

# Lecture Notes

Vyacheslav Lyubchich

2020-11-01



# Contents

1	A short path to GLM and GAM	5
---	-----------------------------	---



# Chapter 1

## A short path to GLM and GAM

This is a brief overview of popular regression models, based on Lyubchich et al. (2019).

Consider a general regression framework:

$$\mathbf{Y} = \mu + \epsilon, \quad (1.1)$$

where  $\mathbf{Y}$  is an  $n \times 1$  column vector comprising observations of the variable of interest (response variable);  $\mu$  is an  $n \times 1$  column vector of expected values  $E(Y_i) \equiv \mu_i$ ;  $\epsilon$  is an  $n \times 1$  column vector of zero-mean random deviations from the expected values,  $i = 1, \dots, n$ ; and  $n$  is the sample size.

In case of a multiple linear regression, the mean response takes the form

$$\mu = \mathbf{X}\beta, \quad (1.2)$$

where  $\mathbf{X}$  is an  $n \times (d + 1)$  matrix with one column of 1's for fitting an intercept in the model and the remaining  $d$  columns for  $d$  correlates (that is, explanatory variables) associated with the response variable;  $\beta$  is a  $(d + 1) \times 1$  column vector of regression coefficients.

The estimation of regression model (1.2) and further inference are based on a *number of assumptions* about the validity of the form of the model (i.e., linearity of relationships between  $Y$  and each  $X$  variable), linear independence of the variables in  $\mathbf{X}$ , relatively equal importance of all the  $n$  observations, as well as uncorrelatedness, homoscedasticity, and normality of errors  $\epsilon$  (Chatterjee and Hadi, 2006). However, the assumption of normality is often violated, and model (1.2) in its classical formulation cannot be used in majority of applied problems.

Generalized linear models (GLMs) help to overcome the violation of normality assumption by extending the applicability of model (1.2) to exponential-type

distributions, such as Poisson, binomial, and gamma (Wood, 2006). In GLMs, distribution of  $Y_i$  belongs to a family of exponential distributions, and a smooth monotonic link function  $g(\cdot)$  is applied to transform the response variable:

$$g(\mu) = \mathbf{X}\beta. \quad (1.3)$$

Canonical link functions are identity,  $\ln$ , and inverse for normal, Poisson, and gamma distributions, respectively. After such transformation, however, model (1.3) still assumes linear relationships between each of the original variables in  $\mathbf{X}$  and the transformed response. Model (1.3) is applicable when the link function successfully linearizes the relationship between the risk variable and a predictor. In other cases, especially if there are multiple predictors, additional work on re-specifying the model may be required. For example, relationships between the response variable and different predictors may require different linearizing transformations, the relationships may be non-monotonic, and many of them may be thresholded (i.e., the effect of a covariate  $X$  is pronounced only when  $X$  takes on values from a certain range, such as the effect of daily precipitation on sediment concentrations in the streams is not noticeable below certain precipitation threshold).

One way we can capture highly non-linear relationships is by inclusion of additional transformed  $X$ -variables, such as power transformed or thresholded variables (e.g.,  $X_i^2$ ;  $\max(0, X_j - a)$ ). However, adding tightly linked variables into the design matrix  $\mathbf{X}$  may introduce multicollinearity and affect the inference. An alternative way of modeling non-linearities is replacing the original variables with those individually transformed using smooth (nonparametric) functions, such as in a generalized additive model (GAM):

$$g(\mu) = \mathbf{X}^*\beta^* + f_1(X_1) + f_2(X_2) + f_3(X_3, X_4) + \dots, \quad (1.4)$$

where  $Y_i$  still follows one of the exponential-family distributions;  $\mathbf{X}^*$  and  $\beta^*$  are the remaining variables and associated coefficients in strictly parametric formulation;  $f(\cdot)$  are smooth functions, often represented by regression splines (Wood, 2006). Model (1.4) can easily deal with deviations from normality and can accommodate non-linearity and non-monotonicity of individual relationships, however, the model still fails to address the issue of remaining dependencies in the errors, e.g., see Kohn et al. (2000).

An extension of model (1.4) by Stasinopoulos and Rigby (2007) to  $k = 1, 2, 3, 4$  parameters  $\theta_k$  of a distribution (not just the location parameter  $\mu_i$ , but also scale  $\sigma_i$ , and shape – skewness and kurtosis; can be generalized for  $k > 4$ ) allows fitting  $k$  individual models

$$g_k(\theta_k) = h_k(\mathbf{X}_k, \beta_k) + \sum_{j=1}^{J_k} h_{jk}(\mathbf{x}_{jk}), \quad (1.5)$$

where  $k = 1$  produces model for the mean;  $h_k(\cdot)$  and  $h_{jk}(\cdot)$  are non-linear functions;  $\beta_k$  is a parameter vector of length  $J_k$ ;  $\mathbf{X}_k$  is an  $n \times J_k$  design matrix;  $\mathbf{x}_{jk}$

are vectors of length  $n$ . The additive terms in this generalized additive model for location scale and shape (GAMLSS) provide a flexible framework to specify random effects and correlation structure as in mixed effects models (Zuur et al., 2009); see Table 3 by Stasinopoulos and Rigby (2007) for other possible specifications of the additive terms. Hence, models of the form (1.5) may be a good choice for insurance problems, because such models accommodate non-normal distributions, possibly highly non-linear relationships, and spatiotemporal dependencies in the data.

Another group of models, called generalized autoregressive moving average (GARMA), was developed by Benjamin et al. (2003) as a combination of GLM (1.3) with Box–Jenkins approach of modeling temporal dependence:

$$g(\mu_t) = \eta_t = \mathbf{X}_t\beta + \sum_{j=1}^p \phi_j \{g(y_{t-j}) - \mathbf{X}_{t-j}\beta\} + \sum_{j=1}^q \theta_j \{g(y_{t-j}) - \eta_{t-j}\}, \quad (1.6)$$

where  $t = 1, \dots, n$  is the time index;  $\phi_j$ ,  $j = 1, \dots, p$ , are autoregressive coefficients;  $\theta_j$ ,  $j = 1, \dots, q$ , are moving average coefficients, and  $p$  and  $q$  are the autoregressive and moving average orders, respectively. Model (1.6) is efficient for dealing with individual time series.

Notice that the issue of different reliability of individual measurements can be solved in models (1.2)–(1.6) by introducing pre-defined weights in the estimation process. An automatic tuning of weights for improved model performance is possible with a number of boosting algorithms, such as AdaBoost.M1 (Hastie et al., 2009).

Overall, model (1.5) is a powerful and flexible choice for a variety of applied problems, when data exhibit complex spatiotemporal dependence and do not adhere to commonly used distributions, such as normal or Poisson.

The challenges of using the above statistical models include the choice of predictors, their transformations, distribution of the response variable, and model specification, which can be attempted with a variety of criteria (for example, Akaike and Bayesian information criteria – AIC and BIC) ubiquitous in statistical literature. Machine learning approaches offer more flexibility by relaxing the assumptions about distributions and forms of relationships, and providing automated solutions for learning meta-features from large amounts of data. At the same time, the large number of tuning parameters that inhere in a machine learning (especially in deep learning) method and their ability of changing the output or extending the computing time dramatically put out a warning for cautious implementation and interpretation of those methods.





# Bibliography

- Benjamin, M. A., Rigby, R. A., and Stasinopoulos, D. M. (2003). Generalized autoregressive moving average models. *Journal of the American Statistical Association*, 98(461):214–223.
- Chatterjee, S. and Hadi, A. S. (2006). *Regression Analysis by Example*. John Wiley & Sons, Hoboken, New Jersey.
- Hastie, T. J., Tibshirani, R. J., and Friedman, J. H. (2009). *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*. Springer, New York, 2 edition.
- Kohn, R., Schimek, M. G., and Smith, M. (2000). *Spline and kernel regression for dependent data*, chapter 6, pages 135–158. John Wiley & Sons, Inc., New York.
- Lyubchich, V., Newlands, N. K., Ghahari, A., Mahdi, T., and Gel, Y. R. (2019). Insurance risk assessment in the face of climate change: integrating data science and statistics. *Wiley Interdisciplinary Reviews: Computational Statistics*, 11:e1462.
- Stasinopoulos, D. M. and Rigby, R. A. (2007). Generalized additive models for location scale and shape (GAMLSS) in R. *Journal of Statistical Software*, 23(7):1–46.
- Wood, S. N. (2006). *Generalized Additive Models: An Introduction with R*. Chapman and Hall/CRC, New York.
- Zuur, A., Ieno, E. N., Walker, N. J., Saveliev, A. A., and Smith, G. M. (2009). *Mixed Effects Models and Extensions in Ecology with R*. Springer, New York.