

Week 5 - Model fitting

The goal of this exercise is to:

- Get familiar with linear regression, GLMs and their applications to neuronal tuning.
- Compute the optimal parameters of the models through maximum likelihood estimation.
- Understand how to evaluate the goodness of fit.

1 Linear regression

One of the fundamental problems in Neuroscience is the modeling of the functional relationship between sensory stimuli and neural activity. Consider the example of seeing an object, which generates specific neural activity in the visual brain areas. Is it possible to construct a model that accurately captures this relationship?

One possible approach consists in assuming a linear relationship between the stimulus and the neural activity. Therefore, we can use linear regression to model the relationship between the neural activity (dependent variable) and the sensory stimulus (independent variable). For instance, we can measure the neural activity of an animal's visual cortex while presenting different visual stimuli. By fitting the neural activity with linear regression, we can obtain an equation that describes how the neural activity y changes in response to the characteristics of different stimuli \mathbf{X} . How can we learn the parameters of the linear regression? How much of the neural activity can we reconstruct with linear regression? Are there better model that can be used to reconstruct the neural activity?

We start addressing these questions by building our first linear regression model. Given the neural activity y of a neuron, we can retrieve its firing rate using the input stimulus x . Since we cannot achieve perfect reconstruction (e.g. measurement noise), we can include everything that our model cannot capture under noise which we initially assume to be Gaussian distributed. Therefore, we can build a Linear-Gaussian regression model:

$$y = \mathbf{X}\theta + \eta, \quad (1)$$

where $y \in \mathbb{R}^N$ with N the number of observations, $\mathbf{X} \in \mathbb{R}^{N,F}$ represents the stimulus matrix with F the number of features, $\theta \in \mathbb{R}^F$ are the linear model parameters and $\eta \in \mathbb{R}^N$ is a multivariate normal Gaussian distribution ($\eta_i \sim \mathcal{N}(0, \sigma)$). The quality of our fitting can be assessed with the mean square error (MSE), i.e. the average squared difference between the predicted values and the actual ones. Therefore, we can find the optimal parameters for our model by minimizing the MSE. In the following exercises, you will need to retrieve the equation for learning the parameters θ .

Exercise 1.1

Find the optimal θ^* by minimizing^a the mean square error (MSE) between the ground truth neural activity (y) and the one reconstructed by the model (\hat{y}):

$$\theta^* = \arg \min_{\theta} MSE(\mathbf{y}, \hat{\mathbf{y}}) \quad (2)$$

^aHint: Remember that the stationary points of a differentiable function can be obtained by setting the derivative equal to zero.

Exercise 1.2

Find the optimal θ^* by computing the maximum likelihood estimation (MLE):

$$\theta^* = \arg \max_{\theta} \log \mathcal{L}(\theta | \mathbf{X}, \mathbf{y}) \quad (3)$$

2 Generalized linear models (GLM)

Previously, we assumed that the probability of a spike occurring at any given time is a linear function of some input variables. However, it is usually assumed that neural spike counts follow a Poisson distribution. In this case, we can build the following model called the linear nonlinear Poisson (LNP) model:

$$\mathbf{y} \sim \text{Poisson}(f(\mathbf{X}\theta)) \quad (4)$$

where $\mathbf{X} \in \mathbb{R}^{N,F}$ represents the stimulus matrix, $\theta\mathbf{X}$ is the linear projections and $f(\cdot) : \mathbb{R} \rightarrow \mathbb{R}$ is an element-wise rectifying nonlinearity, such as the exponential function ($f(k) = e^k$) or the softplus function ($f(k) = \log(1 + e^k)$), where \mathcal{M} is the space of the conditional expectation of y .

Unfortunately, we cannot compute the optimal parameter as we did before for any non linear function because the derivative of the log likelihood has not always a closed form solution. However, we can take advantage of optimization methods such as gradient descent to find the optimal parameters. To this end, we can specify the loss function to optimize our model as the minimization of the negative log likelihood of the spikes (Poisson loss).

Exercise 2.1

Derive the Poisson loss as the negative log likelihood of the spikes given the expected spike counts ($\lambda = f(\mathbf{X}\theta)$) predicted by the model:

3 Model evaluation

To assess the performance of a model in fitting data, various metrics can be used. There are two popular metrics that can be used the coefficient of determination (R^2) and the fraction of explained variance (EV). The goal of these metrics is to quantify how well a model captures the variance in the observed data or how well the predicted values fit the observed values. To calculate these metrics, we start with a set of observed data points (y) and a corresponding set of predicted data points (\hat{y}).

Exercise 3.1

Derive the fraction of explained variance (EV) by computing the variance of the residuals and the total variance of the observed data. The fraction of explained variance is defined as:

$$EV = 1 - \frac{Var(y - \hat{y})}{Var(y)} \quad (5)$$

Exercise 3.2

Derive now the coefficient of determination (R^2) by computing the total sum of squares (TSS), the residual sum of squares (RSS) and the explained sum of squares (ESS)^a. Then, the coefficient of determination is defined as:

$$R^2 = \frac{ESS}{TSS} \quad (6)$$

When can the fraction of explained variance (EV) be considered the same as the coefficient of determination (R^2)?

^aHint: Consider that you need to use the average neural activity \bar{y} for the TSS and the predicted neural activity \hat{y} for the RSS.