


Regularized target encoding outperforms traditional methods in supervised machine learning with high cardinality features

Florian Pargent  · Florian Pfisterer  ·
Janek Thomas  · Bernd Bischl 

Received: date / Accepted: date

Abstract Since most machine learning (ML) algorithms are designed for numerical inputs, efficiently encoding categorical variables is a crucial aspect in data analysis. A common problem are high cardinality features, i.e. unordered categorical predictor variables with a high number of levels. We study techniques that yield numeric representations of categorical variables which can then be used in subsequent ML applications. We focus on the impact of these techniques on a subsequent algorithm's predictive performance, and – if possible – derive best practices on when to use which technique. We conducted a large-scale benchmark experiment, where we compared different encoding strategies together with five ML algorithms (lasso, random forest, gradient boosting, k-nearest neighbors, support vector machine) using datasets from regression, binary- and multiclass- classification settings. In our study, regularized versions of target encoding (i.e. using target predictions based on the feature levels in the training set as a new numerical feature) consistently provided the best results. Traditionally widely used encodings that make un-


Florian Pargent 

Department of Psychology, Psychological Methods and Assessment, LMU Munich, Leopoldstrasse 13, 80802 Munich, Germany

E-mail: florian.pargent@psy.lmu.de

Florian Pfisterer 

Department of Statistics, Statistical Learning and Data Science, LMU Munich, Ludwigstraße 33, 80539 Munich, Germany

Janek Thomas 

Fraunhofer Institute for Integrated Circuits, Hansastrasse 32, 8068 Munich, Germany

Bernd Bischl 

Department of Statistics, Statistical Learning and Data Science, LMU Munich, Ludwigstraße 33, 80539 Munich, Germany

reasonable assumptions to map levels to integers (e.g. integer encoding) or to reduce the number of levels (possibly based on target information, e.g. leaf encoding) before creating binary indicator variables (one-hot or dummy encoding) were not as effective in comparison.

Keywords supervised machine learning · benchmark · high-cardinality categorical features · target encoding · dummy encoding · generalized linear mixed models ·

1 Introduction

While increasing sample size is usually considered the most important step to improve the predictive performance of a machine learning (ML) model, using effective feature engineering comes as a close second. One remaining challenge is how to handle high cardinality features – categorical predictor variables with a high number of different levels but without any natural ordering. While categorical variables with only a small number of possible levels can often be efficiently dealt with using standard techniques such as one-hot encoding, this approach becomes inefficient as the number of levels increases. Although domain knowledge can sometimes be used to reduce the number of theoretically relevant levels, finding strategies that work well on a large variety of problems is highly important for many applications as well as in automated ML (Feurer et al., 2015; Thomas et al., 2018; Thornton et al., 2013). Optimally, strategies should be model-agnostic because benchmarking encoding methods together with ML algorithms from different classes is often necessary for applications. While a variety of strategies exist, there are very few benchmarks that can be used to decide which technique is expected to yield good predictive performance. Furthermore, there has recently been increasing attention on scientific benchmark studies that compare different methods to provide a clearer picture in light of a large number of methods available to practitioners (Bommert et al., 2020; Fernández-Delgado et al., 2014), as they can provide at least partial answers to such questions. The goal of this study is to provide an overview of existing approaches for encoding categorical predictor variables and to study their effect on a model’s predictive performance. In order to achieve this, we carefully design a benchmark setting allowing us to discern the effect of encoding strategies (and their hyperparameters) along with their interaction with subsequently applied ML methods from an unbiased perspective. Since we do not introduce a new competitor method, we do not risk cherry picking [Dehghani et al(2021)Dehghani, Tay, Gritsenko, Zhao, Houlsby, Diaz, Metzler, and Vinyals] or reporting over-optimistic results [Nießl et al(2021)Nießl, Herrmann, Wiedemann, Casalicchio, and Boulesteix].

1.1 Notation

We consider the classical setting of supervised learning from an *i.i.d.* tabular dataset \mathcal{D} of size N sampled from a joint distribution $\mathcal{P}(\mathbf{x}, y)$ of a set of features \mathbf{x} and an associated target variable y . Here, \mathbf{x} consists of a mix of numeric (real-valued or integer-valued) features and categorical features, the latter of which we seek to transform feature-wise to numeric features using a categorical encoding technique. Let x be a single unordered categorical feature from a feature space \mathcal{X} with cardinality $\text{card}(\mathcal{X}) \leq \text{card}(\mathbb{N})$. It holds either $y \in \mathbb{R}$ (regression), $y \in \mathcal{C}$ from a finite class space $\mathcal{C} = \{c_1, \dots, c_C\}$ with $C = 2$ (binary classification) or $C > 2$ (multiclass classification). We always assume to observe all C classes in our training sample, however we might only observe a subset $\mathcal{L}^{\text{train}} \subseteq \mathcal{X}$ of a feature’s available L levels, $\mathcal{L}^{\text{train}} = \{l_1, \dots, l_L\}$ for categorical features. We denote the observed frequency of class c in the training set with N_c and the observed frequency of a level l in the training set with N_l . We investigate categorical encoding techniques to transform each nominal feature x^{train} into numerical features \hat{x}^{train} which are then used for training. If clear from the context, we use \hat{x}_l as the encoded value for an observation with level l . Although datasets might contain multiple high cardinality features, we encode each feature separately but with the same strategy.

1.2 Related Work

We broadly categorize feature encoding techniques into *target-agnostic* methods and *target-based* methods (Micci-Barreca, 2001). Figure 1 contains a taxonomy of our considered encoding methods. *Target-agnostic* methods do not rely on any information about the target variable and can therefore also be used in unsupervised settings. Simple strategies from this domain e.g. one-hot or dummy encoding are widely used – in the scientific literature (Hancock and Khoshgoftaar, 2020; Kuhn and Johnson, 2019) but also on Kaggle¹ to embed variables for classical ML algorithms as well as (deep) neural networks. Such indicator methods map each level of a categorical variable to a set of dichotomous features encoding the presence or absence of a particular level. An obvious drawback of indicator encoding is that it adds one additional feature per level of a categorical variable. When indicator encoding leads to an unreasonable number of features, levels are often mapped to integer values with random order (integer encoding). Alternatively, the “hashing trick” (Weinberger et al., 2009) can be used to randomly collapse feature levels into a smaller number of indicator variables (Kuhn and Johnson, 2019), or levels can be encoded by using the observed frequency of a given level in the dataset (frequency encoding).

Target-based methods try to incorporate information about the target values associated with a given level. Early strategies aimed to reduce the number of levels by methods like hierarchical clustering or decision trees based on

¹ <https://www.kaggle.com>

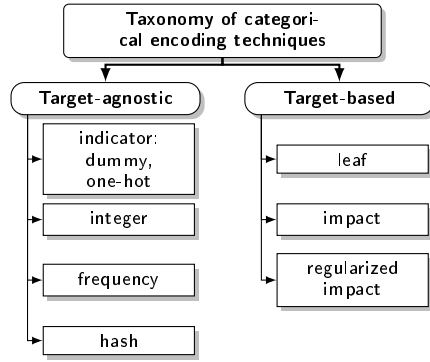


Fig. 1 Taxonomy of common categorical variable encoding techniques.

statistics of the target variable, although this has been rarely described in the scientific literature (for a brief mention, see Micci-Barreca, 2001). The basic idea of more advanced methods called target, impact, mean, or likelihood encoding is to use the training set to make a simple prediction of the target for each level of the categorical feature, and to use the prediction as the numerical feature value \hat{x}_l for the respective level. An early formal description of this strategy is Micci-Barreca (2001). In simple target encoding for regression problems, the mean target value in the training set from all observations with a certain feature level is used to encode that level for all observations: $\hat{x}_l = \frac{\sum_{i: x_i^{train}=l} y_i^{train}}{N_l}$. Simple target encoding often does not perform well with rare levels, where it tends to overfit to the training data and fails to generalize well for new observations. In the extreme case of a categorical feature with unique values (e.g. some hashed ID variable) studied in Prokhorenkova et al. (2018), the mean target for each level of this feature in simple target encoding is similar to the true target value of a single observation. Based on the encoded feature, all observations can be predicted perfectly in the training set, even if the original variable did not contain any useful information. ML models would place a high priority on such an encoded feature during training but would perform badly on test data. To avoid this, practitioners often use regularized target encoding with a smoothing parameter that shrinks those effects towards the global mean (Micci-Barreca, 2001). An alternative strategy is to combine target encoding with cross-validation (CV) techniques (Prokhorenkova et al., 2018).

1.3 Categorical Encoding Benchmarks

Several small-scale studies have been previously conducted. Those works did not yield conclusive results due to narrower scopes or not considering high cardinality variables. One benchmark (6 datasets) on encoding high cardinality features (cardinality between 103 and 9095) in combination with gradient boosting has been published on the Kaggle forums (Prokopen, 2018). Different

versions of target encoding are compared with indicator, integer, and frequency encoding. They recommend combining the smoothed version of target encoding with 4- or 5-fold CV, never using simple target encoding and using indicator encoding only for small datasets. Interestingly, frequency encoding did perform well in many cases. Coors (2018) performed a benchmark (12 datasets) on encoding high cardinality features (maximum number of levels per dataset between 10 and 25847) when developing the automatic gradient boosting (autogb) library (Thomas et al., 2018). They compared different variants of target encoding with integer and indicator encoding. In their benchmarks, target encoding only improved over target-agnostic methods on 2 datasets, while it led to worse results on 4 datasets. For a smaller number of levels, indicator and integer encoding yielded similar results. As those studies only consider gradient boosting and a limited amount of datasets, it is unclear whether results generalize to other datasets and ML algorithms. Several other publications studied encoding text data based on similarity (Cerdeira and Varoquaux, 2020; Cerdeira et al., 2018) and employed indicator or target encoding as baselines. Another line of work studies variable encodings employed within specific ML models. Wright and König (2019) study treatments for categorical variables in random forests together with dummy and integer encoding (18 datasets, cardinality between 3 and 38), concluding that indicator and integer encoding perform subpar in comparison to methods that re-order levels according to the target variable. Prokhorenkova et al. (2018) compared their new CatBoost variant of target encoding to smoothed target encoding without CV, hold-out, and leave-one-out CV on 8 datasets. While the CatBoost method performed best, hold-out came second and target encoding without CV performed worst. An overview of encoding techniques tailored towards neural networks (e.g. the widely adopted *entity embeddings* by Guo and Berkhahn (2016)) is provided in Hancock and Khoshgoftaar (2020). They present a survey of indicator- and embedding-based methods but no benchmark study. In contrast, our work is the first to focus on high cardinality variables and techniques that are agnostic to the subsequent ML method. We study this problem on a larger variety of datasets and settings.

The main goal of our study is to assess the impact different categorical encoding techniques have on subsequent models’ predictive performance. As the optimal encoding might differ depending on the ML algorithm, we consider various state-of-the-art algorithms, *regularized linear models* (LASSO), *random forests* (RF), *gradient tree boosting* (GB), *k-nearest neighbors* (KNN), and *support vector machines* (SVM). To find default settings for high cardinality features, we analyze a variety of datasets with different characteristics including regression, binary classification, and multiclass classification problems. Because our methods vary in runtime and complexity, we are interested in whether more complex methods are to be preferred, or if simpler approaches suffice. We also study the relationship between a feature’s cardinality and the choice of encoding technique by varying the minimum number of levels above which features are transformed.

1.4 Contributions

We survey a broad set of categorical encoding techniques and conduct a comprehensive benchmark study with a focus on high cardinality features. We carefully design a benchmark scenario as well as a preprocessing scheme allowing us to study 7 different encoding techniques in conjunction with 5 commonly used ML algorithms across 24 diverse datasets, both from a classification and a regression regime. We give a detailed description of our study design to highlight important considerations for studying high cardinality features. Our results provide an overview of the performance of various approaches heavily used in the literature. After a discussion of results concerning predictive performance, we provide further analyses seeking to inform practitioners which methods to apply. This includes an important discussion of runtimes.

2 Encodings

Pseudocode is presented in the main text (for non-standard encodings), or the Supplementary Material (for standard methods). An important detail is how encoding techniques treat new levels during the prediction phase.

2.1 Integer Encoding

The simplest strategy for categorical features is integer encoding (also called ordinal encoding). Observed levels from the training set are mapped to the integers 1 to L . Although new levels could be mapped to $L + 1$ or 0, model predictions would be arbitrary as the integer order does not carry information. Thus, we encode new levels as missing values and use mode imputation to obtain the integer which matches the most frequent level in the training set. Integer encoding should only be an acceptable strategy for tree-based models, which can separate all original levels with repeated splits.

2.2 Frequency Encoding

Frequency encoding maps each level to its observed frequency in the training set ($\hat{x}_l = N_l$). This assumes a functional relationship between the frequency of a level and the target. It implicitly reduces the number of levels, and the subsequent model can best differentiate between levels with dissimilar frequencies. This approach is heavily used in natural language processing to encode token or n-gram counts. We encode new levels with a frequency of 1.

2.3 Indicator Encoding

We use indicator encoding as an umbrella term for two common strategies to encode categorical features with a small to moderate number of levels: one-hot

and dummy encoding. One-hot encoding transforms the original feature into L binary indicator columns, each representing one original level. An observation is coded with 1 for the indicator column representing its level ($x_i^{train} = l$) and 0 for all other indicators. Dummy encoding results in only $L - 1$ indicator columns. A reference feature level is chosen that is encoded with 0 in all indicator columns. For one-hot encoding, the zero vector can be used to encode new levels which were not observed during training. For dummy encoding, it is not useful to collapse new levels to the often arbitrary reference category; in our case the first level in alphabetical order. We replace new levels in the prediction phase with the most frequent level in the training set. As constructing all indicator variables is practically infeasible for high cardinality variables, we limit their number by collapsing rare levels beyond a varying threshold to a single *other* category before encoding.

2.4 Hash Encoding

Hash Encoding can be used to compute indicator variables based on a hash function (Weinberger et al., 2009). The basic idea is to transform each feature level l into an integer $hash(l) \in \mathbb{N}$, based on its label. This integer is then transformed into an indicator representation, with 1 in indicator column number $(hash(l) \bmod hash.size) + 1$ and 0 in all remaining columns (Kuhn and Johnson, 2019). Some levels are hashed to the same indicator representation. The smaller the *hash.size*, the higher the number of collapsed levels. The number of indicators is often effectively lower than *hash.size*, as some indicators can be constant in the training set (we remove those columns for both training and prediction). Although we could jointly hash multiple features, we hash each feature separately to improve comparability with the other encoders.

2.5 Leaf Encoding

Leaf encoding fits a decision tree on the training set to predict the target based on the categorical feature. Each level is encoded by the number of the terminal node, in which an observation with the respective level ends up. In that way, leaf encoding combines feature levels with similar target values. We use the *rpart* package in R (Therneau and Atkinson, 2018) which grows CARTs with categorical feature support that can be pruned based on internal performance estimates from 10-fold CV. Thus, our leaf encoder automatically uses an “optimal” number of new levels. To speed up the computation for multi-class classification, our implementation uses the ordering approach presented in Wright and König (2019). New levels are encoded with the arbitrary number of the terminal node with most observations during training. Encoded values are treated as a new categorical feature and encoded by one-hot encoding. Our leaf encoder can be thought of as a simplification of the approach suggested by Grąbczewski and Jankowski (2003).

Algorithm 1 Leaf Encoding

Training: require number of cross-validation folds $K \in \mathbb{N}$
fit CART tree on \mathcal{D}^{train} with complexity pruning based on K -fold cross-validation
for all $x_i^{train} \in \mathcal{X}^{train}$ **do** $\hat{x}_i = t$ with x_i^{train} in terminal node t
Prediction:
for x^{new} **do**
 if $x^{new} \in \mathcal{L}^{train}$ **then** $\hat{x}^{new} = t$ with x^{new} in terminal node t
 else $\hat{x}^{new} = b$ where b indicates the biggest terminal node

2.6 Impact Encoding

An early formal description of a so-called impact, target or James-Stein encoder was provided by Micci-Barreca (2001). The basic idea is to encode each feature level with the conditional target mean (regression) or the conditional relative frequency of one or more target classes (classification). The impact encoder for classification uses a logit link and transforms the original feature into C numeric features, each representing one target class. A smoothing parameter ϵ is introduced to avoid division by zero. This parameter could be used to further regularize towards the unconditional mean. We choose a small $\epsilon = 0.0001$ as we want to compare simple target encoding with the regularized encoder introduced next. Weights of evidence encoding from the credit scoring classification literature (Hand and Henley, 1997) is almost identical to impact encoding, but without regularization or centering.

2.7 GLMM Encoding

Smoothed target encoding (Micci-Barreca, 2001) can be interpreted as a simple (generalized) linear mixed model (glmm) in which the target is predicted by a random intercept for each feature level in addition to a fixed global intercept. This connection is described in Kuhn and Johnson (2019). To achieve regularized impact encoding, we implemented glmm encoders for regression, binary, and multiclass classification which are described in Algorithm 2 (regression) and in the Supplementary Material (classification). The encoded value for each level is based on the spherical conditional mode estimates. In regression, the conditional modes are similar to the mean target value for each level, weighted by the relative observed frequency of that level in the training set (Gelman and Hill, 2006, p. @bates_2018). The estimate of the fixed intercept can be used during the prediction phase to encode new feature levels not observed in the training set. In multiclass classification, we fit C one vs. rest glmms resulting in one encoded feature per class. An important advantage of using a glmm over impact encoding with a smoothing parameter is that a reasonable amount of regularization is determined automatically and tuning the complete ML pipeline is not necessary.

Overfitting can be further avoided through combination with cross-validation (CV) to train the encoder on independent observations without limiting the

data to train the ML model. We provide an implementation which combines target encoding based on glmms with CV. During the training phase, we partition the data using CV into $n.folds$ and fit a glmm on each resulting training set. For each observation, there is exactly one glmm that did not use that observation for model fitting and can be safely used for encoding. Note that the $n.folds$ CV models (for $n.folds > 1$) are only used during the training phase. In the prediction phase, feature values are always encoded by a single glmm fitted to the complete training set. We study this method in three different settings: without CV (noCV), with 5– (5CV) and with 10– fold CV (10CV). In our study, we use the lmer (regression) and glmer (classification) functions from the lme4 package in R (Bates et al., 2015) as an efficient way to fit glmms.

Algorithm 2 GLMM Encoding Regression

Training: require $n.folds \in \mathbb{N}$
fit simple random intercept model: $y_i^{train} = \beta_{0l} + \epsilon_i = \gamma_{00} + u_l + \epsilon_i$ on \mathcal{D}^{train}
with $u_l \stackrel{iid}{\sim} N(0, \tau^2)$, $\epsilon_i \stackrel{iid}{\sim} N(0, \sigma^2)$ and $x_i^{train} = l$, $l \in \mathcal{L}^{train}$
if $n.folds = 1$ **then**
 for all $x_i^{train} \in \mathcal{X}^{train}$ **do**
 $\hat{x}_i^{train} = \hat{\beta}_{0l}^{\mathcal{D}^{train}}$ with $x_i^{train} = l$
else use $n.folds$ cross-validation scheme to make training sets $\mathcal{D}_1^{train}, \dots, \mathcal{D}_{n.folds}^{train}$
 and fit simple random intercept model on each \mathcal{D}_m^{train}
 for all $x_i^{train} \in \mathcal{X}^{train}$ **do**
 $\hat{x}_i^{train} = \hat{\beta}_{0l}^{\mathcal{D}_m^{train}}$ with $x_i^{train} = l$ based on the model m
 with $(x_i^{train}, y_i^{train}) \notin \mathcal{D}_m^{train}$
Prediction:
for x^{new} **do**
 if $x^{new} \in \mathcal{L}^{train}$ **then**
 $\hat{x}^{new} = \hat{\beta}_{0l}^{\mathcal{D}^{train}}$ with $x^{new} = l$ based on full model fitted on \mathcal{D}^{train}
 else $\hat{x}^{new} = \hat{\gamma}_{00}$ based on full model fitted on \mathcal{D}^{train}

2.8 Control Conditions

We include three control conditions to better understand the effectiveness of the investigated encoders: The performance of a featureless learner (**FL** condition) was estimated as a conservative baseline for each dataset. In regression problems, FL predicts the mean of the target variable in the training set for each observation in the test set. In classification problems, the most frequent class of the target within the training set is predicted. For each dataset, we also consider a RF without encoding (**none** condition), to compare the use of encoding methods with a natural categorical splitting approach. The ranger (Wright and Ziegler, 2017) implementation provides efficient categorical feature support by ordering levels once before starting the tree growing algorithm (Wright and König, 2019). In the **remove** high cardinality features control

condition, we omit features with a high number of levels above some threshold and use one-hot encoding (without collapsing rare levels) for the remaining features. This condition reflects on whether including high cardinality features does indeed improve predictive performance. Otherwise, the best encoding might just provide the least impairment compared to not including any high cardinality features. We include an overview of available implementations in widely used ML frameworks for R and python in the Supplementary Material.

3 Benchmark Setup

3.1 Datasets

Table 1 summarizes datasets used in our benchmark. We specifically investigate datasets that contain categorical variables with a large number of levels, including many well-known datasets from previous studies (Cerdeira et al., 2018; Coors, 2018; Kuhn and Johnson, 2019; Prokhorenkova et al., 2018). All datasets can be downloaded from the `OpenML` platform (Vanschoren et al., 2013) based on the name or the displayed `OmlId`. The datasets include 8 regression, 10 binary classification, and 6 multiclass classification problems (between 3 and 12 classes). To assess the imbalance in categorical variables we computed the normalized entropy for each categorical variable. The maximum normalized entropy is 1, which corresponds to a uniform distribution, while a lower number indicates a larger imbalance. Sample sizes range between 736 and 1224158 observations. The total number of features ranges between 5 and 208. Datasets contain between 1 and 20 categorical features with more than 10 levels, with the maximum number of levels for a feature ranging between 14 and 30114. Missing values are present in about half of the datasets.

3.2 High Cardinality Threshold

It is often assumed that advanced encoding methods are only advantageous for variables with a high number of levels, while simple indicator encoding is more appropriate for a small number of levels. To reflect this in our benchmark, a high cardinality threshold (HCT) parameter with values of 10, 25, and 125 was introduced determining varying configurations for the different encoders. For indicator encoding, the $HCT - 1$ most frequent levels are encoded together with a single collapsed category for the remaining levels. For integer, frequency, hash, leaf, impact, and glmm encoders, only features with more than HCT levels in the training set were encoded with the respective strategy, while the remaining categorical features were one-hot encoded. HCT is used as the hash size in hash encoding. In the remove control condition, features with more than HCT levels in the training set are removed from the feature set.²

² Based on the number of categorical features and levels per feature, some HCT settings were removed from the benchmark for some combinations of dataset \times encoder. This ensured

3.3 Machine Learning Pipeline and Algorithms

In total, we investigate 5 ML algorithms. We kept tuning their hyperparameters to a minimum because we were interested in the effect of the encoding techniques instead of a comparison of ML algorithms. LASSOs were fitted with `glmnet` (Friedman et al., 2010), internally tuning the regularization using 5-fold CV. RFs with 500 trees were trained using `ranger` (Wright and Ziegler, 2017) without tuning, because RFs can be expected to give reasonable results with default settings (Probst et al., 2019). GB models were trained using `xgboost` (Chen et al., 2018), setting the learning rate to 0.01 and determining the number of iterations using early stopping on a 20% holdout set. KNN was taken from package `kkn` (Schliep and Hechenbichler, 2016) standardizing features and using a constant $k = 15$ for the number of nearest neighbors, together with an information gain filter (Brown et al., 2012) to limit the number of features to 25. SVMs with radial basis function kernel were trained with `liquidSVM` (Steinwart and Thomann, 2017). The bandwidth and regularization parameters were internally tuned using 5-fold CV. We used a one-vs-all approach for multiclass-settings.

The ML pipeline outlined below was used for all experimental conditions. It was carefully designed to ensure consistent results for extreme conditions (e.g. some levels only existing in the training data).

- **Imputation I:** Create a new factor level for missing values in categorical features with more than two categories. Impute missing values in binary features using the mode and missing values in numerical features using the mean feature value in the training data.
- **Encoding:** Transform the complete categorical features by the respective encoder. In the *no encoding* condition, the encoder simply passes on its input. The leaf and remove conditions still return categorical features, while the remaining encoders return only numerical features. Encoders only affect categorical variables above the specified HCT value.
- **Imputation II:** To handle new levels observed during prediction, impute missing values obtained during encoding.
- **Drop constants:** Drop features that are constant during training. As none of the original datasets includes constant columns, this step only removes constant features that are produced by the encoders or the CV splitting procedure.
- **Final one-hot encoding:** Transform all remaining categorical features via one-hot encoding (skipped for no encoding condition).
- **Learner:** Use the transformed data from each training set to fit the respective ML algorithm. In the prediction phase, transformed feature values (based on the trained encoder) for new observations in each test set are fed into the trained model to compute predictions.

that encoders always affect at least one feature and the remove condition always removes at least one feature. If settings where an encoder would lead to identical encoding strategies for all features of a dataset, we only kept the condition with the smallest HCT value.

Table 1 Benchmark Datasets and Dataset Characteristics

| OmlId | Name | Cl | N | NA% | Num | Bin | Cat | HighCardLevels | Entropy |
|-------|-------------------------------|----|---------|-------|-----|-----|-----|---------------------------------------|---------|
| 41211 | ames-housing | 0 | 2930 | 0.00 | 34 | 2 | 44 | 10,10,16,16,17,28 | |
| 41445 | employee_salaries | 0 | 9228 | 0.03 | 1 | 2 | 3 | 37,385,694 | |
| 41210 | avocado-sales | 0 | 18249 | 0.00 | 8 | 1 | 2 | 54 | |
| 41437 | wine-reviews | 0 | 129971 | 6.22 | 1 | 0 | 6 | 19,43,425,707,1229,16757 | |
| 41444 | medical_charges | 0 | 163065 | 0.00 | 1 | 0 | 5 | 51,100,306,1977,3201 | |
| 41267 | particulate-matter-ukair-2017 | 0 | 394299 | 0.00 | 4 | 0 | 5 | 12,30,53 | |
| 41251 | flight-delay-usa-dec-2017 | 0 | 457892 | 0.00 | 2 | 0 | 7 | 12,31,52,52,294,294 | |
| 41255 | nyc-taxi-green-dec-2016 | 0 | 1224158 | 0.00 | 4 | 3 | 7 | 241,260 | |
| 41283 | churn | 2 | 5000 | 0.00 | 14 | 2 | 3 | 10,51 | |
| 981 | kdd_internet_usage | 2 | 10108 | 0.39 | 0 | 48 | 20 | 10,11,46,77,119,129 | |
| 4135 | Amazon_employee_access | 2 | 32769 | 0.00 | 0 | 0 | 9 | 67,128,177,343,343,449,2358,4243,7518 | |
| 41434 | Click_prediction_small | 2 | 39948 | 0.00 | 3 | 0 | 6 | 6064,19228,19803,22381,25321,30114 | |
| 1590 | adult | 2 | 48842 | 0.95 | 6 | 1 | 7 | 14,16,41 | |
| 1114 | KDDCup09_upselling | 2 | 50000 | 68.50 | 174 | 4 | 30 | 14,...,5073,5713,13990,15415,15415 | |
| 41162 | kick | 2 | 72983 | 6.39 | 14 | 3 | 15 | 12,16,33,37,74,134,153,863,1063 | |
| 41442 | open_payments | 2 | 73354 | 22.61 | 0 | 1 | 4 | 513,1460,2255,4365 | |
| 41447 | road-safety-drivers-sex | 2 | 233964 | 10.78 | 2 | 1 | 3 | 380,20397 | |
| 41224 | porto-seguro | 2 | 595212 | 2.45 | 27 | 23 | 8 | 10,12,18,104 | |
| 188 | eucalyptus | 5 | 736 | 3.20 | 14 | 0 | 5 | 12,14,16,27 | |
| 41446 | Midwest_survey | 10 | 2778 | 1.66 | 0 | 21 | 5 | 1008 | |
| 41212 | hpc-job-scheduling | 4 | 4331 | 0.00 | 5 | 0 | 2 | 14 | |
| 41216 | video-game-sales | 12 | 16598 | 0.25 | 6 | 0 | 2 | 31,578 | |
| 41440 | okcupid-stem | 3 | 50789 | 15.97 | 2 | 1 | 16 | 12,12,15,15,18,32,45,48,184,208,7019 | |
| 4541 | Diabetes130US | 3 | 101766 | 0.00 | 13 | 9 | 25 | 10,10,18,73,717,749,790 | |

Note:

OmlId = Id on OpenML, Name = name on OpenML, Cl = classes (0: regression), N = observations, NA% = percentage of missing values, Num = numeric features, Bin = binary features, Cat = categorical features, HighCardLevels = number of levels for each categorical feature with at least 10 levels (some levels are not displayed for KDDCup09_upselling), Entropy = box-plot of normalized Shannon-entropy across levels (smaller = larger imbalance).

3.4 Performance Evaluation

We perform all analyses in the open-source statistical software R (R Core Team, 2021). To enable a fair and reliable comparison in our study, we implemented all encoding methods on top of the `mlrCP0` package (Binder, 2018). The pre-processing, as well as the final ML algorithm described in 3.3, were trained and resampled using the `mlr` framework (Bischl et al., 2016) together with the `batchtools` package (Lang et al., 2017) to scale the benchmark analysis to HPC compute infrastructure. All materials for this study, including reproducible code for this manuscript and result objects can be downloaded from our **online repository**³.

Throughout our experiments, we use 5-fold CV to obtain estimates of predictive performance. Depending on the target variable, we report root mean squared error (RMSE) for regression, area under the curve (AUC) for binary classification, and its extension AUNU (Ferri et al., 2009) for multiclass problems:

$$RMSE = \sqrt{\frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2}, \quad AUC = \frac{U}{N_{pos} \cdot N_{neg}}, \quad AUNU = \frac{\sum_{c=1}^C AUC_c}{C}$$

N : number of observations, N_{pos} and N_{neg} : number of true labels in both classes, U : test statistic of the Mann-Whitney-U test, and AUC_c : one vs. rest AUC of class c .

In our benchmark study, each encoding method listed in Section 2 is combined with all ML algorithms (c.f. Section 3.3) across high cardinality thresholds (HCT) 10, 25, 125. While we only report results for the best HCT setting for each ML algorithm \times dataset combination in Section 4, we study the effect of the HCT parameter in more detail in section 4.5.

4 Benchmark Results

Our main question is which encoding methods generally work well across various datasets. We report results for regression and classification datasets separately, as the associated metrics differ in scale.

For 6 datasets, some conditions with the SVM led to unexpected crashes due to memory problems or numerical errors. We completely removed those datasets for the SVM when computing ranks or other statistics that compare encodings across datasets.

4.1 Encoder Performance on Each Dataset

Mean performance estimates along with minimum and maximum performance are reported for all datasets in Figures 2 (regression), 3 (binary classification),

³ https://github.com/compstat-lmu/paper_2021_categorical_feature_encodings

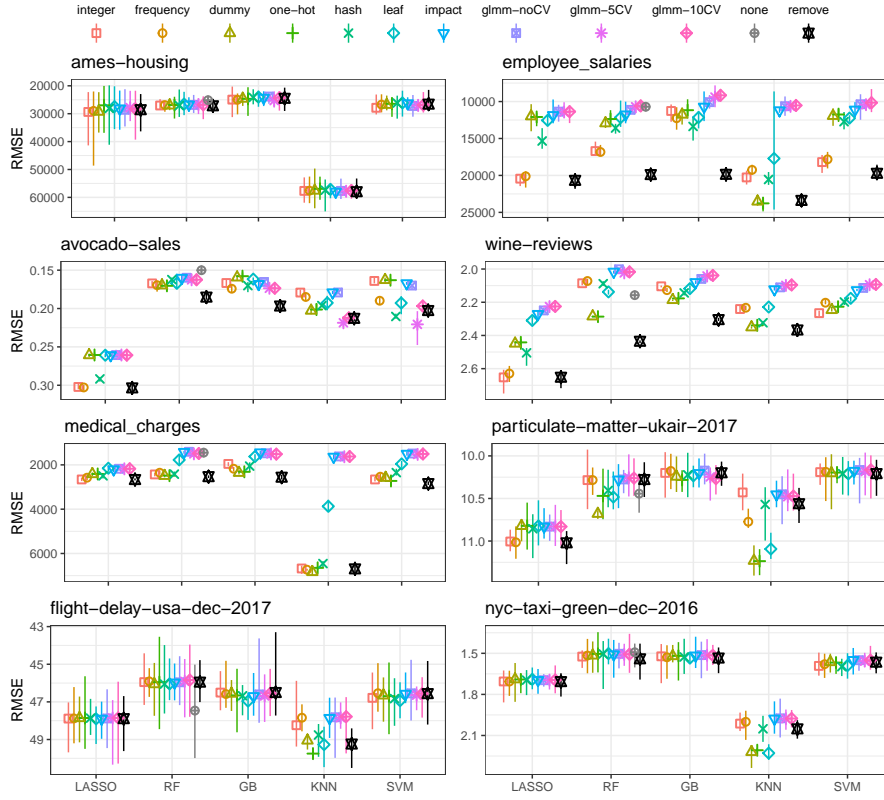


Fig. 2 Performance estimates from 5-CV for regression (mean, min, max). For each combination, only the best HCT condition is displayed. Note the reversed y-axis to ease visual interpretation.

and 4 (multiclass classification). To reduce the complexity induced by the hyperparameter HCT we only display the parameter condition with the best performance for each combination of dataset \times encoding \times ML algorithm. The y-axis differs for all datasets and is reversed for the RMSE for better visual comparison. For some datasets, the remove condition performed very similar to the other encodings (e.g. *ames-housing*, *porto-seguro*), suggesting that categorical features were not informative. Performance in some CV folds was below the FL learner for *flight-delay-usa-dec-2017* ($RMSE_{FL} = 48.8052624$) and *nyc-taxi-green-dec-2016* ($RMSE_{FL} = 2.2152442$).

On datasets with substantive performance differences, target encoding with the glmm encoder was generally **most** effective. The **worst** encoder differed by ML algorithm and dataset. For datasets *Click_prediction_small*, *KDD-Cup09_upselling*, *kick*, and *okcupid-stem* some encoder \times ML algorithm conditions performed worse than simply removing high cardinality features.

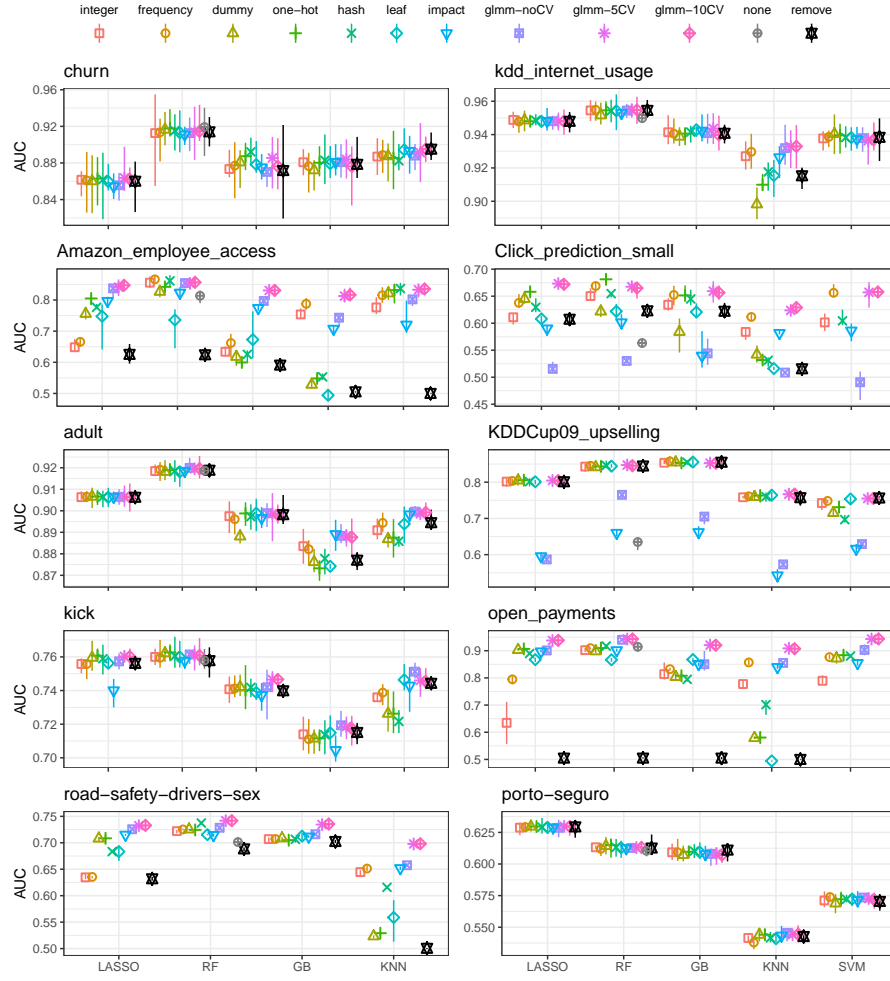


Fig. 3 Performance estimates from 5-CV for binary classification (mean, min and max). For each combination, only the best HCT condition is displayed.

4.2 Meta Rankings and Dataset Clustering

To further summarize, present meta rankings for each ML algorithm in Figure 5. First, we defined an encoder relation within each dataset, based on corrected resample t-tests (Nadeau and Bengio, 2003). An encoding was defined to beat another encoding if the one-sided p-value of the t-test was $< .05$. This allowed us to compute a weak-order consensus ranking R defined by the optimization problem:

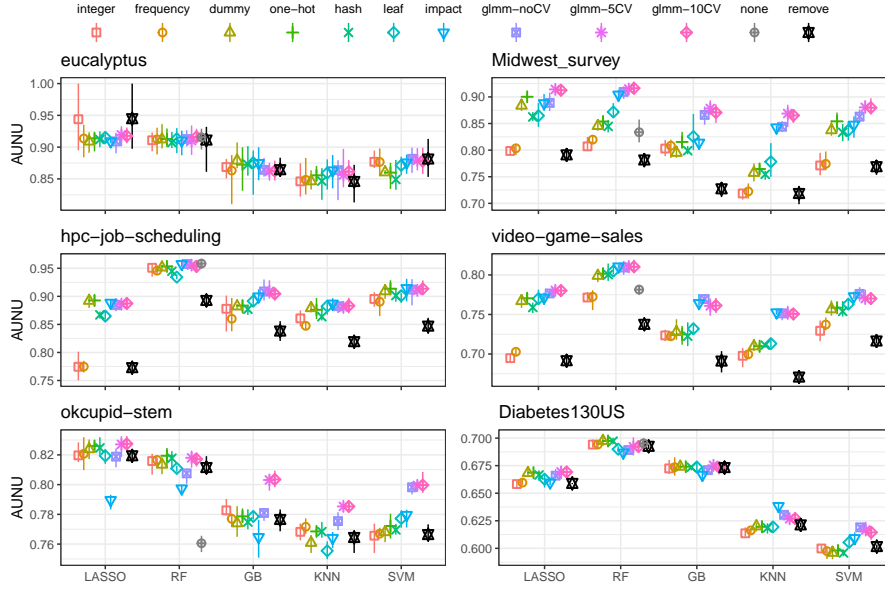


Fig. 4 Performance estimates from 5-CV for multiclass classification (mean, min and max). For each combination, only the best HCT condition is displayed.

$$\arg \min_{R \in \mathcal{C}} \sum_{b=1}^B d(R_b, R)$$

where d is the symmetric difference distance and R_b is the relation for dataset b (Hornik and Meyer, 2007; Meyer and Hornik, 2018). The symmetric difference between two relations is the number of cases one encoding beats another encoding in one relation but not in the other one.

Although the presented solutions of the optimization problem are not unique, rankings were highly stable for the high and low ranks. Meta rankings seem to be highly consistent with the individual patterns of encoder performances reflected in Figures 2 to 4. Looking at meta-rankings, approaches based on GLMM’s in combination with cross-validation outperform all other approaches across all algorithms. A further interesting detail omitted in Figure 5 for clarity is that the *none* encoding strategy for the RF was beaten by all strategies except for the *remove* condition. This implies, that even if the algorithm provides a mechanic for treating categorical variables, it might often be optimal to use a different strategy instead.

In an exploratory analysis, we tried to find clusters of datasets based on systematic patterns of encoder performance (independent of the employed ML algorithm). We computed a partial-order consensus relation across ML algorithms for each dataset and hierarchically clustered the dataset consensus

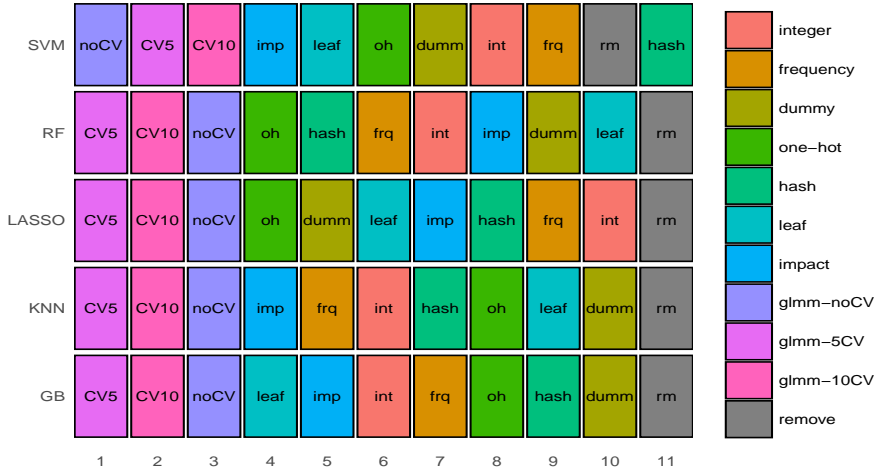


Fig. 5 Consensus rankings across all datasets for each algorithm. Lower ranks indicate better performance. The rank of the *none* control condition of the RF (rank 11) was omitted from the figure.

relations using the symmetric difference distance in combination with the complete linkage agglomeration method. The resulting dendrogram is displayed in Figure 6. Although the cluster structure is somewhat ambiguous, roughly three clusters can be described: The first 9 datasets from the top of the dendrogram are characterized by a medium to high number of levels, low performance of the remove condition (indicating the importance of high cardinality features) and clear performance advantage of target encoders. For the next 13 datasets which contain the smallest number of levels, traditional encodings can compete with target encoding. This biggest cluster also includes 9 datasets with zero distances, in which no significant performance differences could be observed between any encoding conditions (nor with the remove condition, indicating that high cardinality features are less informative). The last two datasets with the highest number of levels formed a separate cluster, in which target encoding without strong regularization (impact, glmm-noCV) showed severe overfitting. Note that clusters were not determined by problem type, again suggesting that encoder rankings are somewhat similar for regression and classification settings.⁴

⁴ We tried to corroborate this observation by comparing meta-rankings between problem types. Unfortunately, problem type specific consensus relations did not converge (probably due to the reduced number of datasets) and were uninterpretable.

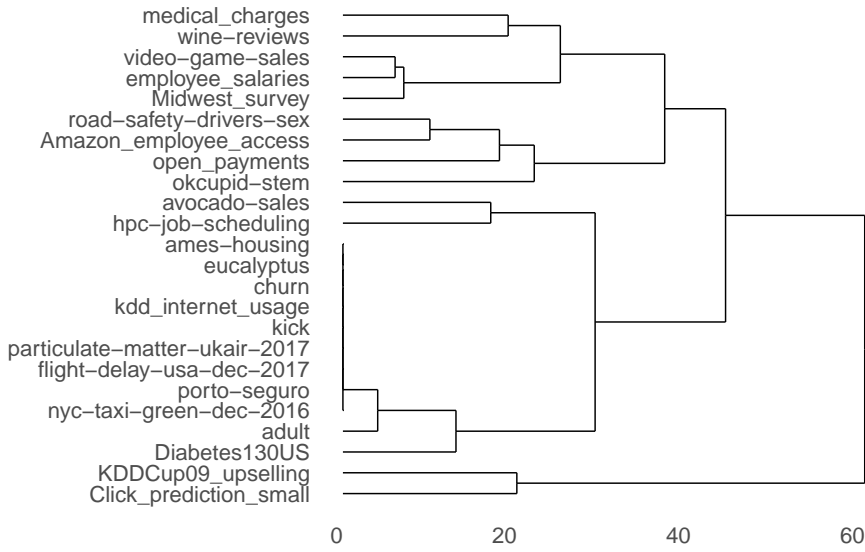


Fig. 6 Hierarchical cluster analysis of benchmark datasets. The symmetric difference distance between two datasets reflects differences in performance patterns between encodings (independent of the employed ML algorithm).

4.3 Summary of Encoder Performance

Regularized target encoding was superior or at least competitive on all datasets. We could not observe a setting in which regularized target encoding was convincingly beaten by target agnostic methods. Especially effective was glmm encoding with 5-fold-CV, which ranked first place for all ML algorithms except SVM. Performance often did not improve with glmm-10CV, suggesting that 5 folds might be a good regularization default in practice. In line with earlier research (Micci-Barreca, 2001; Prokhorenkova et al., 2018), target encoding with regularization (glmm) performed better or equally well in comparison with the unregularized impact encoder. When glmm performed better, impact encoding not only performed worse than glmm with CV but was sometimes also inferior compared to other encoders.

The following observations were also interesting: Surprisingly, integer encoding did not perform well with GB and target-based encoders (especially the glmm encoder) seem to be preferable. For LASSO, previous studies have suggested that indicator encoding works well, even with a very high number of levels (Cerdeira et al., 2018). Although the glmm encoders ranked first in our benchmark for the LASSO, it was the only algorithm where the indicator encoders achieved the next best ranking. Note that for computational reasons we limited the maximum amount of indicator variables per feature to 125 in our experimental design. The $HCT = 125$ setting performed best for LASSO

with indicator encoding in a large number of datasets, indicating that performance might have further improved with higher values. Both KNN and SVM rely on numerical distances in feature space and tend to perform poorly in the presence of high dimensionality. Thus, we expected that target encoding should work well here as it transforms categories into a single, smooth numerical feature. This was backed by our benchmark results. For some datasets (*medical_charges*, *road-safety-drivers-sex*, and *Midwest_survey*) KNN (without tuning of the optimal number of nearest neighbors) could compete with the more sophisticated ML algorithms when combined with target encoding, but performed poorly with other encoders. Although consistent with the big picture, SVM results have to be considered with care, as some experimental conditions resulted in unexpected computational errors. In RF, a widely used strategy to deal with categorical features is to order levels by average target statistics for a given level. For a small number of levels, this approach has been reported superior to indicator and integer encoding (Wright and König, 2019), while we observed poor performance for datasets with a larger number of levels. When looking for a simple default encoding, indicator encoding in combination with collapsing small levels seems a robust alternative, although the glmm encoders performed better. We found that one-hot encoding usually gave a slightly better performance than dummy encoding, which has also been observed by Chiquet et al. (2016) and Tutz and Gertheiss (2016) (p.254). Our results suggest that one-hot encoding is the better standard compared to dummy encoding when applying nonlinear regularized models like RF, GB or SVM with (high cardinal) categorical features. We provide a more detailed comparison for indicator encoding in the Supplementary Material.

4.4 Runtime Analysis

To determine whether traditional encodings are preferred when facing limited computational resources, we further analyzed runtimes of the whole analysis pipeline for different encoders and ML algorithms. Aggregated results are shown in Table 2. To enable a meaningful comparison, we report runtime as the fraction of a full pipeline’s strategy compared to the one-hot encoding condition and then further aggregate across datasets using the median. Again, we only report the **best** HCT setting. Absolute runtimes are hard to interpret and aggregate because runtime distributions across datasets are heavily skewed as proportionally larger runtimes are observed for big datasets. While a pipeline’s runtime can be dominated by the encoder for small datasets, the training of the consecutive ML algorithm is dominating for large datasets, which might render differences during encoding irrelevant. Therefore we aim to report what is important in practice, the time differences for training the **full pipeline**. The results clearly show that regularized target encoding does not consistently yield slower runtimes compared to simple strategies like indicator encoding. Supposedly, the more efficient representations produced by target encoding lead to faster runtimes of subsequent ML algorithms. This

Table 2 Proportional Increase in Runtime Compared to One-Hot Encoding

| Encoding | LASSO | RF | GB | KNN | SVM |
|-----------|------------------------------------|-------------------------------------|-------------------------------------|-------------------------------------|------------------------------------|
| integer | 0.41 ^{1.4} ₀ | 0.59 ^{1.1} ₀ | 0.67 ^{0.9} ₀ | 0.85 ^{1.5} _{0.4} | 1.01 ^{1.6} _{0.6} |
| frequency | 0.34 ^{2.1} ₀ | 0.54 ^{1.7} ₀ | 0.54 ^{1.5} ₀ | 0.79 ^{1.1} _{0.2} | 1.11 ^{1.7} _{0.5} |
| dummy | 1.13 ^{2.8} _{0.8} | 0.91 ^{1.7} _{0.5} | 0.97 ^{5.7} _{0.4} | 1.05 ^{1.6} _{0.6} | 1.02 ^{1.5} _{0.7} |
| one-hot | 1 ¹ ₁ | 1 ¹ ₁ | 1 ¹ ₁ | 1 ¹ ₁ | 1 ¹ ₁ |
| hash | 0.87 ^{2.6} _{0.2} | 1.03 ^{3.3} _{0.4} | 0.98 ^{7.3} _{0.4} | 0.93 ^{2.2} _{0.5} | 1.16 ² _{0.5} |
| leaf | 0.46 ^{1.7} ₀ | 0.67 ^{2.3} ₀ | 0.85 ^{9.1} ₀ | 0.97 ^{2.2} _{0.5} | 1.01 ^{1.9} _{0.6} |
| impact | 0.51 ^{1.9} _{0.1} | 0.61 ^{1.5} _{0.1} | 0.82 ^{120.5} ₀ | 1.11 ²⁴ _{0.3} | 1.06 ^{1.8} _{0.7} |
| glmm-noCV | 0.58 ^{2.8} ₀ | 0.57 ^{3.8} _{0.1} | 1.65 ¹² ₀ | 1.09 ^{10.3} _{0.3} | 1.12 ¹ _{0.7} |
| glmm-5CV | 0.83 ^{2.2} _{0.1} | 0.66 ^{15.3} _{0.1} | 7.21 ^{69.3} ₀ | 2.77 ^{50.1} _{0.7} | 1.41 ^{4.6} _{0.6} |
| glmm-10CV | 1.07 ^{4.2} _{0.1} | 0.96 ^{28.4} _{0.1} | 12.71 ^{127.6} ₀ | 5.37 ^{97.9} _{0.8} | 1.55 ^{4.6} _{0.5} |
| none | | 0.21 ^{1.4} ₀ | | | |
| remove | 0.41 ^{1.7} ₀ | 0.48 ^{1.1} ₀ | 0.58 ^{1.3} ₀ | 0.82 ^{1.2} _{0.2} | 1.09 ² _{0.6} |

Note:

Median $_{min}^{max}$ across datasets of the proportional increase in runtime from 5-CV, when comparing the respective encoder with one-hot encoding. Only the best HCT conditions are reported.

suggests that a possible runtime vs. predictive performance trade-off might also be in favour of target encoding, especially for large datasets where a high number of indicator variables increases computational load. We did observe a substantial increase in runtimes when using regularized glmm with GB (with little tuning), where runtimes are relatively short in comparison to the time required to fit categorical encoders. In other settings (LASSO, RF), the glmm encoders have been observed to be even faster than indicator encoding.

4.5 Analysing High Cardinality Thresholds

Until now, we ignored the HCT parameter by reporting only the condition with the best performance. An interesting question is whether target encoding is only useful for features with a very high number of levels or also for fewer levels, where most practitioners would routinely use indicator encoding. We tested this using HCT thresholds of 10, 25 and 125, but the results are not easy to interpret. In general, the optimal threshold seemed to strongly depend on the dataset, but we also observed some weak patterns: LASSO improved for larger HCT, indicating that its internal regularization can efficiently deal with the sparseness induced by indicator encoding. In comparison, other methods generally yielded better performance if features above the very low HCT of 10 were encoded using one of the target encoding strategies. Differences in regularization for the target based encoders seemed not to prefer different HCT values.

5 Discussion

In our benchmark, we compared encoding strategies for high cardinality categorical features on a variety of regression, binary, and multiclass classification datasets with different ML algorithms. Regularized target encoding was superior across most datasets and ML algorithms. Although the performance of other encoding strategies was comparable in some conditions, target encoding was never outperformed. In general, our results suggest that regularized target encoding based on glmms with 5-fold CV (glmm-5CV) works well for all kinds of algorithms and should be a reasonable default strategy. It sometimes leads to slightly longer runtimes (in comparison to indicator encoding), but especially for larger datasets this is often offset by the more efficient representation produced by target encoding. The glmm encoder has a clear advantage over target encoding with a smoothing hyperparameter (Micci-Barreca, 2001), as costly tuning the whole ML pipeline with different smoothing values is not required.

What constitutes a “high” cardinality problem is a difficult question that might not only depend on the number of levels but also on other characteristics of the dataset and its features. Supposedly, the number of datasets in our study was too small to discover consistent patterns between encoder performance and dataset characteristics. Target encoding features with only 10 levels seemed to be effective in a substantive number of conditions, but for other datasets, higher HCT values performed better. Thus, some form of hyperparameter tuning seems necessary to decide the level threshold for target encoding at this point, as no suitable defaults seem to be available. Note that we compared different HCT values but then used the same encoding strategy for all affected features alongside one-hot encoding for the remaining ones. A further improvement could be to decide whether to use target encoding on a feature by feature basis. ML pipelines could introduce a categorical hyperparameter for each feature that represents which encoding is used. To make this complicated meta optimization problem feasible, only a small number of encoders can be included. Our study can help to decide which encoders could be safely omitted from consideration.

5.1 Limitations

Several design decisions were necessary to make answering our research questions computationally feasible: We used minimal tuning for our ML algorithms, as we were not interested in comparing their performance against each other. This probably led to suboptimal performance for the GB, KNN, and SVM learners. When interpreting our results, we assume that the encoder rankings are comparable when more extensive tuning is used. This is plausible because the meta rankings between algorithms were highly stable. We only used 5-fold-CV without repetitions to estimate predictive performance, but due to a small variance between folds, we could confidently detect performance dif-

ferences of encoders for many datasets. Because we are interested in model-agnostic methods that can be combined with any supervised ML algorithm, we did not investigate recent model-specific strategies (Guo and Berkhahn, 2016; Prokhorenkova et al., 2018). We do not differentiate between nominal and ordinal categorical features: Most publicly available datasets simply do not contain any high cardinal ordinal features. For many applications, ordinal feature information is not available as metadata which makes it less relevant for some applications like automatic ML. Thus, our benchmark does not include specific ordinal encoding methods like Helmert and Polynomial contrasts (Chambers and Hastie, 1992) or ordinal penalization methods (Tutz and Gertheiss, 2016). We only investigate traditional categorical variables and do not extend our analysis to multi-categorical variables or text-strings (Cerdeira and Varoquaux, 2020; Cerdeira et al., 2018), where other encoding techniques can harvest additional information. We only investigate univariate encodings, i.e. we always encode each variable separately. Levels with comparable main effects can not be distinguished based on the transformed feature, which prevents the consecutive ML algorithm from learning interactions for specific levels. Methods that jointly encode features and therefore leverage correlational structures between features are an interesting avenue for future research: Recent approaches use target-based encoding strategies to learn a vector-valued representation of a level, similar to entity embeddings (Guo and Berkhahn, 2016). Those should be compared to classical approaches from the optimal scaling literature (e.g. De Leeuw et al., 1976) or the subsequent *aspect* framework (Mair and Leeuw, 2010), which are unfortunately unfamiliar to many ML researchers.

5.2 Conclusion

This benchmark study compared the predictive performance of a variety of strategies to encode categorical features with a high number of unordered levels for different supervised ML algorithms. Regularized versions of target encoding, which uses predictions of the target variable as numeric feature values, performed better than traditional strategies like integer or indicator encoding. Most effective, with a consistently superior performance across ML algorithms and datasets, was a target encoder which combines simple generalized linear mixed models with cross-validation and does not require hyperparameter tuning.

5.3 Declarations

5.3.1 Funding

This work has been funded by the German Federal Ministry of Education and Research (BMBF) under Grant No. 01IS18036A and by the Bavarian Ministry of Economic Affairs, Regional Development and Energy through the Center

for Analytics – Data – Applications (ADA-Center) within the framework of „BAYERN DIGITAL II“ (20-3410-2-9-8).

5.3.2 Conflicts of Interest

The authors declare no conflict of interest.

5.3.3 Code Availability

All analysis code is publicly available at https://github.com/compstat-lmu/paper_2021_categorical_feature_encodings.

5.3.4 Availability of Data

All benchmark datasets are publicly available at <https://www.openml.org/>.

References

- Bates, D., 2020. Computational methods for mixed models. Vignette for lme4.
- Bates, D., Mächler, M., Bolker, B., Walker, S., 2015. Fitting linear mixed-effects models using lme4. *Journal of Statistical Software* 67, 1–48. doi:10.18637/jss.v067.i01
- Binder, M., 2018. MlrCPO: Composable preprocessing operators and pipelines for machine learning.
- Bischl, B., Lang, M., Kotthoff, L., Schiffner, J., Richter, J., Studerus, E., Casalicchio, G., Jones, Z.M., 2016. mlr: Machine learning in r. *Journal of Machine Learning Research* 17, 1–5.
- Bommert, A., Sun, X., Bischl, B., Rahnenführer, J., Lang, M., 2020. Benchmark for filter methods for feature selection in high-dimensional classification data. *Computational Statistics & Data Analysis* 143. doi:10.1016/j.csda.2019.106839
- Brown, G., Pocock, A., Ming-Jie, Z., Luján, M., 2012. Conditional likelihood maximisation: A unifying framework for information theoretic feature selection. *Journal of Machine Learning Research* 13, 27–66.
- Cerda, P., Varoquaux, G., 2020. Encoding high-cardinality string categorical variables. *IEEE Transactions on Knowledge and Data Engineering* 1–1. doi:10.1109/TKDE.2020.2992529
- Cerda, P., Varoquaux, G., Kégl, B., 2018. Similarity encoding for learning with dirty categorical variables. *Machine Learning* 107, 1477–1494. doi:10.1007/s10994-018-5724-2
- Chambers, J., Hastie, T., 1992. Statistical models. Chapter 2 of statistical models in s. Wadsworth & Brooks/Cole.

- Chen, T., He, T., Benesty, M., Khotilovich, V., Tang, Y., Cho, H., Chen, K., Mitchell, R., Cano, I., Zhou, T., Li, M., Xie, J., Lin, M., Geng, Y., Li, Y., 2018. Xgboost: Extreme gradient boosting.
- Chiquet, J., Grandvalet, Y., Rigai, G., 2016. On coding effects in regularized categorical regression. *Statistical Modelling* 16, 228–237. doi:10.1177/1471082X16644998
- Coors, S., 2018. Automatic gradient boosting (Master’s thesis). LMU Munich.
- De Leeuw, J., Young, F.W., Takane, Y., 1976. Additive structure in qualitative data: An alternating least squares method with optimal scaling features. *Psychometrika* 41, 471–503.
- Fernández-Delgado, M., Cernadas, E., Barro, S., Amorim, D., 2014. Do we need hundreds of classifiers to solve real world classification problems? *Journal of Machine Learning Research* 15, 3133–3181.
- Ferri, C., Hernández-Orallo, J., Modroiu, R., 2009. An experimental comparison of performance measures for classification. *Pattern Recognition Letters* 30, 27–38. doi:10.1016/j.patrec.2008.08.010
- Feurer, M., Klein, A., Eggenberger, K., Springenberg, J., Blum, M., Hutter, F., 2015. Efficient and robust automated machine learning, in: Cortes, C., Lawrence, N.D., Lee, D.D., Sugiyama, M., Garnett, R. (Eds.), *Advances in Neural Information Processing Systems* 28. Curran Associates, Inc., pp. 2962–2970.
- Friedman, J., Hastie, T., Tibshirani, R., 2010. Regularization paths for generalized linear models via coordinate descent. *Journal of Statistical Software* 33, 1–22. doi:10.18637/jss.v033.i01
- Gelman, A., Hill, J., 2006. *Data analysis using regression and multilevel/hierarchical models*. Cambridge university press.
- Grabczewski, K., Jankowski, N., 2003. Transformations of symbolic data for continuous data oriented models, in: Kaynak, O., Alpaydin, E., Oja, E., Xu, L. (Eds.), *Artificial Neural Networks and Neural Information Processing — Ican/Iconip 2003*. Springer Berlin Heidelberg, Berlin, Heidelberg, pp. 359–366.
- Guo, C., Berkahn, F., 2016. Entity embeddings of categorical variables. arXiv preprint arXiv:1604.06737.
- Hancock, J.T., Khoshgoftaar, T.M., 2020. Survey on categorical data for neural networks. *Journal of Big Data* 7, 1–41. doi:10.1186/s40537-020-00305-w
- Hand, D.J., Henley, W.E., 1997. Statistical classification methods in consumer credit scoring: A review. *Journal of the Royal Statistical Society: Series A (Statistics in Society)* 160, 523–541. doi:10.1111/j.1467-985X.1997.00078.x
- Hornik, K., Meyer, D., 2007. Deriving consensus rankings from benchmarking experiments, in: *Advances in Data Analysis*. Springer, pp. 163–170. doi:10.1007/978-3-540-70981-7_19
- Kuhn, M., Johnson, K., 2019. *Feature engineering and selection: A practical approach for predictive models*. Chapman; Hall/CRC.

- Lang, M., Bischl, B., Surmann, D., 2017. Batchtools: Tools for r to work on batch systems. *The Journal of Open Source Software* 2. doi:10.21105/joss.00135
- Mair, P., Leeuw, J. de, 2010. A general framework for multivariate analysis with optimal scaling: The r package aspect. *Journal of Statistical Software* 32, 1–23. doi:10.18637/jss.v032.i09
- Meyer, D., Hornik, K., 2018. Relations: Data structures and algorithms for relations.
- Micci-Barreca, D., 2001. A preprocessing scheme for high-cardinality categorical attributes in classification and prediction problems. *SIGKDD Explor. Newsl.* 3, 27–32. doi:10.1145/507533.507538
- Nadeau, C., Bengio, Y., 2003. Inference for the generalization error. *Machine Learning* 52, 239–281. doi:10.1023/A:1024068626366
- Probst, P., Wright, M.N., Boulesteix, A.-L., 2019. Hyperparameters and tuning strategies for random forest. *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery* 0, e1301. doi:10.1002/widm.1301
- Prokhorenkova, L., Gusev, G., Vorobev, A., Dorogush, A.V., Gulin, A., 2018. CatBoost: Unbiased boosting with categorical features, in: Bengio, S., Wallach, H., Larochelle, H., Grauman, K., Cesa-Bianchi, N., Garnett, R. (Eds.), *Advances in Neural Information Processing Systems* 31. Curran Associates, Inc., pp. 6638–6648.
- Prokopenko, V., 2018. Mean (likelihood) encodings: A comprehensive study. *Kaggle Forums*.
- R Core Team, 2021. R: A language and environment for statistical computing. R Foundation for Statistical Computing, Vienna, Austria.
- Schliep, K., Hechenbichler, K., 2016. Kknn: Weighted k-nearest neighbors.
- Steinwart, I., Thomann, P., 2017. liquidSVM: A fast and versatile svm package. *ArXiv e-prints* 1702.06899.
- Therneau, T., Atkinson, B., 2018. Rpart: Recursive partitioning and regression trees.
- Thomas, J., Coors, S., Bischl, B., 2018. Automatic gradient boosting.
- Thornton, C., Hutter, F., Hoos, H.H., Leyton-Brown, K., 2013. Auto-weka: Combined selection and hyperparameter optimization of classification algorithms, in: *Proceedings of the 19th Acm Sigkdd International Conference on Knowledge Discovery and Data Mining, KDD '13*. ACM, New York, NY, USA, pp. 847–855. doi:10.1145/2487575.2487629
- Tutz, G., Gertheiss, J., 2016. Rejoinder: Regularized regression for categorical data. *Statistical Modelling* 16, 249–260. doi:10.1177/1471082X16652780
- Vanschoren, J., N. van Rijn, J., Bischl, B., Torgo, L., 2013. OpenML: Networked science in machine learning. *SIGKDD Explorations* 15, 49–60. doi:10.1145/2641190.2641198
- Weinberger, K.Q., Dasgupta, A., Langford, J., Smola, A.J., Attenberg, J., 2009. Feature hashing for large scale multitask learning, in: *ICML*.
- Wright, M.N., König, I.R., 2019. Splitting on categorical predictors in random forests. *PeerJ* 7. doi:10.7717/peerj.6339

Wright, M.N., Ziegler, A., 2017. ranger: A fast implementation of random forests for high dimensional data in C++ and R. *Journal of Statistical Software* 77, 1–17. doi:10.18637/jss.v077.i01

References

- Dehghani et al(2021)Dehghani, Tay, Gritsenko, Zhao, Houlsby, Diaz, Metzler, and Vinyals. Dehghani M, Tay Y, Gritsenko AA, Zhao Z, Houlsby N, Diaz F, Metzler D, Vinyals O (2021) The benchmark lottery. 2107.07002
- Nießl et al(2021)Nießl, Herrmann, Wiedemann, Casalicchio, and Boulesteix. Nießl C, Herrmann M, Wiedemann C, Casalicchio G, Boulesteix AL (2021) Over-optimism in benchmark studies and the multiplicity of design and analysis options when interpreting their results. 2106.02447