

Generalization error for Tweedie models: decomposition and error reduction with bagging

Michel Denuit

Institute of Statistics, Biostatistics and Actuarial Science

UCLouvain

Louvain-la-Neuve, Belgium

Julien Trufin

Department of Mathematics

Université Libre de Bruxelles (ULB)

Brussels, Belgium

January 18, 2021

Abstract

Wüthrich and Buser (2020) studied the generalization error for Poisson regression models. This short note aims to extend their results to the Tweedie family of distributions, to which the Poisson law belongs. In case of bagging, a new condition emerges that becomes increasingly binding with the power parameter involved in the Tweedie variance function.

Keywords: Generalization error, Supervised learning, Exponential dispersion family, Tweedie, Bagging.

1 Introduction and motivation

In many applications, the analyst targets the conditional expectation $\mu(\mathbf{X}) = \mathbb{E}[Y|\mathbf{X}]$ of the response Y given the available information summarized in the vector \mathbf{X} . The function $\mathbf{x} \mapsto \mu(\mathbf{x}) = \mathbb{E}[Y|\mathbf{X} = \mathbf{x}]$ is generally unknown and is approximated by an estimator $\mathbf{x} \mapsto \hat{\mu}(\mathbf{x})$. The goal is to produce the most accurate approximation to the true $\hat{\mu}(\mathbf{x})$. Lack of accuracy for $\hat{\mu}(\mathbf{x})$ is defined by the generalization error

$$Err(\hat{\mu}) = \mathbb{E}[L(Y, \hat{\mu}(\mathbf{X}))], \quad (1.1)$$

where $L(., .)$ is a function measuring the discrepancy between its two arguments, called loss function, and the expected value is over the joint distribution of (Y, \mathbf{X}) . We refer the readers to Hastie et al. (2009) for more details about the generalization error. We aim to find a function $\hat{\mu}(\mathbf{x})$ of the features minimizing the generalization error (1.1).

In practice, loss functions $L(., .)$ often correspond to negative log-likelihood functions associated with distributions that belong to the Tweedie family. The Tweedie class regroups the members of the Exponential Dispersion family having power variance functions $V(\mu) = \mu^\xi$ for some ξ . We refer the readers to Denuit et al. (2019) for an extensive treatment of the Exponential Dispersion family and the Tweedie class in the context of insurance. Specifically, the Tweedie class contains continuous distributions such as the Normal, Gamma and Inverse Gaussian distributions. It also includes the Poisson and compound Poisson-Gamma distributions. Compound Poisson-Gamma distributions can be used for modeling data having a positive probability mass at zero and a continuous distribution on the positive real numbers such as yearly insurance losses or rainfall meteorological data.

The following table gives a list of all Tweedie distributions:

	Type	Name
$\xi < 0$	Continuous	-
$\xi = 0$	Continuous	Normal
$0 < \xi < 1$	Non existing	-
$\xi = 1$	Discrete	Poisson
$1 < \xi < 2$	Mixed, non-negative	Compound Poisson-Gamma
$\xi = 2$	Continuous, positive	Gamma
$2 < \xi < 3$	Continuous, positive	-
$\xi = 3$	Continuous, positive	Inverse Gaussian
$\xi > 3$	Continuous, positive	-

Negative values of ξ gives continuous distributions on the whole real axis. There is no probability distribution in the Tweedie class corresponding to power parameters $0 < \xi < 1$. In this paper, we consider non-negative data and we restrict our analysis to $\xi \geq 1$.

We denote by

$$\mathcal{D} = \{(y_1, \mathbf{x}_1), (y_2, \mathbf{x}_2), \dots, (y_n, \mathbf{x}_n)\} \quad (1.2)$$

the set of observations used to fit the model $\hat{\mu}$, called training set. An estimate $\hat{\mu}(\mathbf{x})$ to $\mu(\mathbf{x})$ is obtained by minimizing the total loss on the training set \mathcal{D} , that is, $\sum_{i=1}^n L(y_i, \hat{\mu}(\mathbf{x}_i))$ that can be seen as the empirical version of (1.1). The loss function corresponding to the Tweedie

deviance is

$$L(y, \hat{\mu}(\mathbf{x})) = \begin{cases} 2 \left(y \ln \frac{y}{\hat{\mu}(\mathbf{x})} - (y - \hat{\mu}(\mathbf{x})) \right) & \text{for } \xi = 1 \\ 2 \left(-\ln \frac{y}{\hat{\mu}(\mathbf{x})} + \frac{y}{\hat{\mu}(\mathbf{x})} - 1 \right) & \text{for } \xi = 2 \\ 2 \left(\frac{y^{2-\xi}}{(1-\xi)(2-\xi)} - \frac{y\hat{\mu}(\mathbf{x})^{1-\xi}}{1-\xi} + \frac{\hat{\mu}(\mathbf{x})^{2-\xi}}{2-\xi} \right) & \text{for } \xi \in]1, +\infty[\setminus \{2\}. \end{cases} \quad (1.3)$$

Notice that the Tweedie loss functions (1.3) are particular cases of Bregman loss functions so that the expected loss is minimum for the mean response, whatever the value of the power parameter $\xi \geq 1$, as shown in Savage (1971).

This paper aims to decompose the generalization error (1.1) for loss functions (1.3) derived from Tweedie deviance into a sum of the error for the true model and an estimation error. This extends the classical decomposition that is known to hold for the squared error loss, corresponding to the Normal distribution ($\xi = 0$) and for the loss function derived from the Poisson deviance ($\xi = 1$) as established by Wüthrich and Büser (2020) to the whole Tweedie class with $\xi \geq 1$. The condition under which bagging reduces the error is then obtained, depending on the power parameter ξ . Interestingly, this condition becomes increasingly binding when ξ increases.

2 Generalization error

The generalization error $Err(\hat{\mu})$ given in (1.1) can also be defined for a fixed value $\mathbf{X} = \mathbf{x}$ as

$$Err(\hat{\mu}(\mathbf{x})) = \mathbb{E} [L(Y, \hat{\mu}(\mathbf{X})) | \mathbf{X} = \mathbf{x}]. \quad (2.1)$$

Notice that averaging the local errors $Err(\hat{\mu}(\mathbf{x}))$ enables to recover the generalization error $Err(\hat{\mu})$, that is,

$$Err(\hat{\mu}) = \mathbb{E} [Err(\hat{\mu}(\mathbf{X}))]. \quad (2.2)$$

The generalization error of $\hat{\mu}$ at $\mathbf{X} = \mathbf{x}$ can be expressed as follows.

Proposition 2.1. *We have*

$$Err(\hat{\mu}(\mathbf{x})) = Err(\mu(\mathbf{x})) + \mathcal{E}_\xi(\hat{\mu}(\mathbf{x})) \quad (2.3)$$

with

$$\mathcal{E}_\xi(\hat{\mu}(\mathbf{x})) = \begin{cases} 2\mu(\mathbf{x}) \left(\frac{\hat{\mu}(\mathbf{x})}{\mu(\mathbf{x})} - 1 - \ln \left(\frac{\hat{\mu}(\mathbf{x})}{\mu(\mathbf{x})} \right) \right) & \text{for } \xi = 1 \\ 2 \left(\frac{\mu(\mathbf{x})}{\hat{\mu}(\mathbf{x})} - 1 - \ln \left(\frac{\mu(\mathbf{x})}{\hat{\mu}(\mathbf{x})} \right) \right) & \text{for } \xi = 2 \\ \frac{2}{(2-\xi)(1-\xi)} \left(\hat{\mu}(\mathbf{x})^{2-\xi} (1-\xi) + \mu(\mathbf{x})^{2-\xi} + (\xi-2)\mu(\mathbf{x})\hat{\mu}(\mathbf{x})^{1-\xi} \right) & \text{for } \xi \in]1, +\infty[\setminus \{2\}. \end{cases} \quad (2.4)$$

Proof. For $\xi = 1$, this result can be found in Wüthrich and Büser (2020) (see Equation

(7.5)). Turning to the case $\xi = 2$, we get

$$\begin{aligned}
Err(\widehat{\mu}(\mathbf{x})) &= 2\mathbb{E} \left[-\ln \left(\frac{Y}{\widehat{\mu}(\mathbf{x})} \right) + \frac{Y}{\widehat{\mu}(\mathbf{x})} - 1 \mid \mathbf{X} = \mathbf{x} \right] \\
&= 2\mathbb{E} \left[-\ln \left(\frac{Y}{\mu(\mathbf{x})} \right) + \frac{Y}{\mu(\mathbf{x})} - 1 \mid \mathbf{X} = \mathbf{x} \right] \\
&\quad + 2\mathbb{E} \left[-\ln \left(\frac{Y}{\widehat{\mu}(\mathbf{x})} \right) + \ln \left(\frac{Y}{\mu(\mathbf{x})} \right) + \frac{Y}{\widehat{\mu}(\mathbf{x})} - \frac{Y}{\mu(\mathbf{x})} \mid \mathbf{X} = \mathbf{x} \right] \\
&= Err(\mu(\mathbf{x})) + 2\mathbb{E} \left[\ln \left(\frac{\widehat{\mu}(\mathbf{x})}{\mu(\mathbf{x})} \right) + Y \left(\frac{\mu(\mathbf{x}) - \widehat{\mu}(\mathbf{x})}{\widehat{\mu}(\mathbf{x})\mu(\mathbf{x})} \right) \mid \mathbf{X} = \mathbf{x} \right] \\
&= Err(\mu(\mathbf{x})) + 2 \left(\ln \left(\frac{\widehat{\mu}(\mathbf{x})}{\mu(\mathbf{x})} \right) + \mathbb{E}[Y \mid \mathbf{X} = \mathbf{x}] \left(\frac{\mu(\mathbf{x}) - \widehat{\mu}(\mathbf{x})}{\widehat{\mu}(\mathbf{x})\mu(\mathbf{x})} \right) \right) \\
&= Err(\mu(\mathbf{x})) + 2 \left(\frac{\mu(\mathbf{x})}{\widehat{\mu}(\mathbf{x})} - 1 - \ln \left(\frac{\mu(\mathbf{x})}{\widehat{\mu}(\mathbf{x})} \right) \right).
\end{aligned}$$

Finally, for the remaining cases, we have

$$\begin{aligned}
Err(\widehat{\mu}(\mathbf{x})) &= 2\mathbb{E} \left[\frac{Y^{2-\xi}}{(1-\xi)(2-\xi)} - \frac{Y\widehat{\mu}(\mathbf{x})^{1-\xi}}{1-\xi} + \frac{\widehat{\mu}(\mathbf{x})^{2-\xi}}{2-\xi} \mid \mathbf{X} = \mathbf{x} \right] \\
&= 2\mathbb{E} \left[\frac{Y^{2-\xi}}{(1-\xi)(2-\xi)} - \frac{Y\mu(\mathbf{x})^{1-\xi}}{1-\xi} + \frac{\mu(\mathbf{x})^{2-\xi}}{2-\xi} \mid \mathbf{X} = \mathbf{x} \right] \\
&\quad - 2\mathbb{E} \left[\frac{Y(\widehat{\mu}(\mathbf{x})^{1-\xi} - \mu(\mathbf{x})^{1-\xi})}{1-\xi} - \frac{\widehat{\mu}(\mathbf{x})^{2-\xi} - \mu(\mathbf{x})^{2-\xi}}{2-\xi} \mid \mathbf{X} = \mathbf{x} \right] \\
&= Err(\mu(\mathbf{x})) + \frac{2\widehat{\mu}(\mathbf{x})^{2-\xi} - 2\mu(\mathbf{x})^{2-\xi}}{2-\xi} - \frac{2\mu(\mathbf{x})\widehat{\mu}(\mathbf{x})^{1-\xi}}{1-\xi} + \frac{2\mu(\mathbf{x})^{2-\xi}}{1-\xi}.
\end{aligned}$$

This ends the proof. \square

The smallest generalization error coincides with the one associated to the true model.

Corollary 2.2.

$$Err(\widehat{\mu}(\mathbf{x})) \geq Err(\mu(\mathbf{x})). \quad (2.5)$$

Proof. For $\xi = 1$ and $\xi = 2$, it suffices to notice that $\mathcal{E}_\xi(\widehat{\mu}(\mathbf{x}))$ is always positive since $y \rightarrow y - 1 - \ln y$ is positive on \mathbb{R}^+ . For the remaining cases, define the function f on \mathbb{R}^+ as

$$f(y) = \frac{2}{(2-\xi)(1-\xi)} (y^{2-\xi}(1-\xi) + \mu(\mathbf{x})^{2-\xi} + (\xi-2)\mu(\mathbf{x})y^{1-\xi}). \quad (2.6)$$

We have

$$f'(y) = 2(y^{1-\xi} - \mu(\mathbf{x})y^{-\xi}) \quad \text{and} \quad f''(y) = 2y^{-\xi} \left(1 - \xi + \xi \frac{\mu(\mathbf{x})}{y} \right).$$

Hence, for $y > 0$, $f'(y) = 0$ if, and only if, $y = \mu(\mathbf{x})$ and

$$f''(\mu(\mathbf{x})) = 2\mu(\mathbf{x})^{-\xi} > 0,$$

so that $f(y) \geq f(\mu(\mathbf{x})) = 0$ for all $y > 0$, which completes the proof. \square

The generalization error of $\hat{\mu}$ can be expressed as the sum of the generalization error of the true model μ and an estimation error that is positive. The generalization error of the true model is called the residual error and is thus irreducible. Corollary 2.2 can be found in Gneiting (2011) (see Theorem 7), where the loss functions derived from Tweedie models correspond to the Patton’s family of homogeneous scoring functions on the positive half line.

3 Bagging models and expected generalization error

Bagging is one of the first ensemble methods proposed in the literature by Breiman (1996), who showed that aggregating multiple versions of an estimator into an ensemble improves the model accuracy. Consider a model fitted to our training set \mathcal{D} , obtaining the prediction $\hat{\mu}_{\mathcal{D}}(\mathbf{x})$ at point \mathbf{x} . Bootstrap aggregation or bagging averages this prediction over a set of bootstrap samples in order to reduce its variability with respect to the data used to build it.

The probability distribution of the random vector (Y, \mathbf{X}) is usually not known. It can be approximated by its empirical version which puts an equal probability $\frac{1}{n}$ on each of the observations $\{(y_1, \mathbf{x}_1), (y_2, \mathbf{x}_2), \dots, (y_n, \mathbf{x}_n)\}$ of the training set \mathcal{D} . Hence, instead of simulating B training sets $\mathcal{D}^1, \mathcal{D}^2, \dots, \mathcal{D}^B$ from the probability distribution of (Y, \mathbf{X}) , which is not possible in practice, the idea of bagging is rather to simulate B bootstrap samples $\mathcal{D}^{*1}, \mathcal{D}^{*2}, \dots, \mathcal{D}^{*B}$ of the training set \mathcal{D} from its empirical counterpart. Specifically, a bootstrap sample is thus a random sample of \mathcal{D} taken with replacement which has the same size as \mathcal{D} .

Let $\mathcal{D}^{*1}, \mathcal{D}^{*2}, \dots, \mathcal{D}^{*B}$ be B bootstrap samples of the training set \mathcal{D} . For each \mathcal{D}^{*b} , $b = 1, \dots, B$, we fit a model, giving prediction $\hat{\mu}_{\mathcal{D}^{*b}, \Theta_b}(\mathbf{x}) = \hat{\mu}_{\mathcal{D}^{*b}}(\mathbf{x})$. The bagging prediction is then defined by

$$\hat{\mu}_{\mathcal{D}, \Theta}^{\text{bag}}(\mathbf{x}) = \frac{1}{B} \sum_{b=1}^B \hat{\mu}_{\mathcal{D}^{*b}, \Theta_b}(\mathbf{x}), \quad (3.1)$$

where $\Theta = (\Theta_1, \Theta_2, \dots, \Theta_B)$. Random vectors $\Theta_1, \Theta_2, \dots, \Theta_B$ fully capture the randomness of the training procedure. For bagging, $\Theta_1, \Theta_2, \dots, \Theta_B$ are independent and identically distributed so that Θ_b is a vector of n integers randomly and uniformly drawn in $\{1, 2, \dots, n\}$. Each component of Θ_b indexes one observation of the training set selected in \mathcal{D}^{*b} .

The generalization error $Err(\hat{\mu}_{\mathcal{D}, \Theta}^{\text{bag}}(\mathbf{x}))$ is evaluated conditional on the bootstrap samples $\mathcal{D}^{*1}, \mathcal{D}^{*2}, \dots, \mathcal{D}^{*B}$ of D , so that it gives an idea of the general accuracy of the bagging training procedure for the particular bootstrap samples of \mathcal{D} . In order to assess the general performance of the bagging training procedure, we use the expected generalization error, which averages the generalization error $Err(\hat{\mu}_{\mathcal{D}, \Theta}^{\text{bag}}(\mathbf{x}))$ over \mathcal{D} and Θ , that is, $E_{\mathcal{D}, \Theta} \left[Err \left(\hat{\mu}_{\mathcal{D}, \Theta}^{\text{bag}}(\mathbf{x}) \right) \right]$. From Proposition 2.1, we get

$$E_{\mathcal{D}, \Theta} \left[Err \left(\hat{\mu}_{\mathcal{D}, \Theta}^{\text{bag}}(\mathbf{x}) \right) \right] = Err(\mu(\mathbf{x})) + E_{\mathcal{D}, \Theta} \left[\mathcal{E}_{\xi} \left(\hat{\mu}_{\mathcal{D}, \Theta}^{\text{bag}}(\mathbf{x}) \right) \right]. \quad (3.2)$$

For $\xi = 1$, Wüthrich and Buser (2020) showed that

$$E_{\mathcal{D}, \Theta} \left[Err(\hat{\mu}_{\mathcal{D}, \Theta}^{\text{bag}}(\mathbf{x})) \right] \leq E_{\mathcal{D}, \Theta_b} [Err(\hat{\mu}_{\mathcal{D}, \Theta_b}(\mathbf{x}))],$$

meaning that the bagging prediction $\widehat{\mu}_{\mathcal{D},\Theta}^{\text{bag}}(\mathbf{x})$ outperforms the individual sample estimate $\widehat{\mu}_{\mathcal{D},\Theta_b}(\mathbf{x})$ in the Poisson case. The next proposition extends this result for $\xi \geq 1$ subject to the additional condition (3.3). **Condition (3.3) implies that if the individual sample estimates $\widehat{\mu}_{\mathcal{D},\Theta_b}(\mathbf{x})$ do not overestimate too much the true prediction $\mu(\mathbf{x})$, then it is beneficial to aggregate them in the sense that the aggregation reduces the local generalization error.**

Proposition 3.1. *If the individual sample estimates satisfy*

$$\frac{\widehat{\mu}_{\mathcal{D},\Theta_b}(\mathbf{x})}{\mu(\mathbf{x})} \leq \frac{\xi}{\xi - 1}, \quad (3.3)$$

then we have

$$\mathbb{E}_{\mathcal{D},\Theta} \left[\text{Err}(\widehat{\mu}_{\mathcal{D},\Theta}^{\text{bag}}(\mathbf{x})) \right] \leq \mathbb{E}_{\mathcal{D},\Theta_b} \left[\text{Err}(\widehat{\mu}_{\mathcal{D},\Theta_b}(\mathbf{x})) \right]. \quad (3.4)$$

Proof. The case $\xi = 1$ is due to Wüthrich and Buser (2020) (see Proposition 7.2). Notice that the upper bound in (3.3) is not binding when $\xi = 1$. For $\xi = 2$, by Proposition 2.1, one sees that inequality (3.4) is fulfilled if, and only if,

$$\mathbb{E}_{\mathcal{D},\Theta} \left[\phi \left(\frac{\widehat{\mu}_{\mathcal{D},\Theta}^{\text{bag}}(\mathbf{x})}{\mu(\mathbf{x})} \right) \right] \leq \mathbb{E}_{\mathcal{D},\Theta_b} \left[\phi \left(\frac{\widehat{\mu}_{\mathcal{D},\Theta_b}(\mathbf{x})}{\mu(\mathbf{x})} \right) \right], \quad (3.5)$$

where $\phi : y > 0 \rightarrow \frac{1}{y} + \ln y$. The latter inequality holds true when the individual sample estimates satisfy

$$\frac{\widehat{\mu}_{\mathcal{D},\Theta_b}(\mathbf{x})}{\mu(\mathbf{x})} \leq 2, \quad (3.6)$$

which, in turn, guarantees that $\frac{\widehat{\mu}_{\mathcal{D},\Theta}^{\text{bag}}(\mathbf{x})}{\mu(\mathbf{x})} \leq 2$. Indeed, the function $\phi(y)$ is convex for $y \leq 2$, so that Jensen's inequality implies

$$\begin{aligned} \mathbb{E}_{\mathcal{D},\Theta} \left[\phi \left(\frac{\widehat{\mu}_{\mathcal{D},\Theta}^{\text{bag}}(\mathbf{x})}{\mu(\mathbf{x})} \right) \right] &= \mathbb{E}_{\mathcal{D},\Theta_1,\dots,\Theta_B} \left[\phi \left(\frac{1}{B} \sum_{b=1}^B \frac{\widehat{\mu}_{\mathcal{D},\Theta_b}(\mathbf{x})}{\mu(\mathbf{x})} \right) \right] \\ &\leq \mathbb{E}_{\mathcal{D},\Theta_1,\dots,\Theta_B} \left[\frac{1}{B} \sum_{b=1}^B \phi \left(\frac{\widehat{\mu}_{\mathcal{D},\Theta_b}(\mathbf{x})}{\mu(\mathbf{x})} \right) \right] \\ &= \mathbb{E}_{\mathcal{D},\Theta_b} \left[\phi \left(\frac{\widehat{\mu}_{\mathcal{D},\Theta_b}(\mathbf{x})}{\mu(\mathbf{x})} \right) \right] \end{aligned} \quad (3.7)$$

provided that condition (3.6) holds.

In the remaining cases for ξ , by Proposition 2.1, inequality (3.4) is satisfied if, and only if, $\mathbb{E}_{\mathcal{D},\Theta} \left[f \left(\frac{\widehat{\mu}_{\mathcal{D},\Theta}^{\text{bag}}(\mathbf{x})}{\mu(\mathbf{x})} \right) \right] \leq \mathbb{E}_{\mathcal{D},\Theta_b} \left[f \left(\frac{\widehat{\mu}_{\mathcal{D},\Theta_b}(\mathbf{x})}{\mu(\mathbf{x})} \right) \right]$, where the function $f(\cdot)$ is defined in (2.6). Now, we have $f''(y) \geq 0$ if, and only if, $\frac{y}{\mu(\mathbf{x})} \leq \frac{\xi}{\xi-1}$. Therefore, when $\frac{\widehat{\mu}_{\mathcal{D},\Theta_b}(\mathbf{x})}{\mu(\mathbf{x})} \leq \frac{\xi}{\xi-1}$ and hence $\frac{\widehat{\mu}_{\mathcal{D},\Theta}^{\text{bag}}(\mathbf{x})}{\mu(\mathbf{x})} \leq \frac{\xi}{\xi-1}$, we have $f''(\widehat{\mu}_{\mathcal{D},\Theta_b}(\mathbf{x})) \geq 0$ and $f''(\widehat{\mu}_{\mathcal{D},\Theta}^{\text{bag}}(\mathbf{x})) \geq 0$, so that Jensen's inequality leads to

$$\mathbb{E}_{\mathcal{D},\Theta} \left[f \left(\frac{\widehat{\mu}_{\mathcal{D},\Theta}^{\text{bag}}(\mathbf{x})}{\mu(\mathbf{x})} \right) \right] \leq \mathbb{E}_{\mathcal{D},\Theta_b} \left[f \left(\frac{\widehat{\mu}_{\mathcal{D},\Theta_b}(\mathbf{x})}{\mu(\mathbf{x})} \right) \right].$$

This ends the proof. \square

References

- Breiman, L. (1996). Bagging predictors. *Machine Learning* 24, 123-140.
- Denuit, M., Hainaut, D., Trufin, J. (2019). *Effective Statistical Learning Methods for Actuaries I: GLM and Extensions*. Springer Actuarial Lecture Notes Series.
- Gneiting, T. (2011). Making and evaluating point forecasts. *Journal of the American Statistical Association* 106, 746-762.
- Hastie, T., Tibshirani, R., Friedman, J. (2009). *The Elements of Statistical Learning. Data Mining, Inference, and Prediction*. Second Edition. Springer Series in Statistics.
- Savage, L.J. (1971). Elicitation of personal probabilities and expectations. *Journal of the American Statistical Association* 66, 783-810.
- Wüthrich, M. V., Buser, C. (2020). *Data Analytics for Non-Life Insurance Pricing*. Lecture notes. Available at SSRN, <http://dx.doi.org/10.2139/ssrn.2870308>.