



# Appendix

## Nine measurements

Several measurements are used to evaluate the prediction performance of the regularized random forest. Assume we have true PMI, the response variable  $Y = (y_1, y_2, \dots, y_n)$  where  $n$  is the sample size, and predicted value  $\hat{Y} = (\hat{y}_1, \hat{y}_2, \dots, \hat{y}_n)$ . The difference between  $Y$  and  $\hat{Y}$  is called residual, i.e.,  $r_i = y_i - \hat{y}_i, i = 1, 2, \dots, n$ . We compute the following measurements:

**Mean Squared Error (MSE):**  $MSE = \frac{1}{n} \sum_{i=1}^n r_i^2$

**Median Absolute Deviation (MAD):**  $MAD = median_i (|r_i - median_j (r_j)|)$

**Mean Absolute Percentage Error (MAPE):**  $MAPE = \left\{ \begin{array}{l} \frac{1}{n} \sum_{i=1}^n \left| \frac{r_i}{y_i} \right| \text{ if } y_i \neq 0 \\ \frac{1}{n} \sum_{i=1}^n \left| \frac{r_i + 1}{y_i + 1} \right| \text{ if } y_i = 0 \end{array} \right\}$

**Root Mean Squared Relative Error (RRMSE):**  $RRMSE = \left\{ \begin{array}{l} \sqrt{\frac{1}{n} \sum_{i=1}^n \left( \frac{r_i}{y_i} \right)^2} \text{ if } y_i \neq 0 \\ \sqrt{\frac{1}{n} \sum_{i=1}^n \left( \frac{r_i + 1}{y_i + 1} \right)^2} \text{ if } y_i = 0 \end{array} \right\}$

**Symmetric Mean Absolute Percentage Error (SMAPE):**  $SMAPE = \frac{1}{n} \sum_{i=1}^n \frac{2 \times |r_i|}{|y_i| + |\hat{y}_i|}$

**The Average relative Error (AVGRE):**  $AVGRE = \left\{ \begin{array}{l} \frac{1}{n} \sum_{i=1}^n \left| \frac{r_i}{y_i} \right| \text{ if } y_i \neq 0 \\ \frac{1}{n} \sum_{i=1}^n \left| \frac{r_i + 1}{y_i + 1} \right| \text{ if } y_i = 0 \end{array} \right\}$

**The Maximum Relative Error (MAXRE):**  $MAXRE = \left\{ \begin{array}{l} \max_i \left( \left| \frac{r_i}{y_i} \right| \right) \text{ if } y_i \neq 0 \\ \max_i \left( \left| \frac{r_i + 1}{y_i + 1} \right| \right) \text{ if } y_i = 0 \end{array} \right\}$

**The Total Variation Distance (DTV):**  $DTV = \frac{1}{2} \sum_{i=1}^n |r_i|$

**The Average absolute error or the Mean Difference:**  $MeanDiff = \frac{1}{n} \sum_{i=1}^n |r_i|$