

**86-631/42-631 Neural Data Analysis**  
**Lecture 09: Linear Regression**

- 1) Correlation describes linear dependence
- 2) Ordinary linear regression
- 3) Multiple linear regression
- 4)  $R^2$

We've spoken quite a bit about neural codes and information – we now know that neurons can carry information about stimuli in a number of different ways. Now comes perhaps a more important question: **given a particular stimulus, can we predict the firing of the neuron?** Or perhaps the alternate question: **given that a neuron is firing in a particular way, how can we predict the stimulus?** We are going to spend today's lecture talking about one particular prediction method: **linear prediction**. Linear regression is one of the most useful tools in all of statistics, and that remains as true for neuroscience as it does for any other field. Today we'll introduce **linear regression**, its relationship to correlation, how to estimate regression parameters, and how to characterize the goodness of fit.

### Correlation describes linear dependence

Suppose we have two random variables,  $X$  and  $Y$ , and we are interested in building the best linear predictor we can of  $Y$  from the values of  $X$ . That is, by constructing a linear function of  $X$  [ $f(X) = \alpha + \beta X$ ], we'd like to get the best possible prediction of  $Y$ .

$$\hat{f}(x) = \alpha + \beta x$$

One way of restating this is that we would like to choose  $\alpha$  and  $\beta$  such that  $f(X)$  is as close to  $Y$  as possible. Or, stated yet another way, we would like to solve the following:

$$\underset{\alpha, \beta}{\operatorname{argmin}} \{E[(Y - \alpha - \beta X)^2]\} = \underset{\alpha, \beta}{\operatorname{argmin}} \{MSE\}$$

find error across all  $\alpha, \beta$ 
expectation  
 $\uparrow$ 
 $\nwarrow$ 
real
prediction

Note: this notation means that we would like to find the values of  $\alpha$  and  $\beta$  such that the Expectation is minimized. We are **minimizing the mean squared error (MSE)**!

$$\begin{aligned} \underset{\alpha, \beta}{\operatorname{argmin}} \{MSE\} &= \underset{\alpha, \beta}{\operatorname{argmin}} \{E[(Y - \alpha - \beta X)^2]\} \\ &= \underset{\alpha, \beta}{\operatorname{argmin}} \{E[Y^2 + \alpha^2 + \beta^2 X^2 - 2\alpha Y - 2\beta XY + 2\alpha\beta X]\} \\ &= \underset{\alpha, \beta}{\operatorname{argmin}} \{E[Y^2] + \alpha^2 + \beta^2 E[X^2] - 2\alpha\mu_Y - 2\beta E[XY] + 2\alpha\beta\mu_X\} \end{aligned}$$

linearity of expectation property  
 $\downarrow$

Recalling that  $\operatorname{Var}[Y] = E[Y^2] - E^2[Y]$ , we see that  $E[Y^2] = \sigma_Y^2 + \mu_Y^2$  and  $E[X^2] = \sigma_X^2 + \mu_X^2$ .

Also, remember that the correlation,  $\rho$ , is defined as  $\rho = \operatorname{Cov}[X, Y] / \sigma_X \sigma_Y$ , so

$$\rho \sigma_X \sigma_Y = \operatorname{Cov}[X, Y] = E[XY] - E[X]E[Y], \text{ and therefore } E[XY] = \rho \sigma_X \sigma_Y + \mu_X \mu_Y.$$

So...

$$\begin{aligned} \underset{\alpha, \beta}{\operatorname{argmin}} \{MSE\} &= \underset{\alpha, \beta}{\operatorname{argmin}} \{\sigma_Y^2 + \mu_Y^2 + \alpha^2 + \beta^2 \sigma_X^2 + \beta^2 \mu_X^2 - 2\alpha\mu_Y - 2\beta\rho\sigma_X\sigma_Y - 2\beta\mu_X\mu_Y + 2\alpha\beta\mu_X\} \end{aligned}$$

$\operatorname{Var}[Y] = E[Y^2] - E^2[Y]$   
 $E[Y^2] = \sigma_Y^2 + \mu_Y^2$   
 $E[X^2] = \sigma_X^2 + \mu_X^2$   
 $\rho = \frac{\operatorname{Cov}[X, Y]}{\sigma_X \sigma_Y}$   
 $\operatorname{Cov}[X, Y] = E[XY] - E[X]E[Y]$   
 $E[XY] = \rho\sigma_X\sigma_Y + \mu_X\mu_Y$

To solve this, we first take the derivative of the equation with respect to  $\alpha$  and set equal to 0:

$$\frac{\partial MSE}{\partial \alpha} = 0 \rightarrow 2\alpha - 2\mu_Y + 2\beta\mu_X = 0$$

$$\alpha = \mu_Y - \beta\mu_X$$

$$\frac{\partial MSE}{\partial \alpha} = 0 \xrightarrow{\text{yields}} 2\alpha - 2\mu_Y + 2\beta\mu_X = 0$$

Or

$$\alpha = \mu_Y - \beta\mu_X$$

Substituting this back in to the original equation yields:

$$\operatorname{argmin}_{\alpha, \beta} \{MSE\}$$

$$= \operatorname{argmin}_{\alpha, \beta} \{\sigma_Y^2 + \mu_Y^2 + (\mu_Y - \beta\mu_X)^2 + \beta^2\sigma_X^2 + \beta^2\mu_X^2 - 2(\mu_Y - \beta\mu_X)\mu_Y - 2\beta\rho\sigma_X\sigma_Y - 2\beta\mu_X\mu_Y + 2(\mu_Y - \beta\mu_X)\beta\mu_X\}$$

Cancelling some terms yields:

$$\operatorname{argmin}_{\alpha, \beta} \{MSE\}$$

$$= \operatorname{argmin}_{\alpha, \beta} \{\sigma_Y^2 + 2\mu_Y^2 - 2\beta\mu_X\mu_Y + \beta^2\mu_X^2 + \beta^2\sigma_X^2 + \beta^2\mu_X^2 - 2\mu_Y^2 + 2\beta\mu_X\mu_Y - 2\beta\rho\sigma_X\sigma_Y - 2\beta\mu_X\mu_Y + 2\beta\mu_X\mu_Y - 2\beta^2\mu_X^2\}$$

$$= \operatorname{argmin}_{\alpha, \beta} \{\sigma_Y^2 + \beta^2\sigma_X^2 - 2\beta\rho\sigma_X\sigma_Y\}$$

$$\frac{\partial MSE}{\partial \beta} = 0 \rightarrow 2\beta\sigma_X^2 - 2\rho\sigma_X\sigma_Y$$

To solve for  $\beta$ , we take the derivative of the equation with respect to  $\beta$  and set equal to 0:

$$\frac{\partial MSE}{\partial \beta} = 0 \xrightarrow{\text{yields}} 2\beta\sigma_X^2 - 2\rho\sigma_X\sigma_Y = 0$$

Or

$$\beta = \frac{\rho\sigma_Y}{\sigma_X}$$

Correlation  
Slope

$$\begin{aligned} 2\beta\sigma_X^2 &= 2\rho\sigma_X\sigma_Y \\ \beta &= \frac{2\rho\sigma_X\sigma_Y}{2\sigma_X^2} \\ &= \frac{\rho\sigma_Y}{\sigma_X} \end{aligned}$$

OK, so what does this all mean? The slope of the linear fit that gives the best prediction of Y from X is proportional to the correlation between X and Y. So, positive correlation yields positive slopes, and negative correlation means negative slope. *proportion to  $\rho$  (linear relationship)*

Further, recall that to get  $\alpha$  and  $\beta$ , we were minimizing the MSE, which boiled down to:  *$+\rho \rightarrow +\text{slope}$   
 $-\rho \rightarrow -\text{slope}$*

$$MSE = \sigma_Y^2 + \beta^2\sigma_X^2 - 2\beta\rho\sigma_X\sigma_Y$$

Plugging our solution for  $\beta$  back in yields:

$$MSE = \sigma_Y^2 + \left(\frac{\rho\sigma_Y}{\sigma_X}\right)^2 \sigma_X^2 - 2\left(\frac{\rho\sigma_Y}{\sigma_X}\right) \rho\sigma_X\sigma_Y$$

$$MSE = \sigma_Y^2 + \rho^2\sigma_Y^2 - 2\rho^2\sigma_Y^2$$

$$MSE = \sigma_Y^2 - \rho^2\sigma_Y^2$$

$$MSE = \sigma_Y^2(1 - \rho^2)$$

If  $\rho = 0$ ,  $MSE = \sigma_Y^2$

As  $\rho \uparrow$ ,  $MSE \downarrow$

If  $\rho = 1, -1$ ,  $MSE = 0$  (perfect line)

The MSE is zero whenever the correlation is  $\pm 1$ : perfect prediction! Meanwhile, the MSE is maximized when the correlation is zero: it's as bad as if we hadn't used  $X$  at all. This is why we say that correlation measures the linear relationship between two random variables.

### Ordinary linear regression

So, how do we get these terms when we're actually dealing with data? Assume we have a bunch of measurements  $y_i$  (for  $i \in [1, N]$ ) in response to some covariate  $x_i$ . For example,  $y_i$  could be the number of spikes recorded in response to a grating with contrast  $x_i$ , and you record a bunch of responses to gratings of different contrasts.

To set up this problem again, we want to find the parameters  $\alpha$  and  $\beta$  such that we minimize the mean squared error between our measurement  $y_i$  and our predicted measurement  $\hat{y}_i$ , where

estimate

$$\hat{y}_i = \alpha + \beta x_i$$

$x_i$  = orientation on trial  $i$

$y_i$  = fr on trial  $i$

$N$  = # of trials

So, you want to minimize:

MSE when dealing w/ actual data

$$MSE = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2 = \frac{1}{N} \sum_{i=1}^N (y_i - \alpha - \beta x_i)^2$$

Solve this on your own!!!

over  $\alpha$  and  $\beta$ .

Solving this in the standard way (taking derivatives w.r.t.  $\alpha$  and  $\beta$  and setting them equal to 0 yields solutions for  $\alpha$  and  $\beta$  exactly as we derived, except using the sample means and variances as opposed to the theoretical means and variances. Thus:

$$\alpha = \bar{y} - \beta \bar{x}$$

and

$$\beta = \frac{\sum_{i=1}^N (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^N (x_i - \bar{x})^2}$$

$$\bar{y} = \frac{1}{N} \sum_{i=1}^N y_i$$

Sample mean

(To see this last equation, realize that  $\text{Cov}[X,Y]=\rho\sigma_X\sigma_Y$ , so  $\beta=\rho\sigma_Y/\sigma_X$  becomes  $\beta=\text{Cov}[X,Y]/\text{Var}[X]$ .)

(A useful exercise is to see if you can prove these two equations are correct.)

## Multiple linear regression

Often we are in a situation where we have several covariates to regress against. For example, we could regress our measured firing rate of our motor cortical neuron against both the speed of the movement and the distance of the movement. In this case, we call the problem a multiple linear regression problem:

$$y^{(i)} = \alpha + \beta_1 x_1^{(i)} + \beta_2 x_2^{(i)} + \dots + \beta_p x_p^{(i)} + \epsilon^{(i)}$$

If we gather all  $n$  measurements into a single vector  $\underline{y} \in \mathbb{R}^{n \times 1}$ , we can rewrite this as

$$\underline{y} = \mathbf{X}\underline{\beta} + \underline{\epsilon}$$

Where  $\mathbf{X} \in \mathbb{R}^{n \times (p+1)}$  is:

$$X = \begin{pmatrix} 1 & X_{11} & X_{12} & \dots & X_{1p} \\ 1 & X_{21} & X_{22} & \dots & X_{2p} \\ & & \dots & & \\ & & \dots & & \\ 1 & X_{n1} & X_{n2} & \dots & X_{np} \end{pmatrix}$$

$\underline{\beta} \in \mathbb{R}^{(p+1) \times 1}$  is:

$$\beta = \begin{pmatrix} \alpha \\ \beta_1 \\ \dots \\ \dots \\ \dots \\ \dots \\ \beta_p \end{pmatrix}$$

and  $\underline{\epsilon}$  is the noise vector that accounts for deviations from the perfect line.

The solution to this problem is:

$$\hat{\underline{\beta}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \underline{y}$$

Note how this reduces to our solution for  $\alpha$  and  $\beta$  in the one-dimensional case:

$$X = \begin{bmatrix} 1 & X_1 \\ \vdots & \vdots \\ 1 & X_n \end{bmatrix}, \quad \text{so} \quad X^T X = \begin{bmatrix} n & \sum X_i \\ \sum X_i & \sum X_i^2 \end{bmatrix}$$

Recall that the inverse of a 2x2 matrix  $\begin{bmatrix} a & b \\ c & d \end{bmatrix}$  is  $\frac{1}{ad-bc} \begin{bmatrix} d & -b \\ -c & a \end{bmatrix}$ .

That makes

$$(X'X)^{-1} = \frac{1}{n \sum X_i^2 - (\sum X_i)^2} \begin{bmatrix} \sum X_i^2 & -\sum X_i \\ -\sum X_i & n \end{bmatrix}$$

$$(X'X)^{-1} = \frac{n}{n \sum X_i^2 - (\sum X_i)^2} \begin{bmatrix} \sum X_i^2/n & -\sum X_i/n \\ -\sum X_i/n & 1 \end{bmatrix}$$

$$(X'X)^{-1} = \frac{1}{\sum X_i^2 - n \left( \frac{\sum X_i}{n} \right)^2} \begin{bmatrix} \sum X_i^2/n & -\bar{X} \\ -\bar{X} & 1 \end{bmatrix}$$

$$(X'X)^{-1} = \frac{1}{\sum X_i^2 - n(\bar{X})^2} \begin{bmatrix} \sum X_i^2/n & -\bar{X} \\ -\bar{X} & 1 \end{bmatrix}$$

Note that the sample variance  $s^2$  is:

$$s_X^2 = \frac{1}{n-1} \sum (X_i - \bar{X})^2$$

so

$$(n-1)s_X^2 = \sum (X_i^2 + \bar{X}^2 - 2X_i\bar{X}) = \sum (X_i^2) + n\bar{X}^2 - 2\bar{X} \sum X_i = \sum (X_i^2) + n\bar{X}^2 - 2n\bar{X}^2$$

and

$$(n-1)s_X^2 = \sum (X_i^2) - n\bar{X}^2$$

Therefore:

$$(X'X)^{-1} = \frac{1/(n-1)}{s_X^2} \begin{bmatrix} \sum X_i^2/n & -\bar{X} \\ -\bar{X} & 1 \end{bmatrix}$$

$$X'Y = \begin{bmatrix} n\bar{Y} \\ \sum X_i Y_i \end{bmatrix} = \begin{bmatrix} n\bar{Y} \\ (n-1)C_{XY} + n\bar{X}\bar{Y} \end{bmatrix}$$

This latter equation comes about because the sample Covariance  $C_{XY}$  is

$$C_{XY} = \frac{1}{n-1} \sum (X_i Y_i - \bar{X}\bar{Y}) \xrightarrow{yields} (n-1)C_{XY} = \sum (X_i Y_i) - n\bar{X}\bar{Y}$$

Therefore:

$$(X'X)^{-1}X'Y = \frac{1}{(n-1)s_X^2} \begin{bmatrix} \sum X_i^2/n & -\bar{X} \\ -\bar{X} & 1 \end{bmatrix} \begin{bmatrix} n\bar{Y} \\ (n-1)C_{XY} + n\bar{X}\bar{Y} \end{bmatrix}$$

$$\frac{1}{(n-1)s_X^2} \begin{bmatrix} \bar{Y} \sum X_i^2 - \bar{X}((n-1)C_{XY} + n\bar{X}\bar{Y}) \\ -n\bar{X}\bar{Y} + (n-1)C_{XY} + n\bar{X}\bar{Y} \end{bmatrix} = \begin{bmatrix} \frac{\bar{Y}(\sum X_i^2 + n\bar{X}^2)}{(n-1)s_X^2} - \frac{\bar{X}C_{XY}}{s_X^2} \\ \frac{C_{XY}}{s_X^2} \end{bmatrix}$$

And so finally we have

$$(X'X)^{-1}X'Y = \begin{bmatrix} \bar{Y} - \bar{X}C_{XY}/s_X^2 \\ C_{XY}/s_X^2 \end{bmatrix} = \begin{bmatrix} \alpha \\ \beta \end{bmatrix}!$$



## R<sup>2</sup>: The coefficient of determination (goodness of fit)

Often, R<sup>2</sup> (called the coefficient of determination) is used as the measurement of the “goodness of fit” of the regression line. The R<sup>2</sup> is defined as:

$$R^2 = \frac{SSR}{SST} \quad SST = SSE + SSR$$

Here, SSR is equal to

Sum Squared Residuals  $\rightarrow SSR = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2$  mean  
sum squared deviations of estimated data from mean

SST is known as the Total Sum of Squares, and is equal to

Total Sum of Squares  $\rightarrow SST = \sum_{i=1}^n (y_i - \bar{y})^2$  mean  
sum squared deviations of sample data from mean

An interesting (and occasionally useful) fact is that the total sum of squares is equal to the SSR plus another term known as the sum squared errors or SSE:

Sum squared errors  $\rightarrow SSE = \sum_{i=1}^n (y_i - \hat{y}_i)^2$  sum squared deviations of sample data from estimated data

That is,  $SST = SSE + SSR$ . To see this, we have the following proof: Sum of squares

$$\begin{aligned} SST &= \sum_{i=1}^n (y_i - \bar{y})^2 = \sum_{i=1}^n (y_i - \hat{y}_i + \hat{y}_i - \bar{y})^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 + 2 \sum_{i=1}^n (y_i - \hat{y}_i)(\hat{y}_i - \bar{y}) \\ (a-b)^2 &= (a-c+c-b)^2 \\ &= (a-c)^2 + (c-b)^2 + 2(a-c)(c-b) \end{aligned}$$
$$SST = SSE + SSR + 2 \sum_{i=1}^n (y_i - \hat{y}_i)(\hat{y}_i - \bar{y})$$

So, it remains to prove that the last term is zero.

Recall that

$$\hat{y}_i = a + bx_i \quad \text{where} \quad a = \bar{y} - b\bar{x}$$

Substitution yields:

$$2 \sum_{i=1}^n (\hat{y}_i - \bar{y})(y_i - \hat{y}_i) = 0$$

$$2 \sum_{i=1}^n (y_i - \hat{y}_i)(\hat{y}_i - \bar{y}) = 2 \sum_{i=1}^n (\bar{y} - b\bar{x} + bx_i - \bar{y})(y_i - \hat{y}_i) = 2 \sum_{i=1}^n b(x_i - \bar{x})(y_i - \hat{y}_i)$$

Substitution again yields

$$\begin{aligned}
2 \sum_{i=1}^n b(x_i - \bar{x})(y_i - \hat{y}_i) &= 2 \sum_{i=1}^n b(x_i - \bar{x})(y_i - (\bar{y} - b\bar{x} + bx_i)) \\
&= 2 \sum_{i=1}^n b(x_i - \bar{x})(y_i - \bar{y} - b(x_i - \bar{x})) = 2 \sum_{i=1}^n b(x_i - \bar{x})(y_i - \bar{y}) - 2 \sum_{i=1}^n b^2(x_i - \bar{x})^2
\end{aligned}$$

$$2b \left( \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) - b \sum_{i=1}^n (x_i - \bar{x})^2 \right)$$

Now, again recall that

$$b = \frac{\sum_{i=1}^N (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^N (x_i - \bar{x})^2}$$

So

$$\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) = b \sum_{i=1}^n (x_i - \bar{x})^2$$

Substitution yields:

$$2b \left( \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) - b \sum_{i=1}^n (x_i - \bar{x})^2 \right) = 2b \left( b \sum_{i=1}^n (x_i - \bar{x})^2 - b \sum_{i=1}^n (x_i - \bar{x})^2 \right) = 0$$

And we indeed have

$$SST = SSE + SSR$$

This is another example of a variance decomposition. SST is proportional to the variance in the data (it is the sample variance multiplied by  $n-1$ ). SSE is proportional to the noise variance (it is multiplied by  $n-2$ ), and SSR can be interpreted as proportional to signal variance. For this reason,  $R^2$  is sometimes thought of as the fraction of explained variance:

$$R^2 = \frac{SSR}{SST} = \frac{SST - SSE}{SST} = 1 - \frac{SSE}{SST}$$

When  $R^2$  is 1, there is no error in the fit – all data lie perfectly on a line, and all of the original variation in the data have been “explained”. When  $R^2$  is zero, you have as much variation in the residuals as you did in the original data! Thus, you’ve explained nothing. Could  $R^2$  be negative? Yes. If you, for example, fit the regression on a training set of data, and then measure the residual variation on a testing set, it’s possible for your fit to be so bad that it does worse than just guessing the mean. In other words, the variation in the residuals is actually *greater than* the variation in the data.

