

Appendix

Learning From Mistakes: Machine Learning Enhanced Human Expert Effort Estimates

Federica Sarro, Rebecca Moussa, Alessio Petrozziello and Mark Harman



In this appendix we explain the mathematical formulation of the Linear Programming model we used to predict the MisestimationMagnitude (RQ2) in our paper [1].¹

Linear Programming (LP) [2] aims to achieve the best outcome from a mathematical model with a linear objective function subject to linear equality and inequality constraints. The feasible region is given by the intersection of the constraints and the Simplex (linear programming algorithm) is able to find a point in the polyhedron where the function has the smallest value (minimisation) in polynomial time.

Here, we generalize the model proposed for the effort estimation by Sarro and Petrozziello [3]. In the original implementation, the model is subject to an inequality constraint imposing that the value estimated for each of the observations in the training set has to fall in R_0^+ . Here, we remove the inequality constraints allowing the model to use both positive and negative feature values as well as to optimize for both positive and negative values, as follows:

$$\begin{aligned} \text{minimise} \quad & \sum_{i=1}^n \left| \sum_{j=1}^m a_{ij} x_j - \text{ActualValue}_i \right| \\ & x_j \text{ free}, \quad j = 1, \dots, m \end{aligned} \quad (1)$$

where a_{ij} represents the coefficient of the j^{th} feature for the i^{th} project, x_j is the value of the j^{th} feature, and ActualValue_i is the actual effort of the i^{th} project.

Due to the non-linearity of the absolute value function, the above model has been linearised as follows:

$$\begin{aligned} \text{minimise} \quad & \sum_{i=1}^n t_i \\ \text{subject to} \quad & \sum_{i=1}^n \sum_{j=1}^m a_{ij} x_j - \text{ActualValue}_i - t_i \leq 0 \\ & \sum_{i=1}^n \sum_{j=1}^m a_{ij} x_j - \text{ActualValue}_i + t_i \geq 0 \\ & x_j \text{ free}, \quad j = 1, \dots, m \\ & t_i \text{ free}, \quad i = 1, \dots, n \end{aligned} \quad (2)$$

Let $X_i, \forall i$ be the part of Eq. (1) wrapped in the absolute value. $\forall i$, the slack variable t_i and the following two constraints have been added to the model:

$$\begin{aligned} X_i &\leq t_i \\ -X_i &\leq t_i \end{aligned}$$

Therefore we can have one of the following cases:

- $X_i > 0$ The second constraint, $-X_i \leq t_i$, is always fulfilled as $-X_i$ is negative and t_i is implicitly ≥ 0 . Since t_i is minimised by the objective function and $0 \leq X_i \leq t_i$, the first constraint, $X_i \leq t_i$, is satisfied and t_i is $\text{abs}(X)$.
- $X_i < 0$ The first constraint, $X_i \leq t_i$, is always fulfilled as X_i is negative and t_i is implicitly ≥ 0 . Since t_i is minimised by the objective function and $0 \leq -X_i \leq t_i$, the second constraint, $-X_i \leq t_i$, is satisfied and t_i is $\text{abs}(X)$.
- $X_i = 0$ Both constraints are always fulfilled since t_i is implicitly ≥ 0 . Since t_i is minimised by the objective function, $0 = X_i = t_i$. So t_i is $\text{abs}(X)$.

REFERENCES

- [1] F. Sarro, R. Moussa, A. Petrozziello, and M. Harman, "Learning from mistakes: Machine learning enhanced human expert effort estimates," *IEEE Transactions on Software Engineering*.
- [2] J. C. Nash, "The (dantzig) simplex method for linear programming," *Computing in Science and Engineering*, vol. 2, no. 1, pp. 29–31, 2000.
- [3] F. Sarro and A. Petrozziello, "Linear programming as a baseline for software effort estimation," *ACM Transactions on Software Engineering and Methodology (TOSEM)*, vol. 27, no. 3, pp. 12:1–12:28, 2018.

1. <https://solar.cs.ucl.ac.uk/os/lfm>