

An Exponential Factorization Machine with Percentage Error Minimization to Retail Sales Forecasting

CHONGSHOU LI*, National University of Singapore, Singapore

BRENDA CHEANG, Red Jasper Holdings, Singapore

ZHIXING LUO, Nanjing University, China

ANDREW LIM, National University of Singapore, Singapore

This paper proposes a new approach to sales forecasting for new products (stock-keeping units, SKUs) with long lead time but short product life cycle. These SKUs are usually sold for one season only, without any replenishments. An exponential factorization machine (EFM) sales forecast model is developed to solve this problem which not only takes into account SKU attributes, but also pairwise interactions. The EFM model is significantly different from the original Factorization Machines (FM) from two-fold: (1) the attribute-level formulation for explanatory/input variables and (2) exponential formulation for the positive response/output/target variable. The attribute-level formation excludes infeasible intra-attribute interactions and results in more efficient feature engineering comparing with the conventional one-hot encoding, while the exponential formulation is demonstrated more effective than the log-transformation for the positive but not skewed distributed responses. In order to estimate the parameters, percentage error squares (PES) and error squares (ES) are minimized by a proposed adaptive batch gradient descent (ABGD) method over the training set. To overcome the over-fitting problem, a greedy forward stepwise feature selection (GFSFS) method is proposed to select the most useful attributes and interactions. Real-world data provided by a footwear retailer in Singapore is used for testing the proposed approach. The forecasting performance in terms of both mean absolute percentage error (MAPE) and mean absolute error (MAE) compares favorably with not only off-the-shelf models but also results reported by extant sales and demand forecasting studies. The effectiveness of the proposed approach is also demonstrated by two external public datasets. Moreover, we prove the theoretical relationships between PES and ES minimization, and present an important property of the PES minimization for regression models; that it trains models to underestimate data. This property fits the situation of sales forecasting where unit-holding cost is much greater than the unit-shortage cost (e.g. perishable products).

CCS Concepts: • **Information systems** → **Data analytics; Data mining.**

Additional Key Words and Phrases: forecasting, percentage error minimization, factorization machine, retail sales

ACM Reference Format:

Chongshou Li, Brenda Cheang, Zhixing Luo, and Andrew Lim. 2020. An Exponential Factorization Machine with Percentage Error Minimization to Retail Sales Forecasting. *ACM Trans. Knowl. Discov. Data.* xx, x, Article xxx (September 2020), 31 pages. <https://doi.org/10.1145/xxxx.xxxx>

*corresponding author: Chongshou Li (iselc@nus.edu.sg)

Authors' addresses: Chongshou Li, iselc@nus.edu.sg, National University of Singapore, 1 Engineering Drive 2, Blk E1A #06-25, Singapore; Brenda Cheang, brendacheang@yahoo.com, Red Jasper Holdings, Westech Building, 237 Pandan Loop, Singapore; Zhixing Luo, luozx.hkphd@gmail.com, Nanjing University, 22 Hankou Road, Gulou District, Nanjing, Jiangsu, China; Andrew Lim, isealim@nus.edu.sg, National University of Singapore, 1 Engineering Drive 2, Blk E1A #06-25, Singapore.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

© 2020 Association for Computing Machinery.

1556-4681/2020/9-ARTxxx \$15.00

<https://doi.org/10.1145/xxxx.xxxx>

1 INTRODUCTION

In view of its tremendous impacts to the corporate bottom line, sales forecasting has become less about guessing and more about circumspective calculation. Today, there are a host of predictive methods available to sales forecasting professionals [11, 23, 31]. However, these methods are typically proposed and implemented in areas where significant relevant demand history is available (examples include assortment planning [13], marketing [10], revenue management [12], and innovation diffusion process [7]). In this vein, the challenge concerns the development of dependable sales forecasting models for activities that do not have useful historical data to extrapolate from. Consequently, one of the most challenging forecasting problems retailers face is sales forecasting for new products [6].

Sales forecasting for new products is intricately tied to a retailer’s ongoing concern [33]. It carries a high degree of uncertainty and, therefore, risk for the retailer [25]. For one, new products typically require a longer lead time to market than existing products because of demand uncertainty; this is counterintuitive as speed to market is essential for the successful launch of new products [6]. On the other hand, retailers have to determine the order quantity well before the products reach the stores. Thus, as a result of uncertain demand and long lead time, new products with short life cycles (i.e. seasonal goods) are mostly sold without any replenishments. Therein, miscalculated forecasts could result in significant inventory excesses or shortages.

With that said, in the absence of product sales history, several forecasting models have been developed using proxy data such as point-of-sales (POS) transaction data of existing stock keeping units (SKUs) and expert opinions on product attributes. One of the earliest models used was the utility-based multinomial logit (MNL) model by Fader and Hardie [10]. Fisher and Vaidyanathan [13] developed the exogenous demand model, while Ferreirar et al. [12] employed the regression tree and Chung et al. [7] developed the diffusion process based sales forecast model. However, these models have their own limitations and the results produced by these models that rely on proxy data have not been excellent to say the least. And the performance can be largely improved. Inspired by the aforementioned complexities and the lack of superior solutions, we collaborated with a multinational footwear retailer in Singapore¹ to tackle sales forecasting for new products with long lead time and short life cycle. Thus, this paper proposes a new method which we have coined the exponential factorization machines (EFM) model to solve this vexing problem.

The products sold by our industrial partner are ladies footwear. Each product is associated with a set of attributes. The attributes can be grouped into three categories: (1) visible attributes such as color, size, height etc, (2) latent attributes like designer, factory (producing the products), (3) marketing behavior such as discount, average price. Based on data type, the attributes can also be classified as (1) categorical attributes and (2) numerical attributes. Most attributes are categorical. For each categorical attribute, there is a fixed set of levels and each product must be associated with one level. For example, the set of levels for categorical attribute “*HeelHeightRange*” is “{H, M, L}”, an illustration is given by Figure 1. In terms of business model, these ladies footwear can be classified as “fashion-basic” products according to Abernathy et al. [1] and Caro and Martinezde Albeniz [2]. It combines characteristics of two classical fast fashion models: (1) fast fashion and (2) basic products. The life cycle of these items is as short as eleven weeks sharing the features of the first, while the lead time is as long as six months which is similar to the second. After being sold for eleven weeks, the products must be removed from the shelves due. This makes these fashion products are “*perishable*”.

As seen from Figure 2, we forecast the sales for new products or SKUs (stock keeping units) well before its release. In our study, we forecast sales for the first eleven weeks for each item six months before its launch. The granularity level of product we considered is the most fine-grained level, stock keeping unit (SKU). It is a distinct type of product with unique size and color for sale and Figure 1 displays six SKUs. Moreover, we take into account the store locations and forecast the sales for each new SKU at a store. The SKU-store sales forecasts

¹The retailer has requested not to release the exact name, which does not influence current study.



Fig. 1. An illustration of attribute “*HeelHeightRange*”, its levels and SKU

are consolidated via the SKU-chain sales which is the sales of each SKU over all stores (the entire chain). The SKU-chain sales forecasts are used as the new buying quantities.

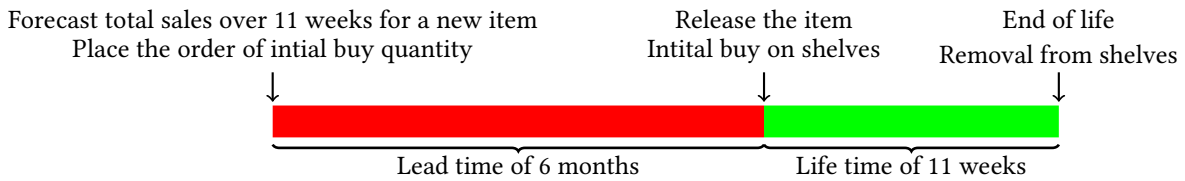


Fig. 2. Sales forecasting process for a new item of the industry partner

We approach this conundrum from a machine learning point of view and consider it as a supervised learning problem. Specifically, the solution is formulated as a regression task where the response (output, target) variable is the sales quantity over 11 weeks for a new SKU at a store (SKU-store sales) and the explanatory variables (factors, features) are product attributes. Our forecast model is based on the factorization machine (FM) which is proposed as a generic predictor and have successfully solved problems in many areas such as recommendation systems [28]. This model has shown excellent prediction capabilities for categorical data. In sales forecasting, most SKU attributes are categorical. And interactions between two attributes are found to be important predictors in three industries by Fisher and Vaidyanathan [13]. Because sales are positive integers in nature, we generalize the original FM [28] of order two with an exponential formulation. In order to estimate parameters of the model, we propose an adaptive batch gradient descent (ABGD) method to minimizing the gap between forecasts and actual sales over training set. The gap is computed by two loss functions: error squares (ES) and percentage error squares (PES). The first is widely used in ordinary least squares (OLS) regression while the second is based on the popular performance indicator, MAPE, in sales and demand forecasting literature.

Our contributions are summarized as follows.

(1) We propose a novel sales forecast model for new SKUs based on historical sales transaction data and attributes data. While previous works have also used historical sales transaction data, none to our knowledge have incorporated SKU attributes and their interactions with the sales transaction data in their approaches. More importantly, inter-attribute-level interactions which were not observed in the training set can still be estimated in our approach; which differs from polynomial regression with the interaction term. We incorporated this element in our approach based on the effectiveness observed by Fisher and Vaidyanathan [13]. However, in their work, instead of attribute-level interactions for new SKUs, they constructed the interactions as a new attribute when implementing their model on three industry categories. At the same time, the SKU-store forecast can be used as buying quantities (after consolidation to the SKU-chain) as well as facilitating initial shipment decisions for the store.

(2) We propose an efficient formulation of the EFM comparing with the formulation of FM [28]. Firstly, an important practical observation from data is that level interactions of the same attribute never appear in both training and test set. For example, attribute “size” has levels “ $\{S, M, L, XL\}$ ” and it is impossible to observe the interaction of “S” and “XL” on any data points. The proposed EFM has removed such infeasible intra-attribute level intersections from the original FM formulation. It results in more efficient feature engineering and reduces the search space comparing than the original FM employing conventional one-hot encoding. Secondly, the EFM employs an exponential formulation for the naturally positive response instead of the original FM with log-transformation. Our computational studies demonstrate that, for positive responses which are not right-skewed distributed, the exponential formulation is more effective than the log-transformation.

(3) We study two loss functions of error square (ES) and percentage error square (PES) for parameter estimation. Although these two are similar in terms of formulation, an important property of PES is found and demonstrated by the computational study that it likely trains the model to underestimate data. This property fits the situation where holding cost is much greater than shortage cost. It is especially true for perishable products such as fruit, vegetables, fresh meat, milk and fashion products, the quality of which deteriorate with time and environmental conditions. For these products, underestimated demand is better than overestimated demand with the same absolute gap. Because overestimation results in buying too much and overstocked items may not be sold before they expire, where the holding costs includes not only conventional inventory cost but also abandonment loss of the expired items. In this case, PES loss function is a good fit for demand/sales forecasting. Moreover, theoretical bounds between objective values of these two optimal solutions are provided (see Theorem 1). We also introduce a new data normalization technique, “instance/sample/row normalization”, which is different from common normalization methods in data mining [18] and is applicable whenever percentage error minimization is the target under the least square framework for linear models.

(4) We propose a practical and effective feature selection method solution for the desired application. Instead of relying on domain knowledge from experts (e.g. working with sales managers in the work of [12]), the proposed greedy forward selection method with cross-validation selects a subset of useful attributes and intra-attribute level iterations from all possible candidates, and is based on the proposed exponential factorization machine (EFM) model and historical training data.

2 LITERATURE REVIEW

The proposed exponential factorization machines (EFM) model is based on the factorization machines (FM) model proposed by Rendle [27]. Likened to support vector machines (SVMs), FMs are said to be excellent general predictors that work for any real-valued feature vector. Nevertheless, FMs are said to be a more superior class of models because they were developed to address the weaknesses of SVMs; by integrating the advantages of SVMs with factorization models [34]. For instance, unlike SVMs, FMs can model all interactions between variables using factorized parameters (by mapping the interactions to a low dimensional space). Therefore, FMs are capable of modeling n -way variable interactions, where n is the number of polynomial order [28]. However, issues such as numerical stability for some optimization methods have obstructed the generalization of n . Thus, the order has been fixed at two as it performs very well in current implementations. As well, FMs not only perform particularly well on sparse data, they are able to compute interactions even in instances with high sparsity [14, 29].

In another instance of the superiority of FMs over SVMs, Rendle [27] established that the model equation of FMs reduce the polynomial computation time to linear complexity, and thus FMs can be optimized directly. At the same time, FMs can mimic specialized factorization models such as parallel factor analysis, matrix factorization, pairwise interaction tensor factorization (PITF) and factorizing personalized Markov chains (FPMC) that are not applicable as generic predictors as these models are not only derived individually for each task, they require

special input data. Thus, the combination of reduced computation time and the ease of application even for users without deep knowledge of factorization models, have made FMs attractive to researchers.

All in all, FMs are popularly applied in collaborative recommendation systems such as recommending movies and music. As they have shown excellent prediction capabilities, they are increasingly being applied to systems such as stock market prediction [3]; and they are recently extended with neural network as deep factorization machines for click-through rate (CTR) prediction and knowledge tracing [16, 19, 35, 36]. What motivated us to propose a FM-based approach to our sales forecasting problem are the common elements observed between our problem and collaborative recommendation systems. Firstly, FMs have shown excellent prediction capabilities for categorical data in collaborative recommendation systems. In sales forecasting, most SKU attributes are categorical too. Second of all, the task from both our problem and recommendation systems are regressions. And as stated earlier, FMs can model all interactions between variables using factorized parameters. And as Fisher and Vaidyanathan [13] found this to be important predictors, we believe it is useful in our case as the model would take interaction between two variables into account and can estimate the effects of new interactions even if the inter-attribute level interactions are not observed. Thus, because sales are positive integers in nature, we employ the exponential formulation of FM and name it as exponential FM (EFM) to it.

FMs can be trained by three algorithms, namely stochastic gradient descent (SGD) [27], alternating least-square (ALS) [30], and Markov Chain Monte Carlo (MCMC) [14]. The performance of SGD depends largely on the learning rate, which is one of the hyper-parameters. The ALS can be applied only when there exists an analytical solution to the minimization problem of parameter estimation. At the same time, both SGD and ALS need to determine regularization hyper-parameters. At the same time, both SGD and ALS need to determine regularization hyper-parameters. Rendle et al. [28] pointed out that not only does MCMC result in fewer hyper-parameters, but those that need to be specified are not as effective to their initial values. Thus, while the ALS cannot be used because no analytical solutions are viable for the parameter minimization problem for the proposed EFM (due to the exponential formulation), the MCMC is also unsuitable as the response variable in our study (“sales”) does not follow a Gaussian distribution because it is a positive integer in nature. Other distributions such as Poisson are also not found to be a good fit for existing sales data. As a result, we propose an adaptive batch gradient descent (ABGD) algorithm to train the EFM. At each iteration, all training data are used as a batch because the data size is not large and around 5,000. When close to the optimal, ABGD slows down the learning rate to make sure that it converges to the optimum.

When training the EFM model, we also need to address the feature selection problem. Features are SKU attributes and attribute interactions. We modeled SKU sales as attributes as it is commonly used in consumer choice literature where each SKU is defined by a set of attribute levels [10]. Attributes used in this area are salient or consumer recognizable such as package size, color. However, we extend the scope and additionally take both latent attributes (e.g. designer) and marketing attributes (e.g. price and discount percentage) into account. In total, we identified 45 attributes in our database. And although increasing the number of attributes likely improves the fitting power for training set, it brought upon the problem of how to select the most useful attributes and interactions. If we consider all of them in the model, it would probably result in over-fitting. This is a typical feature selection problem and is known to be NP-hard [17]. In the literature, Chen et al. [3] and Chen et al. [4] addressed the problem of feature selection for FM. However, the structure of EFM is different from FM in that it disables the existing feature selection methods. As such, we propose a greedy forward stepwise feature selection (GFSFS) algorithm with k -fold cross-validation to tackle this challenge.

Finally, there are many performance indicators in terms of forecasting performance evaluation. The most popular ones are mean absolute error (MAE) (or mean absolute deviation (MAD)), mean squared error (MSE) and mean absolute percentage error (MAPE) [21]. Among other things, the first two indicators are scale-dependent while the third is a scale-independent measure [20]. A survey showed that MAPE and MAE are the top two measures used in the industry; that on average, MAPE of 77% is acceptable in all industries while MAPE of 76%

is acceptable for consumer product industries [21]. Makridakis [24], as well as, Goodwin and Lawton [15] also reported that a property of MAPE is that it treats overestimation differently from underestimation. Therefore, in this study, we develop a percentage error square (PES) loss function based on MAPE and identify its relationships with the classical least square regression.

3 EXPONENTIAL FACTORIZATION MACHINE (EFM) SALES FORECAST MODEL

We model the sales quantity for each stock-keeping unit (SKU) of a store by considering both attributes and the pairwise interactions. Here, we introduce some notations related to the model.

- $N = \{1, 2, \dots, n\}$, index set of all SKUs
- $M = \{1, 2, \dots, m\}$, index set of all stores
- $A = \{1, \dots, a\}$, index set of categorical attributes (e.g, color, size)
- $L^i = \{1, \dots, l_i\}$, $i \in A$, index set of levels of attribute i ($i \in A$) (e.g, "red" of color)
- $I = \{(i, j) \mid i < j, i \in A, j \in A\}$, index set of all possible pairwise interactions between categorical attributes of set A (e.g. ("color", "size") is interaction between "color" and "size")
- d_{is} , response variable, actual sales of SKU i at store s for the first eleven weeks after launching, positive integer ($d_{is} \in \mathbb{N}_{>0}$), $i \in N$, $s \in M$
- x_{cj}^{is} , binary explanatory variable, 1 if SKU i at store s is associated with level j of categorical attribute c , and 0 otherwise, $i \in N$, $s \in M$, $c \in A$ and $j \in L^c$
- \hat{d}_{is} , sales forecast for SKU i at store s for the first eleven weeks after launching, $i \in N$, $s \in M$

Note that for each categorical attribute, each SKU is for one and only one level. The missing value is treated as a synthetic level. The proposed EFM is formulated as follows.

$$\begin{aligned} \hat{d}_{is} = & \exp(\beta_0 + \sum_{c \in A} \sum_{j \in L^c} \beta_{cj} x_{cj}^{is}) \\ & + \sum_{(c, c') \in I} \sum_{j \in L^c} \sum_{j' \in L^{c'}} x_{cj}^{is} x_{c'j'}^{is} \langle \mu_{cj}, \mu_{c'j'} \rangle \end{aligned} \quad (1)$$

The sales forecast is generated by an exponential transformation of a factorization machine because sales data are positive integers, and we name this model as exponential factorization machine (EFM). The intersection term $\sum_{(c, c') \in I} \sum_{j \in L^c} \sum_{j' \in L^{c'}} x_{cj}^{is} x_{c'j'}^{is} \langle \mu_{cj}, \mu_{c'j'} \rangle$ only considers feasible intra-attribute level intersections. It reduces feature space from $\mathcal{O}(C_{l_m \times a}^2)$ of FM's one-hot encoding to $\mathcal{O}(C_a^2 \times l_m^2)$ of the EFM², where $l_m = \max_{i \in A} \{l^i\}$ is the largest number of attribute levels, C_y^2 is to choose 2 from y and $C_y^2 = \frac{y \times (y-1)}{2}$. However, the feasible interactions from these two spaces are the same. The EFM is much more efficient than the FM with one-hot encoding. Below are parameters (or coefficients) to be estimated.

- β_0 : global bias (or intercept)
- β_{cj} : parameter for attribute level j of attribute c , $c \in A$, $j \in L^c$
- $\langle \mu_{cj}, \mu_{c'j'} \rangle$: factorized effect on sales of interaction between level j of attribute c and level j' of attribute c' , $(c, c') \in I$, $j \in L^c$ and $j' \in L^{c'}$. Here $\langle \cdot, \cdot \rangle$ is the dot product of two vectors of size f which is a hyper-parameter and also known as dimensionality of the factorization [28]. It is defined as,

$$\langle \mu_{cj}, \mu_{c'j'} \rangle = \sum_{p=1}^f \mu_{cj,p} \mu_{c'j',p}. \quad (2)$$

The set of attributes associated with parameter $\mu_{cj,p}$ is denoted as

² $C_{l_m \times a}^2 > C_a^2 \times l_m^2$ can be easily proved by expanding them and omitted here.

- $A^I = \{c | \exists c' \in A \text{ such that either } (c, c') \in I \text{ or } (c', c) \in I\}$.

Currently, because set I contains all possible interactions between any two attributes c and c' in set A , set A^I is exactly the same as A . However, this might result in over-fitting. Thus, we address this issue in Section 5.1, where this notation will be used.

Now we introduce new notations for learning algorithm development: $\Theta = \{\beta_0, \beta_{c_j}, \mu_{c'j', p} | c \in A, j \in L^c, c' \in A^I, j' \in L^{j'}, p \in \{1, 2, \dots, f\}\}$, is a set of model parameters while x^{is} is a vector of explanatory variables of size l ($l = \sum_{i=1}^a l_i$) for SKU i at store s and $x^{is} = (x_{11}^{is}, x_{12}^{is}, \dots, x_{1(l_1-1)}^{is}, x_{1l_1}^{is}, x_{21}^{is}, x_{22}^{is}, \dots, x_{2l_2}^{is}, \dots, x_{a1}^{is}, x_{a2}^{is}, \dots, x_{al_a}^{is})$.

Then sales forecast of SKU i at store s , \hat{d}_{is} , can be considered as a function of explanatory variables and parameters, $\hat{d}_{is} = f(x^{is}; \Theta)$. As identified by Rendle et al. [30], an important property of the factorization machine model is multilinearity. The customized results for the EFM are: for any $i \in N$, $s \in M$ and a single parameter $\theta \in \Theta$, sales forecast \hat{d}_{is} is re-written as,

$$\hat{d}_{is} = f(x^{is}; \Theta) = \exp(\theta h_{(\theta)}(x^{is}) + g_{(\theta)}(x^{is})), \quad (3)$$

for every $\theta \in \Theta$. Below are explicit formulations of $h_{(\theta)}(x^{is})$ and $g_{(\theta)}(x^{is})$ for different θ . We can see that and both $h_{(\theta)}(x^{is})$ and $g_{(\theta)}(x^{is})$ are independent from θ .

- Case 1: $\theta = \beta_0$

$$\begin{cases} h_{(\beta_0)}(x^{is}) = 1 \\ g_{(\beta_0)}(x^{is}) = \sum_{c \in A} \sum_{j \in L^j} \beta_{c_j} x_{c_j}^{is} + \sum_{(c, c') \in I} \sum_{j \in L^c} \sum_{j' \in L^{c'}} x_{c_j}^{is} x_{c'_j'}^{is} \sum_{p=1}^f \mu_{c_j, p} \mu_{c'_j', p} \end{cases} \quad (4)$$

- Case 2: $\theta = \beta_{c^*j^*}$, $c^* \in A$, $j^* \in L^{c^*}$

$$\begin{cases} h_{(\beta_{c^*j^*})}(x^{is}) = x_{c^*j^*}^{is} \\ g_{(\beta_{c^*j^*})}(x^{is}) = \beta_0 + \sum_{c \in A, c \neq c^*} \sum_{j \in L^c} \beta_{c_j} x_{c_j}^{is} + \sum_{j \in L^{c^*}, j \neq j^*} \beta_{c^*j} x_{c^*j}^{is} + \\ \sum_{(c, c') \in I} \sum_{j \in L^c} \sum_{j' \in L^{c'}} x_{c_j}^{is} x_{c'_j'}^{is} \sum_{p=1}^f \mu_{c_j, p} \mu_{c'_j', p} \end{cases} \quad (5)$$

- Case 3: $\theta = \mu_{c^*j^*, p^*}$, $c^* \in A^I$, $j^* \in L^{c^*}$, $p^* \in \{1, \dots, f\}$

$$\begin{cases} h_{(\mu_{c^*j^*, p^*})}(x^{is}) = x_{c^*j^*}^{is} \sum_{(c, c^*) \in I \cup (c^*, c) \in I} \sum_{j \in L^c} x_{c_j}^{is} \mu_{c_j, p^*} \\ g_{(\mu_{c^*j^*, p^*})}(x^{is}) = \beta_0 + \sum_{c \in A} \sum_{j \in L^c} \beta_{c_j} x_{c_j}^{is} + \\ \sum_{j \in L^{c^*}, j \neq j^*} \sum_{(c, c^*) \in I \cup (c^*, c) \in I} \sum_{j' \in L^c} x_{c^*j^*}^{is} x_{c'_j'}^{is} \sum_{p=1}^f \mu_{c^*j^*, p} \mu_{c'_j', p} + \\ \sum_{(c, c^*) \in I \cup (c^*, c) \in I} \sum_{j \in L^c} x_{c^*j^*}^{is} x_{c_j}^{is} \sum_{p=1, p \neq p^*}^f \mu_{c^*j^*, p} \mu_{c_j, p} + \\ \sum_{(c, c') \in I, c \neq c^*, c' \neq c^*} \sum_{j \in L^c} \sum_{j' \in L^{c'}} x_{c_j}^{is} x_{c'_j'}^{is} \sum_{p=1}^f \mu_{c_j, p} \mu_{c'_j', p} \end{cases} \quad (6)$$

Now we formalize the data. The database is denoted as $\mathcal{D} = \{(x^{is}, d_{is}) | (i, s) \in S\}$; x^{is} is a vector of all binary explanatory variables for SKU i at store s ; d_{is} represents actual sales and also the target (or response) variable; the observations (or samples, instances) index set is noted as $S = \{(i, s) | i \in N, s \in M, d_{is} > 0\}$. Here we only consider valid samples associated with positive sales. In this study, if a SKU i is not carried at store s ($d_{is} = 0$), we cannot observe any sales transactions and those observations would not be included in \mathcal{D} . Observations with stock-outs are not considered as well. The rationale is two-fold:

- (1) One assumption of the EFM forecast model is that items are always available on shelves for the first eleven weeks after a launch. Including observations with stock-outs violates this assumption.
- (2) Because the time period we considered is the first eleven weeks after launching, we found that few SKUs were stocked-out during that period in the database. Instead, the retailer is usually faced with over-stocked quantities in the end of this period and sells them off through markdowns. Therefore, in our application, there is no need to take stock-outs into account.

In order to estimate parameters, select features and evaluate model performance, we divide the data into training and holdout test sets which are denoted as $\mathcal{D}^T = \{(x^{is}, d_{is}) | (i, s) \in T \subset S\}$ and $\mathcal{D}^E = \{(x^{is}, d_{is}) | (i, s) \in E \subset S\}$, respectively.

4 ESTIMATING THE MODEL PARAMETERS: THE LEARNING PROBLEM

Model parameters are estimated by minimizing the gap between forecasts and actual sales over training observations. The gap is computed by *loss* functions. In order to avoid over-fitting, we apply L2 norm regularization. Resultantly, the model parameters estimation problem becomes an unconstrained minimization problem. Then, a batch gradient descent algorithm is developed to solve it.

4.1 Loss Functions

Two loss functions are used: error squares (ES) and percentage error squares (PES). Given training set $\mathcal{D}^T = \{(x^{is}, d_{is}) | (i, s) \in T\}$ and corresponding sales forecasts $\{\hat{d}_{is} | (i, s) \in T\}$, they are defined as follows.

- Error squares (ES):

$$\mathcal{L}^{\text{ES}}(\Theta) = \frac{1}{2} \sum_{(i,s) \in T} (\hat{d}_{is} - d_{is})^2 = \frac{1}{2} \sum_{(i,s) \in T} [f(x^{is}; \Theta) - d_{is}]^2 \quad (7)$$

- Percentage error squares (PES):

$$\mathcal{L}^{\text{PES}}(\Theta) = \frac{1}{2} \sum_{(i,s) \in T} \left(\frac{\hat{d}_{is} - d_{is}}{d_{is}} \right)^2 = \frac{1}{2} \sum_{(i,s) \in T} \left[\frac{f(x^{is}; \Theta) - d_{is}}{d_{is}} \right]^2 \quad (8)$$

The PES function poses an important property. Given the same absolute error, the PES loss function penalizes data points with small actual sales more heavily than that with large actual sales. For instance, when actual sales is 120 and forecast is 100 which is underestimated and with an absolute error of 20, the PES function value is $(\frac{20}{120})^2 = (16.67\%)^2 = 0.0278$. However, when actual sales is 100 and forecast is 120 (which is overestimated with the same absolute error of 20), the PES function value is $(\frac{20}{100})^2 = (20\%)^2 = 0.04$ which is greater than 0.0278. Therefore, in the fitting process of the PES loss function, data points with small actual values dominate those with large values. In this light, the PES likely trains the model to underestimate the data, which is demonstrated by our computational results in Section 6. This intuition can also be integrated into the classic newsvendor problem, where the overall loss is usually measured as:

$$\text{overall loss} = \text{unit holding cost} \times (\hat{d}_{is} - d_{is})^+ + \text{unit shortage cost} \times (d_{is} - \hat{d}_{is})^+ \quad (9)$$

Here $(\hat{d}_{is} - d_{is})^+ = \max\{(\hat{d}_{is} - d_{is}), 0\}$. The PES loss function can be used when the unit holding cost is much greater than the unit shortage cost (e.g. perishable products). Different from the above newsvendor loss function which is non-differentiable at $\hat{d}_{is} = d_{is}$ and not applied in current study, the PES is differentiable with respect to the forecast and eases the solution development.

Now we investigate the relationship between optimalities for these two losses. The following Theorem 4.1 demonstrates that it is impossible to minimize the PES and ES loss simultaneously for any training data.

THEOREM 4.1. *Given training set $\mathcal{D}^T = \{(x_{is}, d_{is}) | (i, s) \in T\}$, forecasts $\{\hat{d}_{is} | (i, s) \in T\}$, ES minimizer*

$$\Theta_{\text{ES}}^* = \arg \min_{\Theta \in \mathbb{R}^K} \mathcal{L}^{\text{ES}}(\Theta) = \arg \min_{\Theta \in \mathbb{R}^K} \sum_{(i,s) \in T} [\hat{d}_{is} - d_{is}]^2 \quad (10)$$

and PES minimizer

$$\Theta_{PES}^* = \arg \min_{\Theta \in \mathbb{R}^K} \mathcal{L}^{PES}(\Theta) = \arg \min_{\Theta \in \mathbb{R}^K} \sum_{(i,s) \in T} \left[\frac{\hat{d}_{is} - d_{is}}{d_{is}} \right]^2, \quad (11)$$

the following results hold:

(a) $\mathcal{L}^{ES}(\Theta_{PES}^*)$ and $\mathcal{L}^{ES}(\Theta_{ES}^*)$ satisfy:

$$\mathcal{L}^{ES}(\Theta_{ES}^*) \leq \mathcal{L}^{ES}(\Theta_{PES}^*) \leq \frac{d_{\max}^2}{d_{\min}^2} \mathcal{L}^{ES}(\Theta_{ES}^*), \quad (12)$$

(b) $\mathcal{L}^{PES}(\Theta_{PES}^*)$ and $\mathcal{L}^{PES}(\Theta_{ES}^*)$ satisfy:

$$\mathcal{L}^{PES}(\Theta_{PES}^*) \leq \mathcal{L}^{PES}(\Theta_{ES}^*) \leq \frac{d_{\max}^2}{d_{\min}^2} \mathcal{L}^{PES}(\Theta_{PES}^*). \quad (13)$$

Here $\Theta = \{\beta_0, \beta_{c,j}, \mu_{c',j',p} | c \in A, j \in L^c, c' \in A^I, j' \in L^{c'}, p \in \{1, 2, \dots, f\}\}$, $K = 1 + \sum_{i \in A} l_i + \sum_{i \in A} l_i \times f$, l_i is the number of attribute levels of attribute i ($i \in A$) and d_{\max} is the maximum actual sales in training set \mathcal{D}^T , $d_{\max} = \max\{d_{is} | (i, s) \in T\}$ while d_{\min} is the minimum, $d_{\min} = \min\{d_{is} | (i, s) \in T\}$.

See Appendix A for the proof. An interesting observation is the symmetry structure in above bounds. Minimizing ES would lead to the non-minimum PES and the objective value is upper bounded by $\frac{d_{\max}^2}{d_{\min}^2}$ times of the optimal. So does minimizing the percentage error squares. Although above bounds are quite loose, they are general results for any data \mathcal{D}^T and forecasts $\{\hat{d}_{is} | (i, s) \in T\}$ where no distribution assumptions are required by the data. Moreover, in the proof, we do not use any properties of EFM model and learning algorithms, and the results can be extended to other forecast models and learning algorithms. With Theorem 4.1, the difference between PES and ES minimization is determined by ratio $\frac{d_{\max}^2}{d_{\min}^2}$ which can be viewed as an indicator of data variance. Particularly, it indicates the distance between the minimal and maximal boundaries of responses, $\{d_{is} | (i, s) \in T\}$, of data \mathcal{D}^T . For a trivial case, if $d_{\max} = d_{\min}$, all responses in training set are the same. And the ES minimizer, Θ_{ES}^* , is the same as the PES minimizer, Θ_{PES}^* . If $\frac{d_{\max}^2}{d_{\min}^2}$ is large, the ES minimizer can be significantly different from the PES minimizer. We illustrate the theoretical relationships between PES and ES estimation on multiple models and datasets in Section 6.

4.2 Optimization Tasks

We apply L2 regularization and the parameter estimation becomes the following optimization problem.

$$\Theta^* = \operatorname{argmin}_{\Theta \in \mathbb{R}^K} J(\Theta) = \operatorname{argmin}_{\Theta \in \mathbb{R}^K} \mathcal{L}(\Theta) + \frac{1}{2} \sum_{\theta \in \Theta} \lambda_{\theta} \theta^2 \quad (14)$$

Here $l(\Theta)$ is either the ES or the PES loss function. Below are formulations for the two cases.

- Case 1: optimization task of the ES loss function:

$$\Theta_{ES}^* = \operatorname{argmin}_{\Theta \in \mathbb{R}^K} J^{ES}(\Theta) = \operatorname{argmin}_{\Theta \in \mathbb{R}^K} \frac{1}{2} \sum_{(i,s) \in T} (\hat{d}_{is} - d_{is})^2 + \frac{1}{2} \sum_{\theta \in \Theta} \lambda_{\theta} \theta^2. \quad (15)$$

- Case 2: optimization task of the PES loss function:

$$\Theta_{PES}^* = \operatorname{argmin}_{\Theta \in \mathbb{R}^K} J^{PES}(\Theta) = \operatorname{argmin}_{\Theta \in \mathbb{R}^K} \frac{1}{2} \sum_{(i,s) \in T} \left(\frac{\hat{d}_{is} - d_{is}}{d_{is}} \right)^2 + \frac{1}{2} \sum_{\theta \in \Theta} \lambda_{\theta} \theta^2. \quad (16)$$

Here $\lambda_\theta \in \mathbb{R}^+$ is the regularization hyper-parameter for θ which controls the regularization importance compared with the loss term.

4.3 Adaptive Batch Gradient Descent

In this subsection, we describe the adaptive batch gradient descent (ABGD) algorithm to solving optimization tasks. The ABGD computes the optimal parameters iteratively. At each iteration, it looks at all training samples and performs updates simultaneously. For $\theta \in \Theta$, the rule for parameter update is formulated as,

$$\theta^{\text{new}} \leftarrow \theta^{\text{old}} - \eta \left[\frac{\partial l(\Theta)}{\partial \theta} \Big|_{\theta=\theta^{\text{old}}} + \lambda_\theta \theta^{\text{old}} \right], \quad (17)$$

where $\eta \in \mathbb{R}^+$ is the learning rate (or step size) for gradient descent method while $l(\Theta)$ is a general loss function and either the ES or the PES loss. Below are gradients for the two cases.

- Case 1: ES loss function

$$\frac{\partial l^{\text{ES}}(\Theta)}{\partial \theta} = \sum_{(i,s) \in T} (\hat{d}_{is} - d_{is}) \hat{d}_{is} h_{(\theta)}(x^{is}) \quad (18)$$

- Case 2: PES loss function

$$\frac{\partial l^{\text{PES}}(\Theta)}{\partial \theta} = \sum_{(i,s) \in T} \frac{(\hat{d}_{is} - d_{is}) \hat{d}_{is}}{d_{is}^2} h_{(\theta)}(x^{is}) \quad (19)$$

Here \hat{d}_{is} is forecast of training sample (i, s) based on current parameters and it is computed by equation (1). We also utilize the multilinearity of factorization machine for deriving the above results. An observation is that there is a common part in the above derivative for any θ . It is $(\hat{d}_{is} - d_{is}) \hat{d}_{is}$ and $\frac{(\hat{d}_{is} - d_{is}) \hat{d}_{is}}{d_{is}^2}$ for the cases of ES and PES loss function, respectively, and does not vary with θ . At each iteration, we pre-compute the common term for each sample (i, s) before updating parameters. As such, the efficiency is improved. For sample (i, s) , we define v_{is} as follows.

$$v_{is} = \begin{cases} \eta (\hat{d}_{is} - d_{is}) \hat{d}_{is} & \text{if } l(\cdot) = l^{\text{ES}}(\cdot) \\ \eta \frac{(\hat{d}_{is} - d_{is}) \hat{d}_{is}}{d_{is}^2} & \text{if } l(\cdot) = l^{\text{PES}}(\cdot) \end{cases} \quad (20)$$

Note that learning rate is taken as the common term as well. Then the parameter update rule is expressed as

$$\theta^{\text{new}} \leftarrow \theta^{\text{old}} - \sum_{(i,s) \in T} v_{is} h_{(\theta)}(x^{is}) - \eta \lambda_\theta \theta^{\text{old}}. \quad (21)$$

Algorithm 1 is pseudo-code of ABGD procedure. At iteration itr , parameters are updated first and then training error TE_{itr} is computed for the parameters. The definition of training error varies with the loss functions; it is mean absolute error (MAE) for the ES loss function and mean absolute percentage error (MAPE) for the PES loss function. Below is the mathematical formulation.

$$\text{TE}_{\text{itr}} = \begin{cases} \frac{1}{|T|} \sum_{(i,s) \in T} |\hat{d}_{is}(\theta^{\text{itr}}) - d_{is}| & \text{if } l(\cdot) = l^{\text{ES}}(\cdot) \\ \frac{1}{|T|} \sum_{(i,s) \in T} \left| \frac{\hat{d}_{is}(\theta^{\text{itr}}) - d_{is}}{d_{is}} \right| & \text{if } l(\cdot) = l^{\text{PES}}(\cdot) \end{cases} \quad (22)$$

The ABGD adopts learning rate η if two conditions are satisfied: (1) updated solution is close to the optimal, $\text{TE}_{\text{itr}+1} < \epsilon$, and (2) training error is increased, $\text{TE}_{\text{itr}+1} > \text{TE}_{\text{itr}}$. Note that the threshold value ϵ is 0.1 for PES loss function while it is 1.0 for ES loss function. The first condition is necessary because the gradient descent jumps and can produce solutions with increased training error at the initial iterations. The second condition is to make sure that current learning rate is too large to reduce the training error. Note that the stopping criterion in the

algorithm is set by the maximum number of iterations, maxInteractions , which is a hyper-parameter and are generated by the grid search.

The ABGD solves both $J^{\text{ES}}(\Theta)$ and $J^{\text{PES}}(\Theta)$ minimization depending on the input loss function type. If the input loss function is the ES loss $J^{\text{ES}}(\Theta)$, it produces the optimal parameters to J^{ES} minimization; if it is the PES loss $J^{\text{PES}}(\Theta)$, it generates the optimal parameters to J^{PES} minimization. Finally, note that L2 regularization is applied to prevent outlier values of paramters. However, this does not sufficiently overcome the problem. As a result, we propose a forward stepwise selection algorithm with greedy heuristics and k -fold cross validation to select useful features. This is elaborated in the following section.

Algorithm 1 Adaptive Batch Gradient Descent

```

//Inputs:
//training data,  $\mathcal{D}^T$ ; subset of attributes, sA ( $\text{sA} \subset A$ ); subset of interactions, sI ( $\text{sI} \subset I$ )
//loss function  $l(\cdot)$ , either ES or PES loss
//Output: the optimal parameters,  $\Theta^*$ , which minimizes either  $J^{\text{ES}}$  or  $J^{\text{PES}}$ , depending on the input loss.
1: procedure ABGD( $\mathcal{D}^T$ , sA, sI,  $l(\cdot)$ )
2:   construct attribute subset of current interactions,  $\text{sAI} \leftarrow \{i \mid \exists j \in A \text{ such that either } (i, j) \in \text{sI} \text{ or } (j, i) \in \text{sI}\}$ 
3:   initialize a step counter,  $\text{itr} \leftarrow 0$ 
4:   initialize parameters,  $\beta_0^{\text{itr}} \leftarrow 0$ ,  $\beta_{ij}^{\text{itr}} \leftarrow 0 \forall i \in \text{sA}, j \in L^i$ ,  $\mu_{ij,p}^{\text{itr}} \sim \mathcal{N}(0, \sigma) \forall i \in \text{sAI}, j \in L^i, p \in \{1, 2, \dots, f\}$ 
5:   initialize training error of iteration 0,  $\text{TE}_0$ , to large number,  $\text{TE}_0 \leftarrow \text{Large Number}$ 
6:   while stopping criterion is not satisfied ( $\text{iter} \leq \text{maxInteractions}$ ) do
7:     for each  $(i, s) \in T$ , compute  $v_{is}$  by equation (20)
8:      $\beta_0^{\text{itr}+1} \leftarrow \beta_0^{\text{itr}} - \sum_{(i,s) \in T} v_{is} h_{(\beta_0)}(x^{is}) - \eta \lambda \beta_0^{\text{itr}}$ 
9:     for  $c \in \text{sA}, j \in L^i$  do
10:       $\beta_{cj}^{\text{itr}+1} \leftarrow \beta_{cj}^{\text{itr}} - \sum_{(i,s) \in T} v_{is} h_{(\beta_{cj})}(x^{is}) - \eta \lambda \beta_{cj}^{\text{itr}}$ 
11:     end for
12:     for  $c \in \text{sAI}, j \in L^c, p \in \{1, 2, \dots, f\}$  do
13:       $\mu_{cj,p}^{\text{itr}+1} \leftarrow \mu_{cj,p}^{\text{itr}} - \sum_{(i,s) \in T} v_{is} h_{(\mu_{cj,p})}(x^{is}) - \eta \lambda \mu_{cj,p}^{\text{itr}}$ 
14:     end for
15:     compute training error,  $\text{TE}_{\text{itr}+1}$ , based on parameters  $\{\theta^{\text{itr}+1}\}$  following equation (22)
16:     if  $\text{TE}_{\text{itr}+1} < \epsilon$  (solution is close to the optimal) and  $\text{TE}_{\text{itr}+1} > \text{TE}_{\text{itr}}$  then
17:        $\eta \leftarrow \eta/2$ 
18:     end if
19:      $\text{TE}_{\text{itr}} \leftarrow \text{TE}_{\text{itr}+1}$ ,  $\text{itr} \leftarrow \text{itr} + 1$ 
20:   end while
21:   Return  $\Theta^* = \{\beta_0^{\text{itr}}, \beta_{cj}^{\text{itr}}, \mu_{c'j',p}^{\text{itr}} \mid c \in \text{sA}, j \in L^c, c' \in \text{sAI}, j' \in L^{c'}, p \in \{1, 2, \dots, f\}\}$ 
22: end procedure

```

5 TRAINING AND FORECASTING

5.1 Feature Selection

Feature selection is traditionally known as a challenging problem in machine learning. In the EFM model, there are two kinds of features: (1) single attributes A and (2) pairwise interactions between attributes I . In a training set for a class of around 5,000 observations, there are 45 attributes and 990 attribute interactions ($C_{45}^2 = \frac{45 \times 44}{2} = 990$). Sizes of training and test sets are provided by Table 1 of Section 6.1.1. Putting all features into the model can result in minimal losses over training set but likely produces large forecasting error for the hold-out test set. Not all of them are useful and only subsets of A and I should be taken into account. In order to address this over-fitting

problem, we propose a greedy forward stepwise feature selection (GFSFS) algorithm. The GFSFS searches feature subset space iteratively. At each iteration, it augments features by greedily adding a subset of either attributes or pairwise interactions which has the greatest *potential* to improve forecasting accuracy over the test set. Once a new setting of features is built, we investigate how good it is. Although the aim is to find the setting which produces smallest forecasting error over the hold-out test set, the test set is unseen in the training stage. Thus, a k -fold cross-validation (CV) is employed to estimate the error. And the goodness is quantified by validation errors of the CV. In the end, a feature setting of either a subset of attributes A or a subset of interactions I is generated by the GFSFS. This resulting subset of attributes is denoted as sA^* while the subset of interactions is represented by sI^* , and the corresponding minimal validation errors of the CV is noted as $\{cvJ_i^* | i \in \{1, 2, \dots, k\}\}$. The GFSFS is displayed by Algorithm 2 and its sub-procedures are present in the following subsections.

Algorithm 2 Greedy Forward Stepwise Feature Selection

```

//Inputs:
//training data,  $\mathcal{D}^T$ ; set of all possible attributes,  $A$ ; set of all possible interactions,  $I$ 
//loss function  $l(\cdot)$ , either ES or PES loss
//Output: ( $sA^*$ ,  $sI^*$ ), the optimal subsets of attributes  $A$  and interactions  $I$ 
1: procedure GFSFS( $\mathcal{D}^T, L, I, l(\cdot)$ )
2:   randomly partition training set  $\mathcal{D}^T$  into  $k$  complementary groups with roughly equal size,  $\Lambda = \{T_i, i \in \{1, 2, \dots, k\} | T_1 \cup T_2 \cup T_3 \cup \dots \cup T_k = T, T_j \cap T_{j'} = \emptyset, j \neq j', j, j' \in \{1, 2, \dots, k\}\}$ 
3:   initialization,  $currentA^0 \leftarrow \emptyset, currentI^0 \leftarrow \emptyset, sd^1 \leftarrow 1, isFeasible_{-1} \leftarrow True, isFeasible_1 \leftarrow True, itr \leftarrow 1$ 
4:   for each  $i \in 1, 2, \dots, k$ , initialize the optimal validation error  $cvJ_i^*$  to be  $cvJ_i^{null}$  given by equation (28)
5:   while current search direction is feasible ( $isFeasible_{sd^{itr}}$  is True) do
6:     feature subset selection,  $\{\Delta A, \Delta I\} \leftarrow FSS(sd^{itr}, currentA^{itr-1}, currentI^{itr-1})$ 
7:      $currentA^{itr} \leftarrow currentA^{itr-1} \cup \Delta A, currentI^{itr} \leftarrow currentI^{itr-1} \cup \Delta I$ 
8:     if either  $\Delta A$  or  $\Delta I$  is non-empty then
9:       perform cross validation,  $\{cvJ_i^{itr} | i \in \{1, 2, \dots, k\}\} \leftarrow CV(currentA^{itr}, currentI^{itr}, k, l(\cdot), \Lambda)$ 
10:      if  $\{cvJ_i^{itr} | i \in \{1, 2, \dots, k\}\}$  is significantly smaller than  $\{cvJ_i^* | i \in \{1, 2, \dots, k\}\}$  then
11:         $sA^* \leftarrow currentA^{itr}, sI^* \leftarrow currentI^{itr}, cvJ_i^* \leftarrow cvJ_i^{itr} \forall i \in \{1, 2, \dots, k\}$ 
12:         $isFeasible_1 \leftarrow True, isFeasible_{-1} \leftarrow True$ 
13:      else
14:         $isFeasible_{sd^{itr}} \leftarrow False, currentA^{itr} \leftarrow currentA^{itr-1}, currentI^{itr} \leftarrow currentI^{itr-1}$ 
15:      end if
16:    else
17:       $isFeasible_{sd^{itr}} \leftarrow False$ 
18:    end if
19:    if search direction  $-sd^{itr}$  is feasible ( $isFeasible_{-sd^{itr}}$  is True) then
20:       $sd^{itr+1} \leftarrow -sd^{itr}$ 
21:    end if
22:     $itr \leftarrow itr + 1$ 
23:  end while
24:  Return ( $sA^*, sI^*$ )
25: end procedure

```

The GFSFS initializes the optimal subsets of attributes, sA^* , and interactions, sI^* , to be empty (see line 3). The model without any features is named as null model where only global bias (or intercept) is considered. There exist analytical solutions for this model and corresponding errors of k validation sets can be directly derived instead of calling sub-procedure CV. The detail is referred to subsection 5.3. Line 4 initializes $\{cvJ_i^* | i \in \{1, 2, \dots, k\}\}$ to

be validation errors of the null model. After that, the GFSFS constructs new subsets of features iteratively. The iteration number is denoted as itr which is initialized to be one. (Null model is investigated at iteration zero.) At iteration itr , the GFSFS augments current feature subsets of either attributes, $\text{currentA}^{\text{itr}-1}$, or interactions, $\text{currentI}^{\text{itr}-1}$. There are two choices: (1) adding attributes and (2) adding interactions. These two choices are noted as search directions at iteration itr and symbolized as sd^{itr} , which is defined as an enumerated variable in the GFSFS and takes two possible values of positive one and negative one. Being positive one means adding attributes while taking value of negative one represents adding interactions. A trick of this setting is that search direction can be easily changed via multiplying search direction by negative one. The search direction for iteration one is initialized to be positive one of adding attributes, $\text{sd}^1 \leftarrow 1$ (see line 3). At iteration itr , depending on search direction sd^{itr} , the sub-procedure FSS greedily selects a new subset of either attributes ΔL or interactions ΔI for augmenting current features at line 6. Subsection 5.4 presents this sub-procedure in detail.

At the same time, there exists a feasibility issue for search directions. The GFSFS cannot always add either attributes or interactions to current features. For instance, at iteration itr , given that search direction, sd^{itr} , is positive one of adding attributes and current subset of attributes, $\text{currentA}^{\text{itr}-1}$, contains all attributes ($\text{currentA}^{\text{itr}-1} = A$), there is no attributes which are out of current subset and it is impossible to augment current features following the given direction. And the direction is not feasible. In order to formalize it, the GFSFS defines boolean variables for labeling the feasibility of two search directions. They are isFeasible_1 and isFeasible_{-1} for positive one of adding attributes and negative one of adding interactions, respectively. A value of True means that it is feasible while a value of False means that it is infeasible. These two boolean variables are initialized to be True (feasible) at line 3.

At iteration itr , based on search direction sd^{itr} , attribute subset $\text{currentA}^{\text{itr}-1}$ and interaction subset $\text{currentI}^{\text{itr}-1}$, (in line 6) the FSS produces one subset of either attributes, ΔA , or interactions, ΔI , only for augmenting features at iteration $\text{itr} - 1$. Then features for current iteration itr are constructed via taking union of features of previous iteration and newly found increments, $\text{currentA}^{\text{itr}} \leftarrow \text{currentA}^{\text{itr}-1} \cup \Delta A$, $\text{currentI}^{\text{itr}} \leftarrow \text{currentI}^{\text{itr}-1} \cup \Delta I$, which is seen in line 7. If both ΔA and ΔI are empty sets, it indicates that current search direction sd^{itr} is not feasible any more, $\text{isFeasible}_{\text{sd}^{\text{itr}}} \leftarrow \text{False}$ (see line 17). If either ΔA or ΔI is not empty, k -fold cross validation is performed for examining the goodness of this new setting at line 9. If the mean of k validation errors, $\{\text{cvJ}_i^{\text{itr}} | i \in \{1, 2, \dots, k\}\}$, is smaller than that of the current optimal validation errors, $\{\text{cvJ}_i^* | i \in \{1, 2, \dots, k\}\}$, at a significance level of 0.05 by a paired t-test, the optimal subsets and validation errors are updated at line 11, $\text{sA}^* \leftarrow \text{currentA}^{\text{itr}}$, $\text{sI}^* \leftarrow \text{currentI}^{\text{itr}}$, $\text{cvJ}_i^* \leftarrow \text{cvJ}_i^{\text{itr}} \forall i \in \{1, 2, \dots, k\}$. Line 12 sets two search directions to be feasible $\text{isFeasible}_1 \leftarrow \text{True}$ and $\text{isFeasible}_{-1} \leftarrow \text{True}$. On the other hand, if average of k validation errors is not significantly smaller, current search direction is set to be infeasible, $\text{isFeasible}_{\text{sd}^{\text{itr}}} \leftarrow \text{False}$, and feature subsets remain the same, $\text{currentA}^{\text{itr}} \leftarrow \text{currentA}^{\text{itr}-1}$ and $\text{currentI}^{\text{itr}} \leftarrow \text{currentI}^{\text{itr}-1}$ (see line 14). The GFSFS employs a depth-first search (DFS) strategy. At each iteration, the search direction is changed from that of last iteration if it is feasible at lines 19 - 21. The GFSFS stops when current search direction is not feasible (see line 5).

The k -fold cross-validation uses the same partition Λ of training set \mathcal{D}^T which is defined at the beginning of the GFSFS algorithm (line 2). As such, k validation sets do not vary with iterations and t -test is employed for the comparison.

5.2 Cross-Validation

In this subsection, we outline the k -fold cross-validation (CV) in Algorithm 3 which is a sub-procedure and used by the GFSFS algorithm. In k -fold cross-validation procedure, the model is fitted and evaluated k times. At iteration i , fold i is held out and used as validation set noted as \mathcal{D}^{cE} . And the remaining $k - 1$ folds are used as training set noted as \mathcal{D}^{cT} . Parameters for selected subset of features are estimated via calling the ABGD procedure. And the validation error is computed and noted as cvJ_i . The error type depends on the loss function

$l(\cdot)$. If it is the ES loss function, validation error cvJ_i is the mean absolute error (MAE) between forecasts and actual sales over validation set \mathcal{D}^{cE} ; otherwise, it is the mean absolute percentage error (MAPE). It is given as follows.

$$cvJ_i = \begin{cases} \frac{1}{|T_i|} \sum_{(j,s) \in T_i} |\hat{d}_{js} - d_{js}| & \text{if } l(\cdot) = l^{ES}(\cdot) \\ \frac{1}{|T_i|} \sum_{(j,s) \in T_i} \left| \frac{\hat{d}_{js} - d_{js}}{d_{js}} \right| & \text{if } l(\cdot) = l^{PES}(\cdot) \end{cases} \quad (23)$$

In the end, we note that, as a rule of thumb, k is set to be five.

Algorithm 3 k -fold Cross-Validation

```

//Inputs:
//attribute level subset, csL; interaction subset, csI
//partition of training data,  $\Lambda = \{T_i, i \in \{1, \dots, k\}\}$ ; loss function  $l(\cdot)$ 
//Outputs:
//objective values of  $k$  validation sets  $\{cvJ_i | i \in \{1, 2, \dots, k\}\}$ 
1: procedure CV(csL, csI,  $l(\cdot)$ ,  $\Lambda$ )
2:   for  $i = 1 \rightarrow k$  do
3:      $\mathcal{D}^{cE} \leftarrow \{(x^{js}, d_{js}) | (j, s) \in T_i\}$ ,  $\mathcal{D}^{cT} \leftarrow \emptyset$ 
4:     for  $g = 1 \rightarrow k$  do
5:       if  $g \neq i$  then
6:          $\mathcal{D}^{cT} \leftarrow \mathcal{D}^{cT} \cup \{(x^{js}, d_{js}) | (j, s) \in T_g\}$ 
7:       end if
8:     end for
9:     estimate parameters by the ABGD,  $\Theta^c \leftarrow \text{ABGD}(\mathcal{D}^{cE}, \text{csL}, \text{csI}, l(\cdot))$ 
10:    based on input loss function  $(\cdot)$ , compute validation error,  $cvJ_i$ , following equation(23)
11:  end for
12:  Return  $\{cvJ_i | i \in \{1, 2, \dots, k\}\}$ 
13: end procedure

```

5.3 Null Model

The null model is used for initialization in feature selection. In the null model, no attributes and interactions are taken into account. Thus, only global bias (or intercept) β_0 is required to be estimated and all other parameters are set to be zero. For training sample $(i, s) \in T$, forecast \hat{d}_{is} by null model is:

$$\hat{d}_{is} = f^{\text{null}}(x^{is}, \Theta) = \exp(\beta_0).$$

Given training set \mathcal{D}^T , we do not consider regularization for estimating β_0 . And two optimization tasks are:

- Case 1: optimization task of the ES loss for null model

$$J_{\text{null}}^{\text{ES}}(\beta_0) = \frac{1}{2} \sum_{(i,s) \in T} [\exp(\beta_0) - d_{is}]^2 \quad (24)$$

- Case 2: optimization task of the PES loss for null model

$$J_{\text{null}}^{\text{PES}}(\beta_0) = \frac{1}{2} \sum_{(i,s) \in T} \left[\frac{\exp(\beta_0)}{d_{is}} - 1 \right]^2. \quad (25)$$

Taking the first derivate and setting it to be zero, we can get analytic solutions to minimizing these two tasks. The optimal parameter $\beta_{0,ES}^{\text{null}*}$ to $J_{\text{null}}^{\text{ES}}(\beta_0)$ minimization is

$$\beta_{0,ES}^{\text{null}*} = \arg \min_{\beta_0 \in \mathbb{R}} J_{\text{null}}^{\text{ES}}(\beta_0) = \log \left[\frac{\sum_{(i,s) \in T} d_{is}}{|T|} \right]. \quad (26)$$

The optimal parameter $\beta_{0,PES}^{\text{null}*}$ to $J_{\text{null}}^{\text{PES}}(\beta_0)$ minimization is

$$\beta_{0,PES}^{\text{null}*} = \arg \min_{\beta_0 \in \mathbb{R}} J_{\text{null}}^{\text{PES}}(\beta_0) = \log \left[\sum_{(i,s) \in T} \frac{1}{d_{is}} \right] - \log \left[\sum_{(i,s) \in T} \frac{1}{d_{is}^2} \right]. \quad (27)$$

Based on these two solutions and partition Λ of training set \mathcal{D}^T , outputs of the CV can be derived and below are validation error $\text{cvJ}_i^{\text{null}}$ for data set T_i for two different loss functions.

$$\text{cvJ}_i^{\text{null}} = \begin{cases} \frac{1}{|T_i|} \sum_{(i,s) \in T_i} |d_{T \setminus T_i}^{\text{ESM}} - d_{is}| & \text{if } l(\cdot) = l^{\text{ES}}(\cdot) \\ \frac{1}{|T_i|} \sum_{(i,s) \in T_i} \left| \frac{d_{T \setminus T_i}^{\text{PESM}}}{d_{is}} - d_{is} \right| & \text{if } l(\cdot) = l^{\text{PES}}(\cdot) \end{cases} \quad (28)$$

Here $d_{T \setminus T_i}^{\text{ESM}} = \frac{\sum_{(i,s) \in T \setminus T_i} d_{is}}{|T \setminus T_i|}$ and $d_{T \setminus T_i}^{\text{PESM}} = \frac{\sum_{(i,s) \in T \setminus T_i} \frac{1}{d_{is}}}{\sum_{(i,s) \in T \setminus T_i} \frac{1}{d_{is}^2}}$.

5.4 Forward Subset Selection

In this subsection, we present the forward subset selection (FSS) algorithm (see Algorithm 4), a key aspect of the GFSFS. The FSS greedily finds a subset of attributes ΔA or interactions ΔI (depending on input search direction sd) in an effort to augment current features. These two sets are initialized to be empty sets in the beginning of the FSS (see line 2). Features are selected at each iteration which can approximately maximally reduce losses between current forecasts and actual sales over training set \mathcal{D}^T . As such, parameters for current attributes csA and interactions csI are estimated first at line 3 by the ABGD method. Fitted sales \hat{d}_{is}^{cf} for training observation (i, s) is computed based on the estimated parameters Θ^{cf} for current features at line 4. Next depending on input search direction sd, a subset set of either attributes or interactions which are not in the current feature set are selected such that losses between forecasts \hat{d}_{is}^{cf} and actual sales d_{is} over all training observations $(i, s) \in T$ is reduced as much as possible. Note that there are two components in the output of Algorithm 4: ΔA and ΔI . And at most one of them is non-empty.

Next, we describe the details of how to select useful features. First, we discuss the case of input search direction sd being equal to positive one of forward attribute subset selection. Considering an out-of-current attribute c ($c \in A \setminus \text{csA}$), adding it to the current attribute subset csA results in new fitted values over training set \mathcal{D}^T . Specifically, for training observation (i, s) , we approximate the new fitted value to be $\hat{d}_{is}^{\text{cf}} \exp(\sum_{j \in L^c} \beta_{cj} x_{cj}^{is})$ based on the assumption that attributes are independent from each other. Adding attribute c does not affect estimation of parameters for existing features. Here $\{\beta_{cj} | j \in L^c\}$ are new parameters to be estimated. Adding attribute c results in new losses as follows.

$$\mathcal{L}^{\text{NEW}} = \begin{cases} \sum_{(i,s) \in T} [\hat{d}_{is}^{\text{cf}} \exp(\sum_{j \in L^c} \beta_{cj} x_{cj}^{is}) - d_{is}]^2 & \text{if } l(\cdot) = l^{\text{ES}}(\cdot) \\ \sum_{(i,s) \in T} \left[\frac{\hat{d}_{is}^{\text{cf}}}{d_{is}} \exp(\sum_{j \in L^c} \beta_{cj} x_{cj}^{is}) - 1 \right]^2 & \text{if } l(\cdot) = l^{\text{PES}}(\cdot) \end{cases} \quad (29)$$

Minimizing the above objective favors attributes with large number of levels, which produces good predictions over training set but is likely to overfit the data. Therefore, we take the number of levels into account and define

Algorithm 4 Forward Subset Selection

```

//Inputs:
//search direction: sd ∈ {1, -1}:
// 1: forward attribute subset selection; -1: forward interaction subset selection
//csA, current attribute subset; csI, current interaction subset
//loss function, l(·)
//Outputs:
//selected subsets of attributes and interactions for adding: {ΔA, ΔI}
1: procedure FSS(sd, csA, csI, l(·))
2:   initialization, ΔA ← ∅, ΔI ← ∅
3:   estimate parameters of current features, Θcf ← ABGD(DT, csA, csI, l(·))
4:   for each training observation (i, s) ∈ T, compute its fitted sales with current features,  $\hat{d}_{is}^{cf} = f(x^{is}; \Theta^{cf})$ 
5:   if sd is positive one of forward attribute subset selection then
6:     for each out-of-current attribute c ∈ A \ csA, compute its minimal loss JFAS,c* following equation (31)
7:     sort {JFAS,c*FAS,c[1]* ≤ JFAS,c[2]* ≤ ⋯ ≤ JFAS,c[|A \ csA|]*
8:     construct subset ΔA to be the top min{b, |A \ csA|} attributes in the above ranking list, ΔA ← {c[i] | 1 ≤ i ≤ min{b, |A \ csA|}}
9:   else
10:    for each out-of-current interaction (c, c') ∈ I \ csI, compute its minimal loss JFIS,(c,c')* following equation (33)
11:    sort {JFIS,(c,c')*FIS,(c,c')[1]* ≤ JFIS,(c,c')[2]* ≤ ⋯ ≤ JFIS,(c,c')[|I \ csI|]*
12:    select the top min{g, |I \ csI|} interactions as ΔI from ranking list subject to two constraints: (1) each two of them can not have common attributes and (2) every interaction can not have common attributes with existing interactions in csI
13:   end if
14:   Return {ΔA, ΔI}
15: end procedure

```

new objective $J_{FAS,c}$ as follows.

$$J_{FAS,c} = \begin{cases} \sum_{(i,s) \in T} [\hat{d}_{is}^{cf} \exp(\sum_{j \in L^c} \beta_{cj} x_{cj}^{is}) - d_{is}]^2 + \lambda_{A,ES} |L^c| & \text{if } l(\cdot) = l^{ES}(\cdot) \\ \sum_{(i,s) \in T} [\frac{\hat{d}_{is}^{cf}}{d_{is}} \exp(\sum_{j \in L^c} \beta_{cj} x_{cj}^{is}) - 1]^2 + \lambda_{A,PES} |L^c| & \text{if } l(\cdot) = l^{PES}(\cdot) \end{cases} \quad (30)$$

Here $|L^c|$ represents the number of levels of attribute c ; $\lambda_{A,ES}$ and $\lambda_{A,PES}$ are hyper-parameters and control the trade-off between the attribute's complexity and its fitting over the training set. Increasing their values leads to a price to pay for selecting attributes with a large number of levels. Minimizing $J_{FAS,c}$ is easily solved by setting the first derivate with respect to $e^{\beta_{cj}}$ to zero. Below is the minimal loss, $J_{FAS,c}^*$, by attribute c for the two loss functions.

$$J_{FAS,c}^* = \text{Min}_{\{\beta_{cj} \in \mathbb{R}, j \in L^c\}} J_{FAS,c}(\beta_{cj}, j \in L^c) = \begin{cases} \sum_{j \in L^c} \sum_{(i,s) \in T_{cj}} (\hat{d}_{is}^{cf} w_{cj}^{ES} - d_{is})^2 + \lambda_{A,ES} |L^c| & \text{if } l(\cdot) = l^{ES}(\cdot) \\ \sum_{j \in L^c} \sum_{(i,s) \in T_{cj}} (r_{is} w_{cj}^{PES} - 1)^2 + \lambda_{A,PES} |L^c| & \text{if } l(\cdot) = l^{PES}(\cdot) \end{cases} \quad (31)$$

Here $T_{cj} = \{(i, s) | x_{cj}^{is} = 1, (i, s) \in T\}$, $w_{cj}^{ES} = \frac{\sum_{(i,s) \in T_{cj}} d_{is} \hat{d}_{is}^{cf}}{\sum_{(i,s) \in T_{cj}} (\hat{d}_{is}^{cf})^2}$, $r_{is} = \frac{\hat{d}_{is}^{cf}}{d_{is}}$ and $w_{cj}^{PES} = \frac{\sum_{(i,s) \in T_{cj}} r_{is}}{\sum_{(i,s) \in T_{cj}} r_{is}^2}$. Note that, for training observation $(i, s) \in T$ and level j of attribute c , either $x_{cj}^{is} = 1$ or $x_{cj}^{is} = 0$ only must be true. For level j ($j \in L^c$), set T_{cj} contains all training observations associated with level j of attribute c . We apply a greedy strategy and select

the top b attributes associated with the smallest b minimal losses as result of ΔA . Note that b is a hyper-parameter and is denoted as the *forward attribute selection depth*. If there are less than b attributes out of current set csA , all of them are put into result subset ΔA . These operations are presented at lines 7 and 8 in Algorithm 4.

Now we proceed to the case of input search direction sd being equal to negative one of forward interaction subset selection. For an out-of-current interaction (c, c') ($(c, c') \in \Gamma \setminus csI$), it is possible that existing interactions in set csI already contains attribute c or/and attribute c' already. The parameters $\{\mu_{c_j, p}\}$ (or/and $\{\mu_{c'_{j'}, p}\}$) for attribute c (or/and attribute c') associated with an interaction exists and have been estimated already. Interactions with common attributes are highly correlated with each other. Adding such interactions probably does not increase the predictive power [5]. For our case, we only take interactions into account if there are no common attributes. We assume that these interactions are independent of each other. Adding new interaction (c, c') does not affect estimation of parameters for existing interactions. We approximate the new fitted value to be $\hat{d}_{is}^{cf} \exp(\sum_{j \in L^c} \sum_{j' \in L^{c'}} x_{c_j}^{is} x_{c'_{j'}}^{is} \sum_{p=1}^f \mu_{c_j, p} \mu_{c'_{j'}, p})$. In order to avoid over-fitting, we also take the interaction's complexity into account; which is the number of interaction levels. The objective FSS minimizes for searching interactions is denoted as $J_{FIS, (c, c')}$ and defined as follows.

$$J_{FIS, (c, c')} = \begin{cases} \sum_{(i, s) \in T} [\hat{d}_{is}^{cf} \exp(\sum_{j \in L^c} \sum_{j' \in L^{c'}} x_{c_j}^{is} x_{c'_{j'}}^{is} \sum_{p=1}^f \mu_{c_j, p} \mu_{c'_{j'}, p}) - d_{is}]^2 + \lambda_{I, ES} |L^c| \times |L^{c'}| & \text{if } l(\cdot) = l^{ES}(\cdot) \\ \sum_{(i, s) \in T} [\frac{\hat{d}_{is}^{cf}}{d_{is}} \exp(\sum_{j \in L^c} \sum_{j' \in L^{c'}} x_{c_j}^{is} x_{c'_{j'}}^{is} \sum_{p=1}^f \mu_{c_j, p} \mu_{c'_{j'}, p}) - 1]^2 + \lambda_{I, PES} |L^c| \times |L^{c'}| & \text{if } l(\cdot) = l^{PES}(\cdot) \end{cases} \quad (32)$$

Here $\{\mu_{c_j, p} | j \in L^c, p = 1, \dots, f\}$ and $\{\mu_{c'_{j'}, p} | j' \in L^{c'}, p = 1, \dots, f\}$ are new parameters to be estimated. Both $\lambda_{I, ES}$ and $\lambda_{I, PES}$ are two hyper-parameters governing the complexity of interaction. We observe that, by equations (30), (32) and ignoring the penalty part for complexity, minimizing $J_{FIS, (c, c')}$ is equivalent to minimizing $J_{FAS, c \times c'}$, where $c \times c'$ is a synthetic attribute and defined as $x_{c \times c', j''}^{is} = x_{c_j}^{is} \times x_{c'_{j'}}^{is}$ and $j'' = (j - 1) \times |L^c| + j'$. The minimal value is denoted as $J_{FIS, (c, c')}^*$ and provided as follows.

$$J_{FIS, (c, c')}^* = \text{Min}_{\{\mu_{c_j}, \mu_{c'_{j'}} \in \mathbb{R}^f, j \in L^c, j' \in L^{c'}\}} J_{FIS, (c, c')}(\mu_{c_j}, \mu_{c'_{j'}}) = \begin{cases} \sum_{j \in L^c} \sum_{j' \in L^{c'}} \sum_{(i, s) \in T_{c_j, c'_{j'}}} (\hat{d}_{is}^{cf} w_{(c, c')}^{ES} - d_{is})^2 + \lambda_{I, ES} |L^c| \times |L^{c'}| & \text{if } l(\cdot) = l^{ES}(\cdot) \\ \sum_{j \in L^c} \sum_{j' \in L^{c'}} \sum_{(i, s) \in T_{c_j, c'_{j'}}} (r_{is} w_{(c, c')}^{PES} - 1)^2 + \lambda_{I, PES} |L^c| \times |L^{c'}| & \text{if } l(\cdot) = l^{PES}(\cdot) \end{cases} \quad (33)$$

Here $T_{c_j, c'_{j'}} = \{(i, s) | x_{c_j}^{is} = 1, x_{c'_{j'}}^{is} = 1, (i, s) \in T\}$, $w_{(c, c')}^{ES} = \frac{\sum_{(i, s) \in T_{c_j, c'_{j'}}} d_{is} \hat{d}_{is}^{cf}}{\sum_{(i, s) \in T_{c_j, c'_{j'}}} (\hat{d}_{is}^{cf})^2}$, $r_{is} = \frac{\hat{d}_{is}^{cf}}{d_{is}}$ and $w_{(c, c')}^{PES} = \frac{\sum_{(i, s) \in T_{c_j, c'_{j'}}} r_{is}}{\sum_{(i, s) \in T_{c_j, c'_{j'}}} r_{is}^2}$.

Similarly to level subset selection, a greedy strategy is applied and the top g interactions associated with the smallest g minimal losses is built as result ΔI . And g is a hyper-parameter and is denoted as the *forward interaction selection depth*. If there are less than g interactions left, all of them are put into ΔI (see lines 11 and 12).

5.5 Overall Procedure

In this section, we present the overall procedure (OP) which is outlined by Algorithm 5. It mainly consists of two stages: training and test stages.

In the training stage, features are selected first in the GFSFS algorithm. The optimal attribute subset sA^* and interaction subset sI^* are produced (line 2 of Algorithm 5). Parameters for the selected attributes and interactions are then estimated by ABGD algorithm (line 3 of Algorithm 5). Note that in the process of feature selection, parameters are needed to be estimated and ABGD is called in CV and FSS. There exist differences between these two estimations in terms of the value of regularization hyper-parameter λ_θ . For feature selection process (line 2), all regularization hyper-parameters $\{\lambda_\theta\}$ in the ABGD are set to zero to estimate parameters to ascertain the importance ranking for attributes and interactions. At the same time, we have also considered the over-fitting

Algorithm 5 Overall Procedure

```

//Inputs:
//training set  $\mathcal{D}^T$ , test set  $\mathcal{D}^E$ , full attribute set  $A$ , full interaction set  $I$ 
//loss function  $l(\cdot)$  of either the ES or the PES
//Outputs:
//forecasts for training set  $\{\hat{d}_{is} | (i, s) \in E\}$ 
//MAPE, MAE and OR at both SKU-chain and SKU-store aggregate levels
1: procedure OP( $\mathcal{D}^T, \mathcal{D}^E, I, l(\cdot)$ )
//Training stage:
2:   select attribute subset  $sA^*$  and interaction subset  $sI^*$  by GFSFS,  $(sA^*, sI^*) \leftarrow \text{GFSFS}(\mathcal{D}^T, l(\cdot), L, I)$ 
3:   estimate parameters by ABGD,  $\Theta^* \leftarrow \text{ABGD}(\mathcal{D}^T, sA^*, sI^*, l(\cdot))$ 

//Test stage:
4:   for each observation  $(i, s) \in E$  in test set, compute forecast  $\hat{d}_{is} = f(x^{is}; \Theta^*)$  following equation (1)
5:   calculate  $\text{MAPE}_{\text{SKU-store}}^E, \text{MAPE}_{\text{SKU-chain}}^E, \text{MAE}_{\text{SKU-store}}^E, \text{MAE}_{\text{SKU-chain}}^E$  following equations from (34) to (37), respectively
6:   Return  $\{\{\hat{d}_{is} | (i, s) \in E\}; \text{MAPE and MAE for both SKU-store and SKU-chain sales forecasting}\}$ 
7: end procedure

```

problem in feature importance measures in equations (30) and (32). However, after features are selected, we re-estimate parameters with the whole training set and they will be used for produce forecasts for the test set. At this step, we should take the regularization into account to avoid over-fitting.

In test stage, forecasts for test set $\{\hat{d}_{is} | (i, s) \in E\}$ are calculated based on the estimated parameters for the EFM model at the training stage. Then we evaluate the performance and compute two indicators: mean absolute percentage error (MAPE) and mean absolute error (MAE). Forecast \hat{d}_{is} is for SKU i at store s . It is at the aggregation level of SKU-store. As discussed in Section 1, forecasts are aimed at facilitating the buying quantity decisions which are usually an aggregate decision for all stores instead of individual stores. Therefore, the SKU-store sales forecasts are consolidated to forecasts for all stores (SKU-chain forecast) and then used as the buying quantity for each SKU. As such, we calculate both MAPE and MAE for both SKU-store and SKU-chain sales forecasting. Given SKU-store forecasts $\{\hat{d}_{is} | (i, s) \in E\}$ and actual sales $\{d_{is} | (i, s) \in E\}$, formulations of the MAPE and MAE for both SKU-store and SKU-chain forecasting are given as follows. We denote all SKUs in test set E as $E^{\text{SKU}} = \{i | \exists s \in M \text{ such that } (i, s) \in E\}$.

- MAPE

$$\text{MAPE}_{\text{SKU-store}}^E = \frac{1}{|E|} \sum_{(i,s) \in E} \left| \frac{\hat{d}_{is} - d_{is}}{d_{is}} \right| \quad (34)$$

$$\text{MAPE}_{\text{SKU-chain}}^E = \frac{1}{|E^{\text{SKU}}|} \sum_{i \in E^{\text{SKU}}} \left| \frac{\sum_{s \in M} \hat{d}_{is} - \sum_{s \in M} d_{is}}{\sum_{s \in M} d_{is}} \right| \quad (35)$$

- MAE

$$\text{MAE}_{\text{SKU-store}}^E = \frac{1}{|E|} \sum_{(i,s) \in E} |\hat{d}_{is} - d_{is}| \quad (36)$$

$$\text{MAE}_{\text{SKU-chain}}^E = \frac{1}{|E^{\text{SKU}}|} \sum_{i \in E^{\text{SKU}}} \left| \sum_{s \in M} \hat{d}_{is} - \sum_{s \in M} d_{is} \right| \quad (37)$$

6 COMPUTATIONAL STUDIES

In this section, we present the computational studies including not only the sales forecasting dataset but also two public datasets. Several important variants of the EFM model are investigated. Comparisons and insights are provided. Note that our proposed algorithms are implemented in Java. Hardware configurations for these studies are given as follows. Computations for hyper-parameter searches are performed on Supermicro Linux server with 4-core Intel Xeon E5-2623V3 (10MB L3, 3.0GHz), 256 RAM and Ubuntu 16.04 operating system. The remaining implementations are conducted on a Dell personal computer with an Intel Core i7-6700 CPU, 3.40 GHz, 32 GB RAM and 64-bit Windows 10 operating system. We investigate different variants and extensions of the proposed EFM model for comparison analysis. For ease of presentation, we denote the method as follows.

- **EFM-PES:** the proposed EFM model with PES loss function and the PES based feature selection. Hyperparameters were got via the grid search. For this method, we denote the selected attribute set as sA_{PES}^* and the selected attribute interaction set as sI_{PES}^* .
- **EFM-ES:** the proposed method with ES loss function and the ES based feature selection. Hyperparameters are set by the grid search. Similarly, we denote the selected attribute set as sA_{ES}^* and the selected attribute interaction set as sI_{ES}^* .
- **logFM-PES:** this method first performs a log transform of training data $\mathcal{D}^T = \{(x^{is}, d_{is}) | (i, s) \in T\}$ into $\mathcal{D}_{\text{LOG}}^T = \{(x^{is}, \log(d_{is})) | (i, s) \in T\}$. Then $\mathcal{D}_{\text{LOG}}^T$ is used to train the original FM model with attribute set sA_{PES}^* , attribute interaction set sI_{PES}^* and PES loss by the ALS method [30]. The fitted model $f^{\text{logFM}}(\cdot, \Theta)$ is used to forecast for the test set $\mathcal{D}^E = \{(x^{is}, d_{is}) | (i, s) \in E\}$, and the final prediction is computed as $\hat{d}_{is}^{\text{logFM}} = \exp(f^{\text{logFM}}(x^{is}, \Theta))$ for any $(i, s) \in E$ which is used to compute performance indicators of MAPE and MAE.
- **logFM-ES:** this method is similar to the above logFM-PES with two differences: (1) the loss function is ES and (2) attributes and attribute interactions are sA_{ES}^* and sI_{ES}^* , respectively, by the EFM-ES method.

6.1 Retail Sales Forecasting

6.1.1 Database Introduction and Preprocessing. The database is provided by our industry partner, the retailer of ladies footwear in Singapore. The raw data mainly consists of two parts: (1) attributes associated with each SKU and (2) sales transactions from the point-of-sales (POS) system. We first merge these two sources into one set (table). Consequently, there are 45 explanatory attributes (variables) and 1 response variable. The raw data is preprocessed and they include removal of errors and outliers, filling missing values, data anonymization without losing its nature, consolidation as well as discretization. The first three operations follow the standard data mining techniques while the consolidation refers to aggregating fine-grained information of each sales transaction over the first eleven weeks after launch. The last preprocessing is data discretization. In the proposed approach, we only take categorical attributes into account. We employ an unsupervised data discretization of equal frequencies and domain knowledge to transform the numerical (continuous) attributes into categorical attributes. The off-the-shelf function *discretization* in *infotheo* [26] package in R is used here.

Our database covers the sales occurring in retailer's stores in Singapore for the time period from 1 January 2012 to 20 July 2014. We divide the data into training and test sets. Observations associated with launch time from 1 January 2012 to 14 April 2013 are constructed as training set \mathcal{D}^T . And those with launch time from 31 December 2013 to 3 May 2014 are composed as test set \mathcal{D}^E . The rationale is illustrated by Figure 3. We take both lead time of six month and product life time of 11 weeks into account. We apply the proposed solution on three classes with the masked names of 69-y6, p-1 and x-9w, respectively. Table 1 lists sizes of training and test for the three classes.

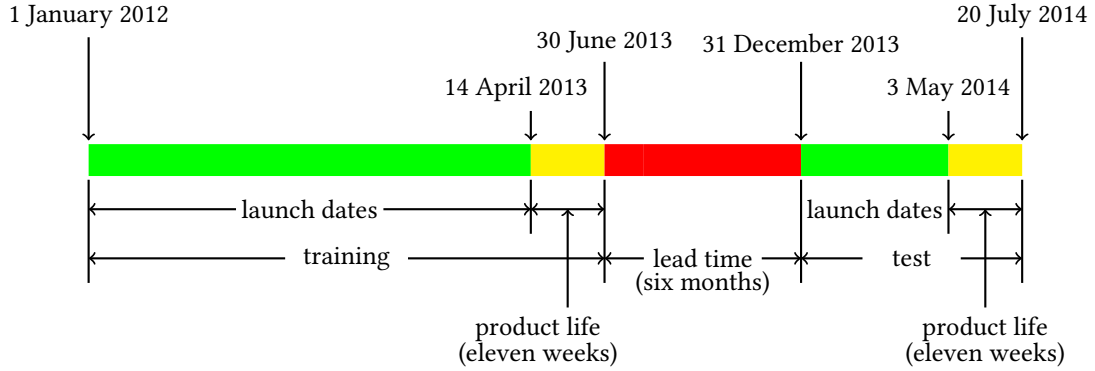


Fig. 3. Partition of data into training and test sets

Table 1. Sizes of training and test for 3 classes

Class (masked name)	Training \mathcal{D}^T	Test \mathcal{D}^E
69-y6	4,984	1,179
p-1	6,222	1,824
x-9w	5,526	1,481

6.1.2 *Methods and Settings.* Besides EFM-PES, EFM-ES, logFM-PES and logFM-ES, we also applied the following methods for current sales forecasting application.

- **SP-EFM-PES:** this method takes the sales quantity of similar SKUs (products) sold at the same store into account with an extended EFM model which is formulated as follows.

$$\begin{aligned} \hat{d}_{is}^E = f^E(x^{is}; \Theta) = & \exp(\beta_0 + \sum_{c \in A} \sum_{j \in L^c} \beta_{cj} x_{cj}^{is} + \sum_{b \in B} \beta_b z_b^{is} + \sum_{(c, c') \in I} \sum_{j \in L^c} \sum_{j' \in L^{c'}} x_{cj}^{is} x_{c'j'}^{is} \langle \mu_{cj}, \mu_{c'j'} \rangle \\ & + \sum_{c \in A} \sum_{j \in L^c} \sum_{b \in B} x_{cj}^{is} z_b^{is} \langle \gamma_{cj}, \gamma_b \rangle + \sum_{b \in B} \sum_{b' < b, b' \in B} z_b^{is} z_{b'}^{is} \langle \zeta_b, \zeta_{b'} \rangle) \end{aligned} \quad (38)$$

Here $\langle \gamma_{cj}, \gamma_b \rangle = \sum_{p=1}^r \gamma_{cj,p} \gamma_{b,p}$ models the intersection between continuous variable z_b^{is} and categorical variable x_{cj}^{is} ; $\langle \zeta_b, \zeta_{b'} \rangle = \sum_{p=1}^t \zeta_{b,p} \zeta_{b',p}$ considers the interactions between two continuous variables z_b^{is} and $z_{b'}^{is}$; B is the index set of continuous variables. We first define the similarity between two SKUs before presenting the detail of B and z_b^{is} . For SKU i and SKU i' at store s , the similarity between them is

$$\text{Similarity}_{i,i'}^s = \frac{\sum_{c \in \{sA_{PES}^* \cup A^{sI_{PES}^*}\}} \sum_{j \in L^c} \mathbb{1}(x_{cj}^{is} == x_{cj}^{i's})}{|sA_{PES}^* \cup A^{sI_{PES}^*}|}. \quad (39)$$

Where $A^{sI_{PES}^*}$ is the attribute set associated with intersection sI_{PES}^* which is $A^{sI_{PES}^*} = \{c | \exists c' \in A \text{ such that } (c, c') \in sI_{PES}^* \text{ or } (c', c) \in sI_{PES}^*\}$, and $\mathbb{1}(x) = 1$ if x is True; 0 otherwise. These continuous variables are sales quantities of similar SKUs which have been sold for eleven weeks at the same store, specifically, for $b \in B = \{1, \dots, |B|\}$, z_b^{is} is the sales quantity of the top b similar SKU in terms of the defined similarity to SKU i at store s . This extended model follows the linearity and the proposed learning method is also applicable here. The loss function is the PES. The input attribute and attribute intersection set for this model are sA_{PES}^* and sI_{PES}^* .

respectively, by the EFM-PES method. The set B is set to be top three similar sales quantities. Interactions between continuous and categorical variables are selected by applying similar feature selection method with the same idea of the proposed approach in Section 5.1 for categorical interactions. We finally note that this method is named as, SP-EFP-PES, where “SP” is short for “similar products”.

- **SP-EFM-ES**: this method is similar to the above SP-EFM-PES with two differences: (1) the loss function is ES and (2) the input attributes and attribute interactions are changed to sA_{ES}^* and sI_{ES}^* , respectively, by the EFM-ES method.
- **Lasso**: off-the-shelf model, `glmnet()`, from R package `glmnet`. The input features are all available attributes and attribute interactions, and the function’s own feature selection mechanism is applied.
- **Random Forest (RF)**: off-the-shelf model, `h2o.randForest()`, from R package `h2o`. The model selects features automatically and all available features are fed.
- **Regression Tree (RT)**: off-the-shelf model, `rpart()`, from R package `rpart`. The features are selected by the function from all available features.
- **Support vector regression (SVR)**: off-the-shelf model, `svm()`, from R package `e1071`. All available features are fed into the function and it selects useful features first.

Table 2. Sales forecasting results of MAPE and MAE for the three classes

	SKU-Store						SKU-Chain					
	MAPE			MAE			MAPE			MAE		
	69-y6	p-1	x-9w	69-y6	p-1	x-9w	69-y6	p-1	x-9w	69-y6	p-1	x-9w
EFM-PES (proposed)	4.55%	5.55%	5.42%	1.59	3.32	3.44	4.57%	5.55%	5.57%	3.81	6.37	6.91
EFM-ES (proposed)	5.15%	7.97%	8.14%	2.49	3.01	3.22	5.32%	8.23%	9.25%	3.61	5.41	5.25
logFM-PES	6.90%	8.00%	10.00%	3.25	4.5	5.5	6.50%	8.52%	11.25%	4.25	8.5	7.5
logFM-ES	7.20%	8.50%	12.00%	2.48	3.4	4.9	7.02%	8.96%	13.00%	3.91	7.2	6.5
SP-EFM-PES	6.50%	8.36%	8.90%	3.62	5.5	5.83	8.20%	10.25%	8.80%	5.06	8.8	7.56
SP-EFM-E-ES	7.23%	9.58%	7.90%	2.3	4.8	4.59	9.50%	11.50%	9.58%	4.05	7.8	6.59
Lasso	85.99%	91.27%	90.24%	28.01	49.69	53.92	86.40%	91.77%	90.63%	67.95	98.2	111.52
RF	116.43%	30.18%	114.49%	31.6	14.26	47.82	109.67%	29.28%	111.01%	76.65	26.69	98.24
RT	142.82%	72.04%	121.04%	39.03	31.54	51.05	135.05%	96.51%	117.29%	94.61	60.69	104.16
SVR	89.77%	47.75%	102.60%	24.56	20.35	43.97	96.40%	42.84%	96.24%	59.58	39.22	90.42

6.1.3 Results. In this subsection, we present forecasting results for the three classes of ladies’ footwear: 69-y6, p-1 and x-9w. Table 2 displays the results, from which we can conclude that (1) the EFM-PES and EFM-ES performs best for most cases in terms of MAPE and MAE, respectively; (2) the feature selection method proposed in current study is effective because other feature selection methods (i.e., Lasso, RF, RT, SVR) result in unacceptable performance; (3) The logFM-PES and logFM-ES cannot perform better than EFM-PES and EFM-ES. We also note that our results of EFM-PES and EFM-ES dominate results of retail literature. The MAPEs for SKU-chain forecasting are around 5% and MAEs are less than 7 units for either each individual class or overall of three classes. It compares favorably with existing studies of sales/demand forecasting [10, 12, 13, 22] among which the best reported performance in terms of MAPE is 16.2% for a single new SKU by Fisher and Vaidyanathan (2014) [13]. For SKU-store forecasts, the overall performance is promising as well.

6.2 Student Performance and Burned Area of Fires Forecasting

6.2.1 Database Introduction and Preprocessing. Here we further apply the proposed methods on two external public datasets and investigate the differences between EFM and logFM, PES and ES loss functions. These two public datasets are:

- **Secondary school student performance (SCTP)** [9]: this dataset is public available at UIC Machine Learning Repository³ and Kaggle⁴, and includes social, demographic and school related attributes and Mathematics and Portuguese language grades of three periods: G1, G2 and G3 (the final grade). There are 15 continuous attributes and 17 categorical attributes; and data sizes are 395 and 649 for Mathematics (SCTP-M) and Portuguese (SCTP-P), respectively. Here we use this dataset to predict G3 the range of which is from 0 to 20. In order to compute MAPE, zeros of G3s are revised to 0.1 in data preprocessing step.
- **Forest fires (FF)** [8]: this dataset is also public available at UIC Machine Learning Repository⁵ and Kaggle⁶. It is on forest fires including meteorological information and can be used to predict the burned area. There are 12 explanatory attributes including 2 categorical and 10 numeric. And there are 517 instances. In order to compute MAPE, zero values of response variable, *burned area*, are preprocessed as 0.1.

Table 3. Hyper-parameters for FM based methods for SCTP and FF Dataset

Loss function	Hyper-parameter	Dataset		
		SCTP-P	SCTP-M	FF
PES	learning rate η	4.95×10^{-6}	3.5×10^{-6}	1.95×10^{-6}
	maxInteractions	4000	4000	15000
	regularization for all attribute selection $\lambda_{A,PES}$	0.005	1.0×10^{-3}	5.0×10^{-4}
	regularization for all interaction selection $\lambda_{I,PES}$	0.10	1.0×10^{-3}	5.0×10^{-4}
	regularization for all attributes λ_{ν}	1×10^{-3}	0.1	0.1
	regularization for all interactions λ_w	10	0	0
	standard deviation σ (initialization)	0.1	0.1	0.1
	attribute selection depth b (categorical and numeric)	3	3	2
	interaction selection depth g (categorical and numeric)	2	2	1
	factorization dimensionalities f, r, t (categorical and numeric)	2	2	2
ES	learning rate η	4.80×10^{-10}	3.15×10^{-10}	2.05×10^{-10}
	maxInteractions	5000	10000	17000
	regularization for all attribute selection $\lambda_{A,ES}$	1000	100	100
	regularization for all interaction selection $\lambda_{I,ES}$	1000	100	500
	regularization for all attributes λ_{ν}	100	0	0
	regularization for all interactions λ_w	0	0	0
	standard deviation σ (initialization)	0.1	0.1	0.1
	attribute selection depth b (categorical and numeric)	3	3	2
	interaction selection depth g (categorical and numeric)	2	2	1
	factorization dimensionalities f, r, t (categorical and numeric)	2	2	2

6.2.2 Methods and Settings. We apply EFM-PES, EFM-ES, logFM-PES, logFM-ES, SVR and RF. Here EFM-PES and EFM-ES are the extended EFM model (38) with numeric variables. The evaluation metrics used are MAPE and MAE. In order to measure the predictive power of these methods, five-fold cross-validation is applied here; each dataset is randomly divided into five parts of equal size; and then one subset is used as test and the remaining is used as training; five runs are performed. Finally, average MAPE and average MAE of test set over the five runs are computed as the results. Now we present settings of hyper-parameters and features. The settings of SVR and RF follow the best results reported by Cortez et al. [9] and Cortez et al. [8] for SCTP and FF datasets, respectively.

³<https://archive.ics.uci.edu/ml/datasets/student+performance>

⁴<https://www.kaggle.com/dipam7/student-grade-prediction>

⁵<https://archive.ics.uci.edu/ml/datasets/Forest+Fires>

⁶<https://www.kaggle.com/elikplim/forest-fires-data-set>

For the remaining FM based methods (i.e., EFM-PES, EFM-ES, logFM-PES and logFM-ES), features are selected by an extended approach from Section 5.1 with the capability of selecting continuous variables; the setting between exponential formulations (i.e., EFM-PES and EFM-ES) and log-transformations (i.e., logFM-PES and logFM-ES) are similar to that of sales forecasting in Section 6.1. Hyper-parameters are given in Table 3 and most of them (i.e., learning rate, maxInteractions, regularization parameters) are got by grid search.

6.2.3 Results. Table 4 displays the results and we have several observations: (1) the performance of FM based methods (i.e., EFM-PES, EFM-ES, logFM-PES and logFM-ES) are better than that of SVR and RF; (2) the EFM-PES and EFM-ES perform best among all methods for SCTP dataset in terms of average MAPE and average MAE, respectively, while the logFM-PES and log-ES dominate all other methods for FF dataset in terms of average MAPE and average MAE, respectively.

Table 4. Forecasting test results of average MAPE and MAE over five runs for SCTP and FF datasets

	average MAPE			average MAE		
	SCTP-P	SCTP-M	FF	SCTP-P	SCTP-M	FF
EFM-PES (proposed)	13.5%	11.32%	35.2%	2.05	3.0	17.00
EFM-ES (proposed)	14.26%	13.51%	37.5%	1.02	1.05	16.25
logFM-PES	15.38%	15.24%	34.3%	2.50	3.6	15.25
logFM-ES	16.51%	17.5%	35.5%	2.05	3.3	14.25
SVR	21%	11.5%	39.5%	2.7	3.5	16.85
RF	17.5%	11.5%	41.5%	1.89	2.05	17.25

6.3 Discussions

6.3.1 Log-transformation v.s. Exponential Formulation. Both log-transformation and exponential formation are two effective approaches for positive response variables. Our studies show that the proposed exponential formulation methods (EFM-PES and EFM-ES) dominate the log-transformation based methods (logFM-PES and logFM-ES) on sales forecasting and SCTP datasets, however, for FF dataset, the log-transformation methods perform better than the exponential formulation based methods. Intuitively, log-transformation projects the response variable in another space and changes the variance. We doubt that the performance differences might be explained by the distribution of response variable. We provide the empirical distribution of response variables in Table 5; for each dataset, the range of response variable values in both training and test was divided into four equal width intervals; the percentage of how many values falling into each interval was computed. From this table, we observe that (1) the distribution of Retail Sales are approximately to the uniform distribution; (2) the Portugueses and Math G3 Grades follow approximate normal distributions; (3) the Burned Area in the FF dataset is significantly different from the others, and highly right-skewed (or positively skewed) distributed. Therefore, we suspect that log-transformation is more effective for datasets with responses following highly right-skewed distributions. For positive response variables which are not right-skewed distributed, using exponential formulation as the model is better than using the log-transformation as the data preprocessing step which might changes the variance of response variables resulting in insufficient training.

6.3.2 PES Loss v.s. ES Loss. Now we investigate the difference between PES and ES minimization for the proposed EFM model on all these datasets. By results over test sets reported at Table 2 and Table 4, EFM-PES (or EFM-ES) performs better than EFM-ES (or EFM-PES) in terms of MAPE (or MAE) for most cases (except for 69-y6 where

Table 5. Empirical distribution of response variables of both training and test in all studied datasets

Response Variable and DataSet	Interval			
	[Min, 0.25Max]	[0.25Max, 0.5Max]	[0.5Max, 0.75Max]	[0.75Max, Max]
Retail Sales in 69-y6	20.58%	29.35%	28.53%	21.54%
Retail Sales in p-1	28.11%	22.59%	23.25%	26.05%
Retail Sales in x-9w	23.60%	27.52%	26.85%	22.03%
Portuguese G3 Grade in SCTP	2.46%	12.94%	64.41%	20.19%
Math G3 Grade in SCTP	11.64%	35.44%	42.78%	10.12%
Burned Area in FF	99.43%	0.19%	0.19%	0.19%

EFM-PES performs better than EFM-ES even in terms of MAE). This results support their effectiveness. Now we analyze the difference in training process between these two loss functions. By Theorem 4.1, the difference is proved to be related to the ratio of maximum response square, Y_{\max}^2 , to minimum response square, Y_{\min}^2 , of the training set. (Here we use Y to denote the general response variable; it refers to sales quantity, G3 grade and burned area for sales, SCTP and FF datasets, respectively). We have also claimed that PES loss function tends to train model to underestimate data. Now we provide all related information in Table 6 where training measurements of mean error squares (MES), mean percentage error squares (MPES) and underestimation ratio are provided. Note that, the underestimation ratio is defined as the fraction of data points fitted (trained) of which are less than the ground-truth (actual); for sales forecasting datasets, here training measurements are for data estimation step only (Step 3 in Algorithm 5) and not for the feature selection; for SCTP and FF datasets, five runs of training are performed and here average measurements are presented. From Table 6, we find that EFM-PES minimizes the MPES while EFM-ES minimizes the MES which follows the definition; the underestimation ratio of the EFM-PES increases with the $\frac{Y_{\max}^2}{Y_{\min}^2}$ while that of the EFM-ES is around 50% for all datasets. This confirms Theorem 4.1 and that PES trains model for underestimation.

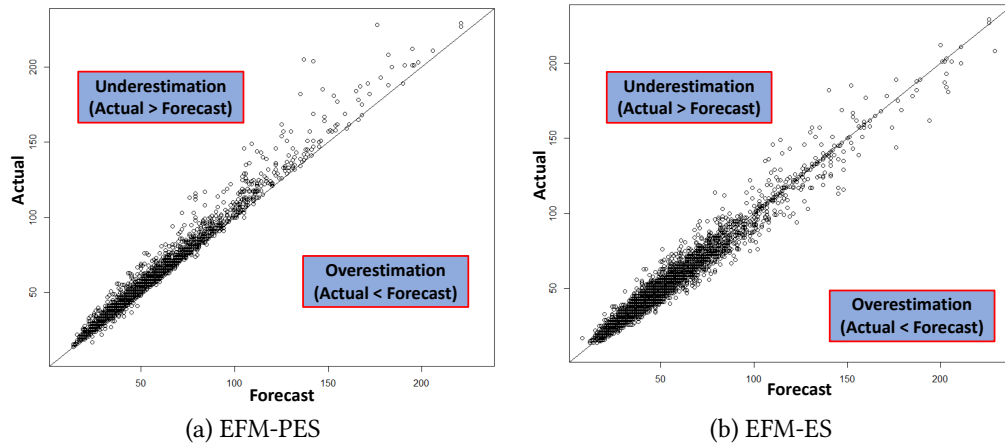


Fig. 4. SKU-store actual sales versus forecast of EFM-PES and EFM-ES for test observations of all three classes

Now, we further provide visualization for SKU-store sales forecasting test results of the three classes at Figure 4 where, for each test instance $(i, s) \in E$, the point (\hat{d}_{is}, d_{is}) is plotted in the Forecast-Actual coordination system and compared with the line Forecast = Actual; points lying exactly on the solid line Actual = Forecast represents

Table 6. Training measurements for PES and ES loss functions for all datasets

Dataset	Training measurement	Method		$\frac{Y_{\max}^2}{Y_{\min}^2}$
		EFM-PES	EFM-ES	
69-y6	mean error squares (MES)	2.5×10^2	1.1×10^2	4.0×10^2
	mean percentage error squares (MPES)	5.2×10^{-5}	6.6×10^{-5}	
	underestimation ratio	0.53	0.42	
p-1	mean error squares (MES)	4.1×10^2	3.8×10^2	1.3×10^3
	mean percentage error squares (MPES)	8.3×10^{-5}	9.7×10^{-5}	
	underestimation ratio	0.64	0.48	
x-9w	mean error squares (MES)	8.5×10^2	7.1×10^2	2.3×10^3
	mean percentage error squares (MPES)	1.2×10^{-4}	3.9×10^{-5}	
	underestimation ratio	0.62	0.46	
SCTP-P	average mean error squares (avg. MES)	2.03×10^2	1.03×10^2	3.61×10^4
	average mean percentage error squares (avg. MPES)	6.2×10^{-4}	9.8×10^{-4}	
	average underestimation ratio	0.85	0.52	
SCTP-M	average mean error squares (avg. MES)	3.04×10^2	2.02×10^2	4×10^4
	average mean percentage error squares (avg. MPES)	5.7×10^{-4}	8.5×10^{-4}	
	underestimation ratio	0.82	0.46	
FF	average mean error squares (avg. MES)	6.5×10^3	4.3×10^3	1.09×10^9
	average mean percentage error squares (avg. MPES)	8.9×10^{-3}	4.5×10^{-2}	
	underestimation ratio	0.93	0.45	

a perfect prediction; points over the line mean underestimation (Actual > Forecast); points under the line denote overestimation (Actual < Forecast). From this Figure, we again observe that the EFM-PES tends to underestimate data while the EFM-ES does not; another interesting observation is that, overall, forecasting accuracy of EFM-PES is better than that of EFM-ES for instances with small actual responses d_{is} , while, for instances with large actual responses d_{is} , the EFM-ES performs better. This is also explainable: training errors of instances with large response values contributes larger for the ES loss than that of instances with small response values; on the other hand, the PES loss increases more heavily with training errors of instances with small response values than instances with large response values.

In conclusion, the above analysis evidences the propositions: (1) the EFM-PES model poses an important property of favouring with underestimation, and (2) PES loss function is dominated by instances with small response values while the ES loss function focuses on instances with large response values.

6.4 Instance Normalization for PES Minimization of Linear Models

In closing this section, we propose a new simple yet effective data normalization for PES minimization for linear models which is motivated by the popular and classic demand-price relationship for pricing and revenue management in supply chain management and marketing area [32]. The linear model $d = \beta_0 + \beta_1 x$ is assumed with two constraints: (1) $d \geq 0$ and (2) $d(x)$ is strictly decreasing in x , $\beta_1 < 0$. Without loss of generality, we set β_0 and β_1 to be 1200 and -10, respectively. A synthetic data set of 100 points $\{(x_i, d_i) | i = 1, \dots, 100\}$ are created: x_i is sampled uniformly on $[1, 50]$ while $d_i = 1200 - 10x_i + \epsilon_i$ and noise ϵ_i follows Gaussian distribution with mean 0 and standard deviation σ , $\epsilon_i \sim \mathcal{N}(0, \sigma)$. The variance of data points is subject to σ .

We minimize the ES and PES loss function to estimate parameters of the linear model and study the differences between them. For the ES minimization, there exists an analytical solution and parameters can be directly derived

by least square (LS). For the PES minimization, note that $\sum_{i=1}^{100} (\frac{\hat{d}_i - d_i}{d_i})^2$ is equivalent to $\sum_{i=1}^{100} (\frac{\hat{d}_i}{d_i} - 1)^2$. We propose the least percentage square (LPS) method for parameter estimation for PES minimization and it consists of two steps as follows.

- (1) Normalization: $\{(x_i, d_i) | i = 1, 2, \dots, 100\}$ is normalized to $\{(x_i/d_i, 1) | i = 1, 2, \dots, 100\}$
- (2) Apply least square solution on the normalized data set.

This normalization, “*instance/sample/row normalization*”, unifies the response to be 1 and divides explanatory variables of each data point by its by the response. It is a new normalization technique and different from popular and classic normalization methods in data mining literature [18] where data are usually normalized per feature/column/variable.

Now we demonstrate the effectiveness and differences of the LPS and LS methods on two synthetic data sets of $\sigma = 10$ and $\sigma = 200$. Figure 5 displays the fit results. In Figure 5 (a), the standard deviation σ is 10 and $\frac{d_{\max}^2}{d_{\min}^2}$ is 3; no big difference between LS and LPS regression lines are observed. But if the standard deviation σ is increased to 200 and $\frac{d_{\max}^2}{d_{\min}^2}$ becomes 60.7 as shown in Figure 5 (b), there is a big difference between LPS line which is dotted in blue and the LS line which is solid red; and clearly the blue LPS regression line underestimate the overall data. By Theorem 4.1, the indicator $\frac{d_{\max}^2}{d_{\min}^2}$ can effectively project the difference between LS and LPS regression. To investigate it, we generate 200 data sets with σ varying from 1 to 200 with increase of 1, and compute MES, MPES, underestimation ratio and $\frac{d_{\max}^2}{d_{\min}^2}$. Figure 6 displays the results. Obviously, the gap of MES, MPES and underestimation ratio between LS and LPS increases with indicator $\frac{d_{\max}^2}{d_{\min}^2}$. This visualization result demonstrates that $\frac{d_{\max}^2}{d_{\min}^2}$ can effectively indicate the difference between LS and LPS regression for linear models.

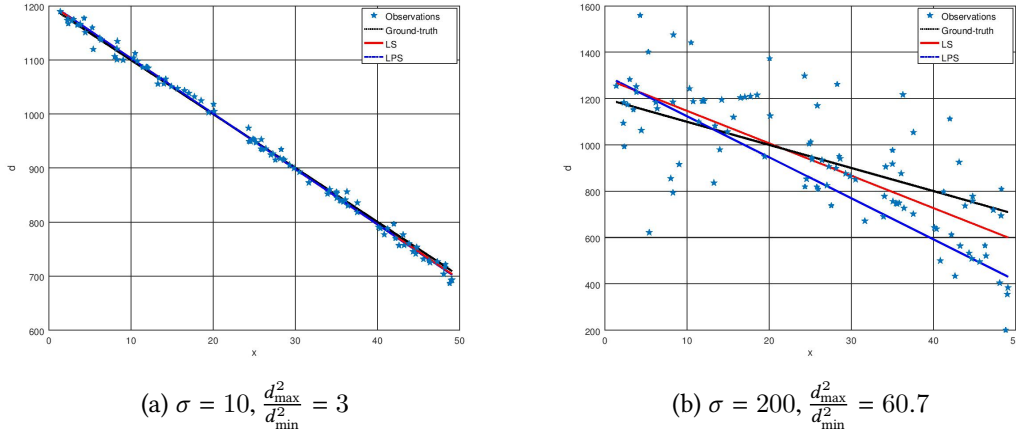


Fig. 5. LS and LPS regression on two data sets.

7 CONCLUSIONS

In this paper we study the sales forecasting problem for new products with long lead time but short product life cycle. The challenge of this problem is two-fold: (1) long lead time and short product life cycle and (2) no historical sales data for the new items. An EFM sales forecast model is developed to address it by taking attributes and

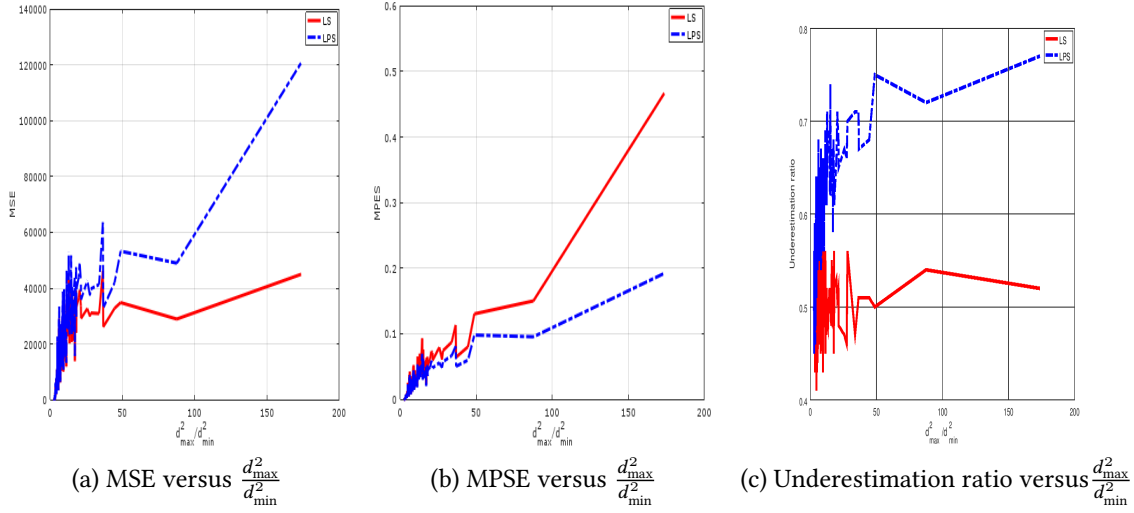


Fig. 6. MSE, MPSE, Underestimation Ratio versus $\frac{d_{\max}^2}{d_{\min}^2}$

their interactions into account. In order to estimate parameters, we minimize two loss functions of percentage error squares (PES) and error squares (ES) over the training set. A real-world data set provided by a footwear retailer in Singapore and two public datasets are used to test the proposed approach and the effectiveness has been demonstrated. Several important takeaways from this study are provided below.

The first takeaway is that sales forecasting for new items should take attribute interactions into account. Modelling sales by attributes is not a new idea in marketing literature. Yet, few studies take this into account. For instance, although Fisher and Vaidyanathan [13] took interactions into account in the implementation, their model somehow ignore them. Our study considers not only attributes but also their interactions.

Another takeaway from this study is the identified differences between PES and ES minimization for parameter estimation. Minimizing the PES for parameter estimation of regression models results in underestimation which fits the situation where unit-holding cost is much greater than unit-shortage cost (e.g. perishable products). Moreover, minimizing the PES for linear regression models can be easily solved by a “*instance/sample/data point/row*” normalization technique with classic least squares solution which differs from the popular “*feature/column/variable*” normalization.

The last takeaway is that exponential formulation is better than conventional log-transformation for modeling positive but not skewed distributed response variables, while log-transformation may be employed when the positive responses follow highly right-skewed distributions.

Last but not least, we close this section by pointing out avenues for future research. Integrating replacement behavior into the EFM mode for substitutable products is a potential extension. This allows replacement behavior considerations to happen not only between attribute levels but also between the interactions. Another potential extension is the investigation of probabilistic explanation for the EFM model and theoretically quantifying the training gap between exponential formulation and the log-transformation of the distributional skewness of positive response variables.

ACKNOWLEDGEMENTS

The authors thank three anonymous referees for their constructive comments. This work is supported by NRF Singapore [Grant NRF-RSS2016-004], MOE-AcrRF-Tier 1 [Grants R-266-000-096-133, R-266-000-096-731, R-266-000-100-646 and R-266-000-119-133], and National Natural Science Foundation of China (Nos. 71801124).

REFERENCES

- [1] Frederick H Abernathy, John T Dunlop, Janice H Hammond, and David Weil. 1999. *A stitch in time: Lean retailing and the transformation of manufacturing—lessons from the apparel and textile industries*. Oxford University Press.
- [2] Felipe Caro and Victor Martínez-de Albéniz. 2015. Fast fashion: Business model overview and research opportunities. In *Retail supply chain management*. Springer, 237–264.
- [3] Chen Chen, Wu Dongxing, Hou Chunyan, and Yuan Xiaojie. 2014. Exploiting social media for stock market prediction with factorization machine. In *Proceedings of the 2014 IEEE/WIC/ACM International Joint Conferences on Web Intelligence (WI) and Intelligent Agent Technologies (IAT)-Volume 02*. IEEE Computer Society, 142–149.
- [4] Tianqi Chen, Hang Li, Qiang Yang, and Yong Yu. 2013. General Functional Matrix Factorization Using Gradient Boosting. In *ICML (1)*. 436–444.
- [5] Chen Cheng, Fen Xia, Tong Zhang, Irwin King, and Michael R Lyu. 2014. Gradient boosting factorization machines. In *Proceedings of the 8th ACM Conference on Recommender systems*. ACM, 265–272.
- [6] Chern Ching-Chin, Ao Ieong Ka Ieng, Wu Ling-Ling, and Kung Ling-Chieh. 2010. Designing a decision-support system for new product sales forecasting. *Expert Systems with applications* 37, 2 (2010), 1654–1665.
- [7] Casey Chung, Shun-Chen Niu, and Chelliah Sriskandarajah. 2012. A Sales Forecast Model for Short-Life-Cycle Products: New Releases at Blockbuster. *Production and Operations Management* 21, 5 (2012), 851–873.
- [8] Paulo Cortez and Anibal de Jesus Raimundo Morais. 2007. A data mining approach to predict forest fires using meteorological data. (2007).
- [9] Paulo Cortez and Alice Maria Gonçalves Silva. 2008. Using data mining to predict secondary school student performance. (2008).
- [10] Peter S Fader and Bruce GS Hardie. 1996. Modeling consumer choice among SKUs. *Journal of marketing Research* (1996), 442–452.
- [11] Peter S Fader and Bruce GS Hardie. 2005. The value of simple models in new product forecasting and customer-base analysis. *Applied Stochastic models in business and industry* 21, 4-5 (2005), 461–473.
- [12] Kris Johnson Ferreira, Bin Hong Alex Lee, and David Simchi-Levi. 2015. Analytics for an online retailer: Demand forecasting and price optimization. *Manufacturing & Service Operations Management* 18, 1 (2015), 69–88.
- [13] Marshall Fisher and Ramnath Vaidyanathan. 2014. A demand estimation procedure for retail assortment optimization with results from implementations. *Management Science* 60, 10 (2014), 2401–2415.
- [14] Christoph Freudenthaler, Lars Schmidt-Thieme, and Steffen Rendle. 2011. Bayesian factorization machines. (2011).
- [15] Paul Goodwin and Richard Lawton. 1999. On the asymmetry of the symmetric MAPE. *International journal of forecasting* 15, 4 (1999), 405–408.
- [16] Huifeng Guo, Ruiming Tang, Yunming Ye, Zhenguo Li, and Xiuqiang He. 2017. DeepFM: a factorization-machine based neural network for CTR prediction. In *Proceedings of the 26th International Joint Conference on Artificial Intelligence*. 1725–1731.
- [17] Isabelle Guyon and André Elisseeff. 2003. An introduction to variable and feature selection. *Journal of machine learning research* 3, Mar (2003), 1157–1182.
- [18] Jiawei Han, Jian Pei, and Micheline Kamber. 2011. *Data mining: concepts and techniques, Chapter 3*. Elsevier.
- [19] Xiangnan He and Tat-Seng Chua. 2017. Neural factorization machines for sparse predictive analytics. In *Proceedings of the 40th International ACM SIGIR conference on Research and Development in Information Retrieval*. 355–364.
- [20] Rob J Hyndman and Anne B Koehler. 2006. Another look at measures of forecast accuracy. *International journal of forecasting* 22, 4 (2006), 679–688.
- [21] Kenneth B Kahn. 1998. Benchmarking sales forecasting performance measures. *The Journal of Business Forecasting* 17, 4 (1998), 19.
- [22] Jonathan Lee, Peter Boatwright, and Wagner A Kamakura. 2003. A Bayesian model for prelaunch sales forecasting of recorded music. *Management Science* 49, 2 (2003), 179–196.
- [23] Henrik Madsen. 2007. *Time series analysis*. CRC Press.
- [24] Spyros Makridakis. 1993. Accuracy measures: theoretical and practical concerns. *International journal of forecasting* 9, 4 (1993), 527–529.
- [25] John T Mentzer and Mark A Moon. 2004. *Sales forecasting management: a demand management approach*. Sage.
- [26] Patrick E Meyer. 2014. infotheo: Information-Theoretic Measures (2014). URL <http://CRAN.R-project.org/package=infotheo>. *R package version 1*, 1 (2014).
- [27] Steffen Rendle. 2010. Factorization machines. In *Data Mining (ICDM), 2010 IEEE 10th International Conference on*. IEEE, 995–1000.
- [28] Steffen Rendle. 2012. Factorization machines with libfm. *ACM Transactions on Intelligent Systems and Technology (TIST)* 3, 3 (2012), 1–22.

- [29] Steffen Rendle. 2012. Social network and click-through prediction with factorization machines. In *KDD-Cup Workshop*.
- [30] Steffen Rendle, Zeno Gantner, Christoph Freudenthaler, and Lars Schmidt-Thieme. 2011. Fast context-aware recommendations with factorization machines. In *Proceedings of the 34th international ACM SIGIR conference on Research and development in Information Retrieval*. ACM, 635–644.
- [31] Nada R Sanders and Karl B Manrodt. 1994. Forecasting practices in US corporations: survey results. *Interfaces* 24, 2 (1994), 92–100.
- [32] David Simchi-Levi, Xin Chen, and Julien Bramel. 2005. The logic of logistics. *Theory, Algorithms, and Applications for Logistics and Supply Chain Management* (2005).
- [33] Robert J Thomas. 1993. Method and situational factors in sales forecast accuracy. *Journal of Forecasting* 12, 1 (1993), 69–77.
- [34] Vladimir Vapnik. 2013. *The nature of statistical learning theory*. Springer science & business media.
- [35] Jill-Jenn Vie. 2018. Deep factorization machines for knowledge tracing. *arXiv preprint arXiv:1805.00356* (2018).
- [36] Zhenxing Xu, Junyi Zhang, Daoqiang Zhang, and Hanyu Wei. 2019. A New Network Traffic Identification Base on Deep Factorization Machine. In *International Conference on Intelligent Science and Big Data Engineering*. Springer, 209–218.

A PROOF OF THEOREM 3.1

PROOF. We start with proving part (a). (a) First, by definition of the ES minimizer Θ_{ES}^* in equation (10), it is true that:

$$\mathcal{L}^{ES}(\Theta_{ES}^*) \leq \mathcal{L}^{ES}(\Theta_{PES}^*). \quad (40)$$

Now we prove the second part in (a) which is $\mathcal{L}^{ES}(\Theta_{PES}^*) \leq \frac{d_{\max}^2}{d_{\min}^2} \mathcal{L}^{ES}(\Theta_{ES}^*)$. By definition of Θ_{PES}^* in equation (11), we have:

$$\mathcal{L}^{PES}(\Theta_{PES}^*) \leq \mathcal{L}^{PES}(\Theta_{ES}^*). \quad (41)$$

Extend the left part and we get:

$$\begin{aligned} \mathcal{L}^{PES}(\Theta_{PES}^*) &= \sum_{(i,s) \in T} \left[\frac{\hat{d}_{is} - d_{is}}{d_{is}} \right]^2 \\ &= \sum_{(i,s) \in T} \frac{[f(x^{is}; \Theta_{PES}^*) - d_{is}]^2}{d_{is}^2} \\ &\geq \sum_{(i,s) \in T} \frac{[f(x^{is}; \Theta_{PES}^*) - d_{is}]^2}{d_{\max}^2} \\ &= \frac{1}{d_{\max}^2} \sum_{(i,s) \in T} [f(x^{is}; \Theta_{PES}^*) - d_{is}]^2 \\ &= \frac{1}{d_{\max}^2} \mathcal{L}^{ES}(\Theta_{PES}^*), \end{aligned} \quad (42)$$

where the inequality follows from the fact that, for each $(i, s) \in T$, $d_{\max} \geq d_{is} \geq 1$, $\frac{1}{d_{is}^2} \geq \frac{1}{d_{\max}^2} > 0$ and $[f(x^{is}; \Theta_{PES}^*) - d_{is}]^2 \geq 0$; and the first and the last equivalences are from definitions of $l^{PES}(\cdot)$ and $l^{ES}(\cdot)$.

Extend the right part of equation (40) and we obtain:

$$\begin{aligned}
\mathcal{L}^{\text{PES}}(\Theta_{\text{ES}}^*) &= \sum_{(i,s) \in T} \left[\frac{\hat{d}_{is} - d_{is}}{d_{is}} \right]^2 \\
&= \sum_{(i,s) \in T} \frac{[f(x^{is}; \Theta_{\text{ES}}^*) - d_{is}]^2}{d_{is}^2} \\
&\leq \sum_{(i,s) \in T} \frac{[f(x^{is}; \Theta_{\text{ES}}^*) - d_{is}]^2}{d_{\min}^2} \\
&= \frac{1}{d_{\min}^2} \sum_{(i,s) \in T} [f(x^{is}; \Theta_{\text{ES}}^*) - d_{is}]^2 \\
&= \frac{1}{d_{\min}^2} \mathcal{L}^{\text{ES}}(\Theta_{\text{ES}}^*),
\end{aligned} \tag{43}$$

where, similarly, the inequality follows from the fact that, for each $(i, s) \in T$, $d_{is} \geq d_{\min} \geq 1$, $\frac{1}{d_{\min}^2} \geq \frac{1}{d_{is}^2} > 0$ and $[f(x^{is}; \Theta_{\text{ES}}^*) - d_{is}]^2 \geq 0$; and the first and the last equivalences are from definitions of $l^{\text{PES}}(\cdot)$ and $l^{\text{ES}}(\cdot)$.

By equations (41), (42) and (43), we get:

$$\frac{1}{d_{\max}^2} \mathcal{L}^{\text{ES}}(\Theta_{\text{PES}}^*) \leq \mathcal{L}^{\text{PES}}(\Theta_{\text{PES}}^*) \leq \mathcal{L}^{\text{PES}}(\Theta_{\text{ES}}^*) \leq \frac{1}{d_{\min}^2} \mathcal{L}^{\text{ES}}(\Theta_{\text{ES}}^*), \tag{44}$$

which is simplified by $d_{\max}^2 > 0$:

$$\mathcal{L}^{\text{ES}}(\Theta_{\text{PES}}^*) \leq \frac{d_{\max}^2}{d_{\min}^2} \mathcal{L}^{\text{ES}}(\Theta_{\text{ES}}^*). \tag{45}$$

Therefore, part (a) in Theorem 4.1 is proved.

Now we prove part (b). First, by equation (11), because Θ_{PES}^* is the minimizer of $\mathcal{L}^{\text{PES}}(\Theta)$,

$$\mathcal{L}^{\text{PES}}(\Theta_{\text{PES}}^*) \leq \mathcal{L}^{\text{PES}}(\Theta_{\text{ES}}^*). \tag{46}$$

Now we recall equation (40), that is:

$$\sum_{(i,s) \in T} [f(x^{is}; \Theta_{\text{ES}}^*) - d_{is}]^2 \leq \sum_{(i,s) \in T} [f(x^{is}; \Theta_{\text{PES}}^*) - d_{is}]^2. \tag{47}$$

Dividing both sides of the inequality by positive value $\frac{1}{d_{\min}^2}$:

$$\frac{1}{d_{\min}^2} \sum_{(i,s) \in T} [f(x^{is}; \Theta_{\text{ES}}^*) - d_{is}]^2 \leq \frac{1}{d_{\min}^2} \sum_{(i,s) \in T} [f(x^{is}; \Theta_{\text{PES}}^*) - d_{is}]^2. \tag{48}$$

Now we extend left side of the inequality:

$$\begin{aligned}
\frac{1}{d_{\min}^2} \sum_{(i,s) \in T} [f(x^{is}; \Theta_{\text{ES}}^*) - d_{is}]^2 &= \sum_{(i,s) \in T} \frac{[f(x^{is}; \Theta_{\text{ES}}^*) - d_{is}]^2}{d_{\min}^2} \\
&\geq \sum_{(i,s) \in T} \frac{[f(x^{is}; \Theta_{\text{ES}}^*) - d_{is}]^2}{d_{is}^2} \\
&= \mathcal{L}^{\text{PES}}(\Theta_{\text{ES}}^*),
\end{aligned} \tag{49}$$

where the first inequality follows from that the domain of d_{is} is $\mathbb{N}_{>0}$; for each d_{is} , $(i, s) \in T$, $\frac{1}{d_{\min}^2} \geq \frac{1}{d_{is}^2} > 0$ and $[f(x^{is}; \Theta_{\text{ES}}^*) - d_{is}]^2 \geq 0$. The last equality comes from definition of the PES loss function $\mathcal{L}^{\text{PES}}(\cdot)$.

Now we extend the right side of inequality (48):

$$\begin{aligned}
 \frac{1}{d_{\min}^2} \sum_{(i,s) \in T} [f(x^{is}; \Theta_{\text{PES}}^*) - d_{is}]^2 &= \frac{d_{\max}^2}{d_{\min}^2} \frac{1}{d_{\max}^2} \sum_{(i,s) \in T} [f(x^{is}; \Theta_{\text{PES}}^*) - d_{is}]^2 \\
 &= \frac{d_{\max}^2}{d_{\min}^2} \sum_{(i,s) \in T} \frac{[f(x^{is}; \Theta_{\text{PES}}^*) - d_{is}]^2}{d_{\max}^2} \\
 &\leq \frac{d_{\max}^2}{d_{\min}^2} \sum_{(i,s) \in T} \frac{[f(x^{is}; \Theta_{\text{PES}}^*) - d_{is}]^2}{d_{is}^2} \\
 &= \frac{d_{\max}^2}{d_{\min}^2} \mathcal{L}^{\text{PES}}(\Theta_{\text{PES}}^*),
 \end{aligned} \tag{50}$$

where the inequality is due to that, for $(i, s) \in T$, $d_{\max} \geq d_{is} \geq 1$, $\frac{1}{d_{is}^2} \geq \frac{1}{d_{\max}^2} > 0$ and $[f(x^{is}; \Theta_{\text{ES}}^*) - d_{is}]^2 \geq 0$

By inequalities (48), (49) and (50), we can get:

$$\mathcal{L}^{\text{PES}}(\Theta_{\text{ES}}^*) \leq \frac{d_{\max}^2}{d_{\min}^2} \mathcal{L}^{\text{PES}}(\Theta_{\text{PES}}^*), \tag{51}$$

which completes the proof for part (b). \square