

DS 102 Discussion 8

Monday, April 6, 2020

1. Generalized Linear Models

In this discussion, we'll review some of the time-series models presented in Lecture 22 (April 2) and see how they are examples of **generalized linear models** (GLMs). As their name implies, GLMs generalize the linear regression problem we know and love to model situations where 1) the distribution of the output given the input is not simply Gaussian and 2) the mean of that distribution is not simply a linear function of the input.

A quick refresher on the linear regression model we're familiar with: we have d -dimensional input T , and the scalar output X is given by

$$X = \beta^T T + \epsilon \tag{1}$$

where d -dimensional β are parameters, and $\epsilon \sim \mathcal{N}(0, \sigma^2)$ is a noise term that may capture, for example, measurement noise or variance due to other unobserved factors (we use X to denote the output, to be consistent with Section 22.4 of Lecture 22). Equivalently, we can state that the conditional distribution of X given T is:

$$X | T \sim \mathcal{N}(\beta^T T, \sigma^2). \tag{2}$$

There are two key types of assumptions in (2) that GLMs relax.

1. **Output distribution.** In the above linear regression model, we assumed that given the input, the output variable is distributed as a Gaussian random variable. Furthermore, we assumed that the input only determines the mean of this conditional distribution. That is, we have

$$X | T \sim \mathcal{N}(\mu(T), \sigma^2 I) \tag{3}$$

where $\mu(T) := \mathbb{E}[X | T]$ is called the **mean function** because it describes, as a function of T , the mean of the conditional distribution of X given T .

2. **Link function.** In the above linear regression model, we further assumed that the mean function is simply a linear function of the input: $\mu(T) = \beta^T T$.

To relax these two assumptions and model more interesting phenomena, for a GLM we need to specify two components:

1. **Output distribution.** No longer shackled to the Gaussian, we can pick other distributions to model the conditional distribution of the output X given the input T , depending on what is appropriate for the application. However, as in the linear regression case, we keep the assumption that the input only determines the mean $\mathbb{E}[X | T]$ of that distribution, and does not affect any other parameters.

2. **Link function.** To describe the mean function, we can pick any invertible function of a linear function of the inputs. That is, we have

$$g(\mu(T)) = \beta^T T \quad (4)$$

where g is called the **link function**. (In linear regression, the link function is just the identity.)

We'll now review an application similar to the one in Section 22.4 of Lecture 22, and show how it can be cast as a GLM. Let X_t , the observed output variable, denote the number of COVID-19 hospitalizations on day t . We'd like to model the relationship between the output X_t and the input t .

- (a) As noted in the lecture, epidemiology tells us that in some settings, exponential growth for the mean of X_t is reasonable:

$$\mathbb{E}[X_t] = \alpha \exp(\gamma t). \quad (5)$$

Find an appropriate link function. That is, as shown in (4) find a function of the mean function that is equivalent to just a linear function of the input.

- (b) To complete the specification of the GLM, we need an output distribution. Since we are modeling integer-valued X_t , what are natural choices for the conditional distribution of X_T given T ?
- (c) Show how the GLM we've described is equivalent to the model in Section 22.4 of Lecture 22:

$$X_t \sim \text{Poisson}(Z_t) \quad (7)$$

$$Z_{t+1} = (1 + r)Z_t. \quad (8)$$

That is, express α and γ in our GLM as functions of r and Z_0 .

2. **GLM for continuous data.** The examples in Lecture 22 involved discrete data. That is, the observed output variable X_t was always a positive integer (*e.g.*, number of hospitalizations on day t). However, what happens if X_t can be any real number? Then the Poisson GLM wouldn't make much sense as a model.

Consider the following example of a time series with a continuous output. Suppose a rocket has been launched in Florida, and we in California start observing the rocket at time $t = 0$. We want to measure the rocket's distance from Earth at some time t in the future. At each time step, we obtain a noisy measurement of this distance.

Let $\beta_0 \in \mathbb{R}$ be the initial distance from Earth (in miles) of the rocket when we started observing it at time $t = 0$. Suppose the rocket is moving away from Earth at a constant rate of $\beta_1 \in \mathbb{R}$ miles per time step t . Let X_t denote our observation of the rocket's distance from Earth, which is noisy due to weather, measurement error, etc. Assume that for all t our observation noise is normally distributed with a standard deviation of $\sigma = 50$ miles.

- (a) What is the distribution of X_t , the observed distance of the rocket at time t ? Write this distribution in terms of β_0 , β_1 , and σ .
- (b) Suppose we don't know β_0 or β_1 , and we observe X_0, X_1, \dots, X_T positions of the rocket from California. Our goal is to predict the future positions of the rocket by estimating β_0 and β_1 . First, we'll cast our model as a GLM with output X_t and input t . What is the output distribution and link function?
- (c) Given the data X_0, X_1, \dots, X_T , how might we solve for β_0 and β_1 in the above GLM?