix.

Since we want to predict y , we’ll also need the covariance between y and x :

$$\langle yx \rangle = A \begin{bmatrix} \exp \left[\frac{-(t-t_1)^2}{T^2} \right] \\ \exp \left[\frac{-(t-t_2)^2}{T^2} \right] \end{bmatrix}. \quad (30)$$

In the homework, you are asked to look at the demo “intgauss.m”, which shows the result of choosing different values of the time scale T . As T increases, the interpolation approaches a straight line, and the skill improves. At short time scales, both the estimate and the skill fade rapidly to zero with distance from the data.

Mapped values relax to zero far from data, so unless our background state is really zero, we want to make sure that we’re removing a mean from the data and that we focus on mapping anomalies.

```
%intgauss.m (Dan Rudnick)

%Linear estimate interpolation using gaussian correlation.
%
% PARAMETER INPUT
d=input('Data (2-vector for values at t=[-1 1])? ');
scale=input('e-folding scales? ');
noise=input('Noise? ');

% INITIALIZE ARRAYS
d=d(:);
t=(-5:0.1:5)';
skill=zeros(length(t),length(scale));
x=zeros(size(skill));
```

```

skillt=zeros(size(skill));
xt=zeros(size(skill));

% DEFINE COVARIANCE
for n=1:length(scale)
    % define data-data covariance matrix (2x2)
    cov=[1+noise exp(-(2/scale(n))^2); exp(-(2/scale(n))^2) 1+noise];
    % define data-model covariance matrix (2xN)--2 values for each location to which
    ct=[exp(-((t+1)/scale(n)).^2) exp(-((t-1)/scale(n)).^2)];

% --baseline solution
skill(:,n)=diag(ct/cov*ct');
x(:,n)=ct/cov*d;

% --covariance matrix and solution for time derivative
ctt=-2/(scale(n).^2)*[t+1 t-1].*ct;
skillt(:,n)=diag(ctt/cov*ctt')/(2/(scale(n).^2));
xt(:,n)=ctt/cov*d;
end

% PLOT RESULTS
% --Figure 1: fitted data and skill
figure;
subplot(2,1,1)
plot(t,x),xlabel('t'),ylabel('x');
subplot(2,1,2)
plot(t,skill),xlabel('t'),ylabel('skill')

% --Figure 2: fitting time derivative instead of raw data
figure;
subplot(2,1,1)
plot(t,xt),xlabel('t'),ylabel('dxdt');
subplot(2,1,2)
plot(t,skillt),xlabel('t'),ylabel('skill-dxdt')

```

or

```

import numpy as np
import scipy
import xarray as xr
import cmocean as cmo
import matplotlib.pyplot as plt
from numpy import linalg as LA

### set these parameters to call intgauss
d=np.array([3,5]) ## data points
scale = [.1,1] ## decorrelation scale

```

```

noise = .1  ## noise

def intgauss(d, scale, noise):
    t=np.arange(-5,5.1,.1)
    skill=np.zeros([len(t),len(scale)])
    x=np.zeros([len(t),len(scale)])
    skillt=np.zeros([len(t),len(scale)])
    xt=np.zeros([len(t),len(scale)])

    cov_matrix=np.zeros([2,2])
    ct=np.zeros([len(t),2])
    ctt=np.zeros([len(t),2])
    delta_t=np.zeros([len(t),2])
    delta_t[:,0]=t+1
    delta_t[:,1]=t-1
    for i in range(len(scale)):
        cov_matrix[0,0]=1+noise
        cov_matrix[0,1]=np.exp(-(2/scale[i])**2)
        cov_matrix[1,0]=cov_matrix[0,1]
        cov_matrix[1,1]=cov_matrix[0,0]

        cov_inverse=LA.inv(cov_matrix)

        ct[:,0]=np.exp(-(t+1)/scale[i])**2)
        ct[:,1]=np.exp(-(t-1)/scale[i])**2)

        ctt=-2/(scale[i]**2)*delta_t*ct

        ctcov=np.matmul(ct,cov_inverse)

        skill[:,i]=np.diagonal(np.matmul(np.matmul(ct,cov_inverse),ct.T))
        x[:,i]=np.matmul(np.matmul(ct,cov_inverse),d.flatten())
        skillt[:,i]=np.diagonal(np.matmul(np.matmul(ctt,cov_inverse),ctt.T))/(2/scale[i])
        xt[:,i]=np.matmul(np.matmul(ctt,cov_inverse),d.flatten())

    return t,x,skill,xt,skillt

t,x,skill,xt,skillt = intgauss(d,scale,noise)

plt.subplot(2,1,1)
plt.plot(t,x)
plt.xlabel('t',fontsize=14)
plt.ylabel('y',fontsize=14)

plt.subplot(2,1,2)
plt.plot(t,skill)

```

```
plt.xlabel('t', fontsize=14)
plt.ylabel('skill', fontsize=14)

plt.subplot(2, 1, 1)
plt.plot(t, xt)
plt.xlabel('t', fontsize=14)
plt.ylabel('dxdt', fontsize=14)

plt.subplot(2, 1, 2)
plt.plot(t, skillt)
plt.xlabel('t', fontsize=14)
plt.ylabel('skill-dxdt', fontsize=14)
```

Mapping with two points

In lecture (and in the homework) we looked at a simplified scenario in which we had two data points and mapped values between them. We initially looked at a problem that matched the code, with data at $t_1 = -1$ and $t_2 = +1$, but let's now make this a little more general by using $t_1 = -\delta$ and $t_2 = +\delta$.

In order to estimate values of $y = x(t)$, we'll need to know covariances. To get started we assume a Gaussian covariance of the form:

$$\langle x(t_1)x(t_2) \rangle = A \exp \left[\frac{-(t_1 - t_2)^2}{T^2} \right] \quad (31)$$

but you'll notice that the covariance has no intrinsic dependence on the time t_1 or t_2 , but depends only on their separation $t_2 - t_1 = \Delta t = \tau$.

We can rewrite the covariance as

$$\rho(\tau) = A \exp \left[\frac{-\tau^2}{T^2} \right]. \quad (32)$$

For our two data points, we can write a 2×2 data–data covariance matrix:

$$\langle \mathbf{x}\mathbf{x}^T \rangle = A \begin{bmatrix} 1 & \exp \left[\frac{-4\delta^2}{T^2} \right] \\ \exp \left[\frac{-4\delta^2}{T^2} \right] & 1 \end{bmatrix}. \quad (33)$$

You'll notice that this matrix is square and symmetric, which is critical since we need it to be invertible.

The data–model covariance matrix is:

$$\langle y\mathbf{x} \rangle = A \begin{bmatrix} \exp \left[\frac{-(t+\delta)^2}{T^2} \right] \\ \exp \left[\frac{-(t-\delta)^2}{T^2} \right] \end{bmatrix}. \quad (34)$$

As homework you were asked to look closely at this problem. Here I just want to highlight a few key issues. In the appendix, I'll look at the limit as δ and t are small, when this resembles interpolation.

In setting this up, we have made some key simplifications, that we should examine in detail.

Noise matters

The symmetric data–data covariance matrix has an Achilles’ heel that we’ll need to watch. Suppose that our decorrelation scale T is really long. That would mean that adjacent points are highly correlated. You can readily see that in such a case, the matrix would approach:

$$\langle \mathbf{x}\mathbf{x}^T \rangle = A \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}, \quad (35)$$

which is singular. What do we do about this?

Appendix: Interpolation in the limit of small data separation

In the two point case, we might ask when linear estimation begins to look like linear interpolation. Linear estimation has all the machinery of a covariance matrix. The result of straight line interpolation for small data separation is quite general regardless of the functional form of the autocovariance, as we now prove. Our estimate of a continuous function can be written as

$$\hat{x}(t) = a_1 x(-\delta) + a_2 x(\delta) \quad (36)$$

where a_1 and a_2 are components of the vector \mathbf{a} in (??). Assuming stationary statistics, the autocorrelation is

$$\rho(\tau) = \frac{\langle x(\tau + t_0)x(t_0) \rangle}{\langle x^2 \rangle}. \quad (37)$$

Using (37) in (??), and inserting into (36), the solution is

$$\hat{x}(t) = \frac{\rho(t + \delta) - \rho(2\delta)\rho(t - \delta)}{1 - \rho^2(2\delta)} x(-\delta) + \frac{\rho(t - \delta) - \rho(2\delta)\rho(t + \delta)}{1 - \rho^2(2\delta)} x(\delta) \quad (38)$$

To make progress, we need to know how the autocovariance behaves at small lag. Start by making the Taylor series expansion of $x(t + t_0)$:

$$x(\tau + t_0) = x(t_0) + \dot{x}(t_0)\tau + \frac{1}{2}\ddot{x}(t_0)\tau^2 + \dots \quad (39)$$

where the dots indicate time derivatives. Using (39) the autocovariance is

$$\langle x(\tau + t_0)x(t_0) \rangle = \langle x^2 \rangle + \langle x\dot{x} \rangle\tau + \frac{1}{2}\langle x\ddot{x}(t_0) \rangle\tau^2 + \dots \quad (40)$$

The assumption of stationarity allows simplification of the second two terms on the righthand side of (43). Consider the covariance of a variable with its time derivative

$$\langle x\dot{x} \rangle = \left\langle x(t_0) \frac{\partial x(t_0)}{\partial t_0} \right\rangle = \frac{1}{2} \frac{\partial}{\partial t_0} \langle x^2(t_0) \rangle = 0. \quad (41)$$

The last equality is a result of stationarity, that the variance is independent of time. Consider the covariance of a variable with its second time derivative

$$\langle x\ddot{x} \rangle = \frac{\partial}{\partial t_0} \langle x(t_0)\dot{x}(t_0) \rangle - \langle \dot{x}(t_0)\dot{x}(t_0) \rangle = -\langle \dot{x}^2 \rangle \quad (42)$$

where stationarity is used to get rid of one of the terms. Using (41-42), the covariance (43) becomes

$$\langle x(\tau + t_0)x(t_0) \rangle = \langle x^2 \rangle - \frac{1}{2} \langle \dot{x}^2 \rangle \tau^2 + \dots \quad (43)$$

For small t , the higher order terms are negligible, and the autocorrelation is

$$\rho(\tau) = 1 - \frac{1}{2} \frac{\langle \dot{x}^2 \rangle}{\langle x^2 \rangle} \tau^2. \quad (44)$$

Substituting (44) into (38) and some simplifying produces

$$\hat{x}(t) = \frac{\delta - t}{2\delta} x(-\delta) + \frac{\delta + t}{2\delta} x(\delta) \quad (45)$$

$$= \frac{x(\delta) + x(-\delta)}{2} + \frac{x(\delta) - x(-\delta)}{2\delta} t \quad (46)$$

which is old-fashioned straight line interpolation.