1                                        The title

2                            First Author[1] & Ernst-August Doelle[1,2]

3                                  [1] Wilhelm-Wundt-University

4                                  [2] Konstanz Business School

5                                       Author Note

<sub>14</sub>                                                    Abstract

<sub>15</sub> One or two sentences providing a **basic introduction** to the field, comprehensible to a

<sub>16</sub> scientist in any discipline. Two to three sentences of **more detailed background**,

<sub>17</sub> comprehensible to scientists in related disciplines. One sentence clearly stating the **general**

<sub>18</sub> **problem** being addressed by this particular study. One sentence summarizing the main

<sub>19</sub> result (with the words "**here we show**" or their equivalent). Two or three sentences

<sub>20</sub> explaining what the **main result** reveals in direct comparison to what was thought to be

<sub>21</sub> the case previously, or how the main result adds to previous knowledge. One or two

<sub>22</sub> sentences to put the results into a more **general context**. Two or three sentences to

<sub>23</sub> provide a **broader perspective**, readily comprehensible to a scientist in any discipline.

<sub>24</sub>     *Keywords:* keywords

<sub>25</sub>     Word count: X

26      The title

27      **Introduction**

28      Employee turnover is a critical concern for organizations, as it impacts productivity,

29  performance, and overall organizational effectiveness (Blau & Boal, 2020; Griffeth et al.,

30  2000). Accurate prediction of turnover is crucial for proactive human resource management

31  and the implementation of effective retention strategies (Mitchell et al., 2001; Mobley et

32  al., 1979). In recent years, the application of predictive modeling techniques has gained

33  prominence in addressing this challenge. Specifically, the debate arises as to whether

34  regression-based models or machine learning models are more effective in predicting

35  turnover, particularly when working with small sample sizes.

36      The high amount of computer power in the cloud environment nowadays and the

37  developments in the field of machine learning are providing easy access to high-performance

38  services. Machine learning-supported tools are enabling companies to analyze and evaluate

39  information in a quick and effective way (Tambe et al., 2019). We see this in the form of

40  applications, software, and solutions that are common in business or that automate the

41  different decision-making processes, such as programs that create and post job descriptions,

42  application tracking systems that identify key words to place candidates in the right

43  openings, tools for scheduling interviews with your online calendar, chatbots for screening,

44  etc (Rąb-Kettler & Lehnervp, 2019). Human resources practices are not being oblivious to

45  these developments. Experts in these practices are realizing the advantages of data-driven

46  decision making (Fallucchi et al., 2020). Large amounts of human resources data can be

47  analyzed in a short time and empirical inferences can be made, enabling experts to better

48  understand employees and help anticipate issues and patterns (Merlin & Jayam, 2018).

49  Being able to predict the best suited personnel for positioning or that will turn over is of

50  particular interest to human resource departments and companies in general. Making the

51  wrong decision when giving a promotion or demotion can cause waste of time and energy,

52    as well as compromise the perceived organizational justice and support, resulting in more

53    turnover. This is why personnel placement processes are some of the most pivotal in

54    human resources (Merlin & Jayam, 2018).

55         When making decisions for placement using traditional methods, there is a high

56    probability of this being affected by subjective factors that can cause biased choices from

57    time to time (Fallucchi et al., 2020). With machine learning, on the other hand, these

58    decisions are based on a bit more objective foundation than most other recommended

59    methods, since it is based solely on the patterns the algorithm finds on the data (though

60    there can still be bias in the development or implementation of particular algorithms, this

61    is minimized in comparison with traditional methods). Not only this, but the decisions

62    made in the personnel placement process can be explained to the candidates with their

63    reasons , providing them with confidence in the results and diminishing the chances of low

64    perceived organizational justice/support and high turnover (Jha et al., 2020).

65         The aim of this dissertation is to investigate and compare the predictive capabilities

66    of regression-based models and machine learning models in the context of turnover

67    prediction, focusing specifically on sample sizes below 200,000. By examining the strengths

68    and limitations of these modeling approaches, this study seeks to shed light on which

69    method offers greater accuracy and reliability in predicting turnover within

70    resource-constrained environments.

71         Regression-based models, including linear regression, logistic regression, and Cox

72    proportional hazards regression, have long been established as prominent tools in predictive

73    modeling (Hosmer et al., 2013; Cox, 1972). These models are characterized by their

74    simplicity, interpretability, and assumption of linearity between predictors and the outcome

75    variable. The straightforward nature of regression-based models allows for the

76    identification of significant predictors and estimation of their individual effects, facilitating

77    an understanding of the underlying mechanisms driving turnover (Meyer et al., 2004; Hom

78  et al., 2009).

79      Contrarily, machine learning models have garnered significant attention due to their

80  ability to handle complex relationships and patterns in large datasets (Breiman, 2001;

81  Hastie et al., 2009). Algorithms such as random forests, support vector machines, and

82  artificial neural networks offer the potential to capture non-linear and interactive effects,

83  making them valuable tools in predictive modeling (Niculescu-Mizil & Caruana, 2005;

84  Zhang & Singer, 2010). Machine learning models have been increasingly applied to

85  turnover prediction, displaying promising results in various studies (Biswas et al., 2019;

86  Jaradeh & Dehghan, 2021).

87      While the application of machine learning models has gained momentum, their

88  performance in the context of small sample sizes remains an open question. The literature

89  suggests that machine learning models may face challenges, such as overfitting, when

90  trained on limited data (Jiang et al., 2020; Varoquaux et al., 2018). Consequently, the

91  predictive performance of these models might be compromised when sample sizes are below

92  a certain threshold. Thus, it becomes imperative to evaluate whether regression-based

93  models, with their simplicity and interpretability, outperform machine learning models

94  when the sample size is below 200,000.

95      Additionally, the performance of machine learning models may be affected when the

96  number of independent variables is small. In such scenarios, these models may face

97  challenges such as overfitting or difficulty in identifying meaningful patterns (Varoquaux et

98  al., 2018; Guyon & Elisseeff, 2003). On the other hand, regression-based models, with their

99  simplicity and interpretability, may offer advantages in situations where the number of

100 independent variables is limited, as they are less prone to overfitting and can provide

101 transparent insights into the relationships between predictors and turnover (Pedersen &

102 Skogstad, 2020; Hosmer et al., 2013).

103     Traditional variables in the demographic and biodata domain, such as age, gender,

104  and education, have been commonly used in turnover prediction models (Hom et al., 2009;

105  Meyer et al., 2004). However, the utilization of antecedent variables typically studied in

106  I-O psychology, such as job satisfaction, organizational commitment, and work-life balance,

107  may provide deeper insights into the underlying factors contributing to turnover (Lee et al.,

108  2019; Hom et al., 2009).

109      The application of I-O psychology antecedent variables in machine learning-based

110  models holds promise for improving turnover prediction accuracy. These variables capture

111  psychological and organizational aspects that directly impact employees' turnover

112  intentions and behaviors (Lee et al., 2019; Meyer et al., 2004). By considering these

113  variables in predictive models, organizations can gain a more comprehensive understanding

114  of the complex dynamics that drive turnover and develop targeted interventions to mitigate

115  it (Griffeth et al., 2000; Hom et al., 2009).

116      Contrarily, models relying solely on demographics and biodata may overlook critical

117  factors contributing to turnover. While these variables provide basic demographic

118  information, they may lack the depth and specificity necessary to capture the nuances and

119  complexities of turnover behavior (Hom et al., 2009). Incorporating I-O psychology

120  antecedent variables can offer a more nuanced and accurate prediction of turnover by

121  considering individual attitudes, perceptions, and experiences within the organizational

122  context (Meyer et al., 2004; Lee et al., 2019).

123      This study aims to address this research gap by employing a comprehensive dataset

124  from multiple organizations. By leveraging turnover data and a restricted set of

125  independent variables, along with relevant predictor variables such as

126  demographics/biodata, job characteristics, employee engagement, and other I-O psychology

127  antecedent variables of turnover, a comparative analysis will be conducted. The

128  performance of regression-based models and machine learning models will be assessed using

129  various metrics, including accuracy, precision, recall, and the area under the receiver

¹³⁰ operating characteristic curve (AUC-ROC) (Davis & Goadrich, 2006; Saito & Rehmsmeier,

¹³¹ 2015).

¹³²       The findings from this research will contribute to the existing literature on turnover

¹³³ prediction and provide valuable insights for practitioners and researchers alike.

¹³⁴ Understanding the relative performance of regression-based models and machine learning

¹³⁵ models when dealing with small sample sizes can guide decision-making regarding the

¹³⁶ choice of modeling techniques in resource-limited scenarios. Ultimately, this research aims

¹³⁷ to enhance our understanding of turnover prediction and inform effective retention

¹³⁸ strategies to mitigate the negative consequences of employee turnover.

## What Constitutes Machine Learning

¹⁴⁰       Definitions of what constitutes machine learning (ML) and the differences with

¹⁴¹ statistical modeling have been discussed at length in the literature (Breiman, 2001), yet the

¹⁴² distinction is not clear-cut (Moons, 2014). The seminal reference on this issue is Breiman's

¹⁴³ review of the "two cultures" (Breiman, 2001). Breiman contrasts theory-based models such

¹⁴⁴ as regression with empirical algorithms such as decision trees, artificial neural networks,

¹⁴⁵ support vector machines, and random forests.

## Theory-based models

¹⁴⁷       Theory-based models are models that are based on theory and assumptions, such as

¹⁴⁸ traditional linear regression, and benefit from human intervention and subject knowledge

¹⁴⁹ for model specification. The analysis in this approach starts with assuming a stochastic

¹⁵⁰ data model for the inside of the black box. Usually, research that uses this approach starts

¹⁵¹ by assuming that the data are generated by a particular model. This model is used as a

¹⁵² template for statistical analysis. When faced with an applied problem, researchers that use

¹⁵³ this approach come up with a data model by looking at the literature developed by

previous scholars or by their own theorizing, or some combination of both. This enables them to develop a reasonably good parametric class of models for a complex mechanism devised by nature, and then parameters are estimated and conclusions are drawn. However, these conclusions are about the model's mechanism, not about nature's mechanism, and therefore if the model is a poor emulation of nature, the conclusion could be wrong. Breiman (2001) criticized this approach, pointing out that: "A few decades ago, the commitment to data models was such that even simple precaution such as residual analysis or goodness-of-fit tests were not used. The belief in the infallibility of the data models was almost religious. It is a strange phenomenon – once a model is made, then it becomes truth and the conclusions from it are infallible (p. 202)." He concludes by using the following old saying: "If all a man has is a hammer, then every problem looks like a nail (p. 202)." To solve a wide range of problems, such as is the case in the social sciences with the abundance of variables, a larger set of tools is needed. The rapidly increasing ability of computers to store and manipulate data can provide us with more varied tools. Empirical-based models In the mid-1980s, with the development of neural networks and decision trees, a new community of researchers appeared focused on predictive accuracy (Cristianini, 2002). They began using these algorithms on working in complex prediction problems where it was obvious that data models were not applicable, such as speech recognition, image recognition, handwriting recognition, times series analysis, or financial market analysis. The approach is that nature produces data in a black box whose insides are complex and partly unknowable. The goal is not to explain the patterns in this data, but to predict them based on input; not to focus on data models, but on the characteristics of the algorithms (Breiman, 2001). Within psychology, this is the same approach as dustbowl empiricism (Schoenfeldt, 1999). A useful definition of machine learning is that it focuses on models that directly and automatically learn from data (Mitchell & Mitchell, 1997). For example, machine learning performs modeling more automatically than regression regarding the inclusion of nonlinear associations and interaction terms

(Boulesteix, 2014). To do so, machine learning algorithms are often highly flexible

algorithms that require penalization to avoid overfitting (Deo, & Nallamothu, 2016). Some

researchers describe the distinction between statistical modeling and machine learning as a

continuum (Beam, & Kohane, 2018). Other researchers label any method that deviates

from basic regression models as machine learning (He, & Garcia, 2009), such as penalized

regression (e.g., LASSO, elastic net) or generalized additive models (GAM). We note that

these methods do not belong to machine learning by using the "automatic learning from

data" definition, and did not classify these as machine learning in this study

**Data Hungriness**

The concept of data hungriness refers to the sample size needed for a method to

generate a prediction model with a good predictive accuracy (van der Ploeg et al., 2014).

The data hungriness of a predictive modeling technique is defined as the minimum number

of events per variable at which the optimism of the generated model is less than 0.01.

Optimism is defined as the difference between error on our sample data and the error when

applying our model to another dataset. Every machine learning model has some amount of

error in its predictions. This error usually comes from two different sources: bias and

variance. Bias is the tendency of the model to underfit, and variance is the tendency to

overfit. The relationship between these two sources of error is known as the bias-variance

tradeoff, and developers of machine learning models have to find the balance between the

two.

To test if this trade-off has been done in a way that minimizes error, it is a good idea

to measure performance using data that the model has never seen before. The performance

of the model on this "test data" will be a more accurate predictor of the model's

performance in the real world, which is the fundamental basis for cross-validation . The

model's optimism, therefore, is the difference between the training error estimated from the

data used to build the model and the test error estimated by applying the model into

out-of-sample data. The sample size needed to minimize this difference is what we call data

hungriness. Machine learning algorithms require big sample sizes to minimize this

difference (van der Ploeg et al., 2014).

**Parsimony**

Parsimony is defined as the sample size and number of variables that a dataset must

have in order to maximize the predictive accuracy of a model (Sanchez-Pint et al., 2018).

A model is considered parsimonious when it both uses the least amount of variables

possible (sparsity) and has good prediction accuracy. Typically, parsimony is reported as

the performance metric of models that are sparse.

**Hypothesis**

Hypothesis 1: Does regression cross-validate better when the sample size is below

200,000 compared to ML Hypothesis 2: Does regression cross-validate better when the

number of variables is below 200 compared to ML Hypothesis 3: Do ML-based predictive

models have better predictions of turnover when using IO IVs VS other IVs.

## Methods

We report how we determined our sample size, all data exclusions (if any), all

manipulations, and all measures in the study.

## Participants

## Material

## Procedure

## Data analysis

We used R (Version 4.2.2; R Core Team, 2022) and the R-packages *papaja* (Version 0.1.1.9001; Aust & Barth, 2022), and *tinylabels* (Version 0.2.3; Barth, 2022) for all our analyses.

# Results

# Discussion

## References

Aust, F., & Barth, M. (2022). *papaja: Prepare reproducible APA journal articles with R Markdown.* Retrieved from https://github.com/crsh/papaja

Barth, M. (2022). *tinylabels: Lightweight variable labels.* Retrieved from https://cran.r-project.org/package=tinylabels

R Core Team. (2022). *R: A language and environment for statistical computing.* Vienna, Austria: R Foundation for Statistical Computing. Retrieved from https://www.R-project.org/