

Deep learning in business analytics and operations research: Models, applications and managerial implications

*Submitted to special issue on
“Business Analytics: Defining the field and identifying a research agenda”*

Mathias Kraus^a, Stefan Feuerriegel^a, Asil Oztekin^{b,*}

^a*ETH Zurich, Weinbergstr. 56/58, 8092 Zurich, Switzerland*

^b*Department of Operations & Information Systems, Manning School of Business, University of Massachusetts Lowell,
Lowell, MA 01854 USA*

Abstract

Business analytics refers to methods and practices that create value through data for individuals, firms, and organizations. This field is currently experiencing a radical shift due to the advent of deep learning: deep neural networks promise improvements in prediction performance as compared to models from traditional machine learning. However, our research into the existing body of literature reveals a scarcity of research works utilizing deep learning in our discipline. Accordingly, the objectives of this overview article are as follows: (1) we review research on deep learning for business analytics from an operational point of view. (2) We motivate why researchers and practitioners from business analytics should utilize deep neural networks and review potential use cases, necessary requirements, and benefits. (3) We investigate the added value to operations research in different case studies with real data from entrepreneurial undertakings. All such cases demonstrate improvements in operational performance over traditional machine learning and thus direct value gains. (4) We provide guidelines and implications for researchers, managers and practitioners in operations research who want to advance their capabilities for business analytics with regard to deep learning. (5) Our computational experiments find that default, out-of-the-box architectures are often suboptimal and thus highlight the value of customized architectures by proposing a novel deep-embedded network.

Keywords: Analytics, Deep learning, Deep neural networks, Case studies, Managerial implications, Research agenda

1. Introduction

What was once described by Davenport & Harris (2007) as “*firms compet[ing] on analytics*” is now more true than ever. This development is currently propelled by the surge in “big data”, allowing for an efficient access and analysis of large datasets (Chen et al., 2012; Agarwal & Dhar, 2014; Baensens et al., 2016). Innovations in business analytics have become not only desirable but a key necessity for

*Corresponding author. Technical questions should be directed to Mathias Kraus.

Email addresses: mathiaskraus@ethz.ch (Mathias Kraus), sfeuerriegel@ethz.ch (Stefan Feuerriegel), asil_oztekin@uml.edu (Asil Oztekin)

the successful performance of firms (Mortenson et al., 2015; Lim et al., 2013; Ranyard et al., 2015). This holds true for all areas of business operations. Examples include, for instance, supply chain management (Carbonneau et al., 2008), risk modeling (Lessmann et al., 2015), commerce (Scholz et al., 2017), preventive maintenance (Sun et al., 2009), and manufacturing (Hu et al., 2017). The core component for successfully competing with business analytics is the underlying predictive model. This represents the unit responsible for making the actual forecasts and its accuracy contributes directly to the overall value creation.

With recent advances in machine learning, a specific type of predictive model has received great traction lately: *deep learning* (LeCun et al., 2015). The underlying concept is not specific to machine learning or data-analytics approaches from operations research, as it simply refers to *deep* neural networks. However, what has changed from early experiments with neural networks is the dimension of the networks, which now can easily contain up to hundreds of layers, millions of neurons and complex structures of connections between them (e.g. He et al., 2016). This introduces the unprecedented flexibility to model even highly complex, non-linear relationships between predictor and outcome variables, a quality that has allowed deep neural networks to outperform models from traditional machine learning in a variety of tasks. Figure 1 depicts an illustrative sketch of a deep neural network. It should also be noted that big data (George et al., 2014, 2016) is now prevalent in firms, yet traditional machine learning (e.g., support vector machines, random forests) can often not be applied to such large datasets, whereas the optimization routines in deep learning scale efficiently.

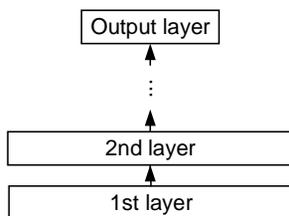


Figure 1: Illustrative deep neural network with four layers. Here the width of a layer corresponds to its dimension, i.e. the number of neurons.

The expected improvement in prediction performance provided by deep learning has led to a selection of showcases.¹ For example, computerized personal assistants, such as Apple’s Siri, Amazon’s Alexa, Google Now or Microsoft’s Cortana, now make heavy use of deep neural networks to recognize, understand and answer human questions.² In this regard, Microsoft unveiled a speech recognition system in 2016 that is capable of transcribing spoken words almost as accurately as professionally trained humans. In October 2016, Google launched an update to its translation system that utilizes deep learning in order to improve translation accuracy, thereby approaching the performance of humans.³ Deep learning

¹ *Harvard Business Review* (2017): “Deep learning will radically change the ways we interact with technology” <https://hbr.org/2017/01/deep-learning-will-radically-change-the-ways-we-interact-with-technology>, accessed February 20, 2018.

² *The Economist* (2017): “Technology quarterly: Finding a voice”. <http://www.economist.com/technology-quarterly/2017-05-01/language>, accessed February 20, 2018.

³ *Washington Post* (2016): “Google Translate is getting really, really accurate”. Available via <https://www.washingtonpost.com/news/innovations/wp/2016/10/03/google-translate-is-getting-really-really-accurate/>,

has not only shown great success in natural language processing but also in image classification, object detection, object localization and image generation. For instance, Alipay introduced a mobile payment app to more than 120 million people in China that allows them to use face recognition for payments. This technology was ranked by *Technology Review* as one of the ten breakthrough technologies of 2017.⁴ Aside from these applications, deep learning has also been successfully applied to recommendation systems. In this regard, both Amazon and Netflix utilize deep neural networks for personalized product recommendations.

Why deep neural networks have only now become so powerful has several explanations (Goodfellow et al., 2017):

1. *Computational power.* Computational capabilities have increased rapidly, especially due to the widespread use of graphics processing units (GPUs). These are particularly suited to executing the operations from linear algebra necessary for fitting neural networks. For instance, Google DeepMind optimized a deep neural network using 176 GPUs for 40 days to beat the best human players in the game Go (Silver et al., 2017). This would have required far greater computational resources without acceleration through GPUs.
2. *Data.* Second, large datasets are needed to train deep neural networks in order to prevent overfitting and fine-tune parameters (see Figure 2). The performance of deep neural networks generally improves with increasing amounts of data; smaller datasets incorporating only several hundred of datapoints were not sufficient to optimize DNNs in previous years. However, in the era of big data, large datasets are now common in most businesses which are important, as empirical results show that even DNNs still benefit from additional data even when already having millions of datapoints (Goodfellow et al., 2017). As a result, big data can be used effectively, often by mining public content from the Web. One prominent example, the so-called ImageNet dataset, was developed to support tasks in computer vision and comprises more than 14 million images. We later detail this aspect by showing how the size of the dataset affects the overall performance of deep learning.
3. *Optimization algorithms.* Optimizing the parameters in deep neural networks is a challenging undertaking. Several optimization algorithms have been proposed since Hinton & Salakhutdinov (2006) published what is now regarded as the seminal work of deep learning. In this publication, the authors increase the depth of neural networks gradually by alternating between adding a new layer and optimizing the network parameters. This technique, which stabilizes the optimization, paved the way for learning deeper networks. Further innovations deal with the optimizer itself. For instance, the large size of datasets prohibits one from optimizing the overall performance directly; deep neural networks are instead trained by stochastic optimization (LeCun et al., 1998). Here parameters are updated based on an approximation of the objective function (i. e. by evaluating it on a subsample of the whole dataset). This distinguishes deep neural networks from many other machine learning models, which are not designed for big data and, hence, do not scale to millions of datapoints. To improve

accessed February 20, 2018.

⁴ *Technology Review* (2017): “10 breakthrough technologies: Paying with your face”. <https://www.technologyreview.com/s/603494/10-breakthrough-technologies-2017-paying-with-your-face/>, last accessed February 20, 2018.

optimization, there are a number of common optimization methods (e. g. Adam, Adagrad, RMSprop) that implement variants of stochastic gradient descent, often paired with an adaptive regulation of the step size. In addition, deep neural networks with millions of parameters suffered from overfitting the model to the training data. As a remedy, it is now recommended that one integrate regularization into the optimization procedure (i. e. weight decay, dropout or batch normalization) in order to improve the generalizability.

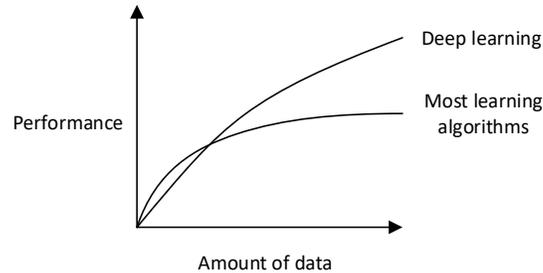


Figure 2: Illustrative comparison between performance of deep learning against that of most other machine learning algorithms. Deep neural networks still benefit from large amounts of data, whereas the performance increase of other machine learning models plateaus (Ng, 2016a)

While deep learning is on the way to becoming the industry standard for predictive analytics within business analytics and operations research, our discipline is still in its infancy with regard to adopting this technology. To support this claim, we conducted an extensive literature review of all papers published by October 2018 across the premier journals in our field, namely, *European Journal of Operational Research*, *Operations Research*, *Production and Operations Management*, *Journal of Operations Management*, and *Decision Sciences*. We specifically searched for papers containing the terms “deep neural network” or “deep learning”. Our search returned 15 matches but a closer inspection led us discard 12 of them, since these only contained the term in passing but without actually using deep learning. The resulting three matches are listed in Table 1. This is an interesting observation as the operations research community has a longstanding tradition of applying neural networks. Examples of 1-hidden layer networks appear in the areas of, for instance, healthcare (Misiunas et al., 2016; Oztekin et al., 2018), demand forecasts (Carbonneau et al., 2008; Venkatesh et al., 2014) and maintenance (Mahamad et al., 2010; Mazhar et al., 2007). We further extended our literature search to other fields of management science, namely, accounting, finance, marketing and information systems; yet these outlets (as per Financial Times 50 list) published a total of only two papers that benefit from deep learning; see Adamopoulos et al. (2018) and Li et al. (2017).

It is not only academia that has yet to fully incorporate deep learning in its decision-making routines. This fact also holds true for a majority of enterprises: a 2016 report on “The Age of Analytics” by the McKinsey Global Institute refers to deep learning as “*the coming wave*” (Henke et al., 2016). This is also borne out by our own expertise when collaborating with large-cap companies, including consulting firms, across Europe and the United States. One member from a top management consulting firm even admitted to us that “*we don’t have a clue how deep learning works, neither do our clients*”. Hence, the above mentioned showcases of deep learning are largely exceptions among a handful of selected firms,

thereby highlighting the dire need for company professionals to better understand deep learning, its applications and value (cf. Lee, 2018).

The objective of this overview is to provide an overview on the use of deep learning in academia and practice from an operations research perspective: (1) we summarize the most relevant mathematical concepts in deep learning for readers who are not familiar with this technique. (2) We demonstrate the use of deep learning across three case studies from operations research. Here we specifically compare its performance to conventional models from machine learning, thereby showing the improvements in both prediction and operational performance. As part of it, we see that default, out-of-the-box architectures are often not sufficient and, following this, we propose a novel deep-embedded network architecture. (3) We provide recommendations regarding which network architecture to choose and how to tune parameters. (4) We derive implications for managers and practitioners who want to engage in the use of deep neural networks for their operations or business analytics. (5) We propose directions for future research, especially with regard to the use of deep learning for business analytics and operations research.

Paper	Problem	Network architecture
Krauss et al. (2017)	Short-term arbitrage trading for the S&P 500	5-layer perceptron
Probst et al. (2017)	Generative neural networks for estimation of distribution algorithms	Stochastic neural network (i. e. restricted Boltzmann machine)
Fischer & Krauss (2018)	Predict directional movement of constituent stocks of the S&P 500	Long short-term memory

Table 1: Literature review for papers from operations research utilizing techniques from deep learning.

The remainder of this overview article is structured as follows. Section 2 recapitulates deep neural networks by rewriting the concept as an optimization problem using conventional terminology of operations research, while Section 3 proposes our tailored deep-embedded network architecture. To demonstrate the potential use of deep learning in practice, we then conduct three experiments with real-world data from actual business problems (see Section 4). All of our experiments reveal superior prediction performance on the part of deep neural networks. Finally, Section 5 provides recommendations for the use of deep learning in operations research and highlights the implications of our work for both research and management. Section 6 concludes with a summary.

2. Mathematical background: From neural networks to deep learning

This section reviews the transformation from “shallow” neural networks to deep learning. For a detailed description, we refer to Russell et al. (2010) and specifically to Goodfellow et al. (2017). In addition, Schmidhuber (2015) presents a chronological review.

2.1. Predictive analytics

Predictive models in business analytics exploit input features $\mathbf{x} \in \mathbb{X} \subseteq \mathbb{R}^m$ to predict an unknown variable $y \in \mathbb{Y}$. In practical settings, the input features could be recent sales figures in order to forecast future production needs, or historic sensor data to anticipate machinery failure. Depending on whether

this is a discrete label ($y \in \{0, \dots, k\}$) or a real value ($y \in \mathbb{R}$), we refer to it as a classification or regression task, respectively. The objective is then to find a mapping $f : \mathbf{x} \mapsto y$.

The choice of such functions is given by the predictive model $f(\cdot; w)$ with additional parameters w ; that is, $y \approx f(\mathbf{x}; w)$. Then, in practice, the objective behind the prediction results in an optimization problem whereby one must find the best parameters w . A variety of models f are common in business analytics: examples from traditional machine learning involve, for instance, linear models (e. g. Bertsimas & King, 2016; Bertsimas & Shioda, 2007), decision trees, support vector machines, neural networks (e. g. Delen et al., 2012; Oztekin et al., 2016), or even deep neural networks as motivated by this work. Each of them is often accompanied by a tailored optimization strategy; see Bertsimas & Kallus (2014).

The above optimization requires a performance measure that assesses the model, i. e. the error between the predicted value $\hat{y} = f(\mathbf{x}; w)$ and the true observation y . This is formalized by a loss function $\mathcal{L} : \mathbb{Y} \times \mathbb{Y} \rightarrow \mathbb{R}$. The actual choice depends on the desired objective (e. g. whether one wants to penalize false-positives or true-negatives) and the prediction task, i. e. classification or regression. Predictive modeling now simplifies to minimizing the loss $\mathcal{L}(\hat{y}, y) = \mathcal{L}(f(\mathbf{x}; w), y)$, usually summed over a set of samples i . Hence, we yield the optimization problem

$$w^* = \arg \min_{w \in \mathbb{W}} \sum_i \mathcal{L}(\hat{y}_i, y_i) = \arg \min_{w \in \mathbb{W}} \sum_i \mathcal{L}(f(\mathbf{x}_i; w), y_i), \quad (1)$$

where \mathbb{W} denotes the weight space. With the increasing dimensionality of the weight space, the predictive model gains flexibility and thus becomes able to adapt to complex relationships between features and outcome. At the same time, a high-dimensional space heightens the computational requirements for solving the optimization problem. In addition, a large number of pairs (\mathbf{x}, y) are required to ensure that there are sufficient samples with each combination of value, due to the curse of dimensionality (Bellman, 1972).

The following sections address three issues: how to define the function f in deep neural networks (as compared to traditional neural networks); how to find the parameters w ; how to prevent overfitting.

2.2. Neural networks

2.2.1. Single-layer neural networks

Similar to a biological neural network, an artificial neural network is a collection of connected units called neurons. An artificial neuron receives inputs from other neurons, computes the weighted sum of the inputs (based on weights w), and maps the sum via an activation function to the output. Figure 3 illustrates how a neuron receives inputs $x_1^{(i)}, \dots, x_m^{(i)}$ and processes them with weights w_1, \dots, w_m .

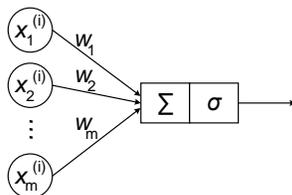


Figure 3: Illustrative neuron in an artificial network with inputs $x_1^{(i)}, \dots, x_m^{(i)}$, weights w_1, \dots, w_m and activation function σ . Σ denotes the summation over all inputs.

Finally, the output is passed to connected neurons as their inputs. If all connections follow a forward structure (i. e. from input neurons that process features of \mathbf{x} to the output neurons), we call this a feedforward neural network. Accordingly, the network is free of cycles or feedback connections that pass information backwards. The latter are called recurrent neural networks and are discussed in Section 2.4.2.

The simplest neural network follows the above concept and is thus represented by a single-layer perceptron where the network f_{1NN} is computed via a linear combination embedded in an activation function σ . We yield

$$f_{\text{1NN}}(\mathbf{x}; W, b) = \sigma(W \mathbf{x} + b) \tag{2}$$

with parameters W (the weight matrix) and b (an intercept called bias). Here the flexibility of the network stems from choosing an activation that is non-linear.

Common choices of activation functions are as follows:

- the sigmoid function $\sigma(x) = \frac{1}{1 + e^{-x}} \in (0, 1)$,
- the hyperbolic tangent $\sigma(x) = \frac{e^{2x} - 1}{e^{2x} + 1} \in [-1, 1]$, and
- the rectified linear unit (ReLU) given by $\sigma(x) = \max(0, x) \in [0, \infty)$.

It is apparent that they all share a particular characteristic, namely, that a certain threshold must be exceeded in order for the activation functions to pass through values. Hence, the idea of activation functions is biologically inspired in the sense that they resemble neurons in the human brain, which also have to receive a certain stimuli in order to be activated (Stachenfeld et al., 2017). In practice, the different choices entail their own (dis-)advantages; for instance, the ReLU activation function can lead to neurons which output zero for essentially all inputs, thereby losing flexibility.

Because of the activation functions, neural networks can model non-linear and non-convex functions. However, this also makes the optimization problem denoted by Equation (1) difficult to optimize. In particular, there exists no direct closed-form solution that gives the global optimum of the optimization problem. Instead, one uses gradient descent to find local optima. This technique has been shown to provide solutions that are close to the global optimum (Choromanska et al., 2015). Mathematically, one updates parameter $w \in \mathbb{W}$ by calculating the partial derivative of w with respect to the loss \mathcal{L} of the training samples, i. e.

$$\sum_i \frac{\delta}{\delta w} \mathcal{L}(f(\mathbf{x}_i, w), y_i). \tag{3}$$

Accordingly, one updates the parameter w such that the loss \mathcal{L} decreases via

$$w \leftarrow w - \eta \sum_i \frac{\delta}{\delta w} \mathcal{L}(f(\mathbf{x}_i, w), y_i), \tag{4}$$

where η denotes the step size or learning rate of the gradient descent optimization. One proceeds similarly to update the bias b .

Single-layer neural networks have many limitations. Most famously, for monotonic activation functions, they cannot learn the XOR function given by $f_{\text{1NN}}([0, 0], w) = 0$, $f_{\text{1NN}}([0, 1], w) = 1$, $f_{\text{1NN}}([1, 0], w) =$

1 and $f_{\text{INN}}([1, 1], w) = 0$. As a remedy, multi-layer neural networks consist of many layers that first transform their input into higher-dimensional representations and then into the output. In theory, the universal approximator theorem guarantees that neural networks with three layers are sufficient to represent arbitrary functions $f : \mathbb{X} \rightarrow Y$ (Cybenko, 1989). Nevertheless, practical experience suggests that deeper models can better reduce the generalization error (Goodfellow et al., 2017).

2.2.2. Deep neural networks

While the previous neural networks consisted of only a single layer, one can extend the mathematical specification to multi-layered perceptrons. We refer to these using the term “deep”, which can reflect an arbitrary number of layers.⁵ In such networks, the number of free parameters increases, as well as the flexibility of the network to represent highly non-linear functions. We can formalize this mathematically by stacking several single-layer networks into a deep neural network with k layers, i. e.

$$f_{\text{DNN}}(\mathbf{x}) = \underbrace{f_{\text{INN}}(f_{\text{INN}}(\dots f_{\text{INN}}(\mathbf{x})))}_k = \underbrace{f_{\text{INN}} \circ \dots \circ f_{\text{INN}}}_k(\mathbf{x}). \quad (5)$$

The first layer is referred to as the input layer, the last as the output layer and the remainder are termed hidden layers.

We note that the dimension of each layer is not necessarily equal across all layers, that is to say it can differ. In practice, the depth of deep neural networks varies across applications, ranging between two hidden layers to potentially several hundred. Together with the corresponding dimension of each layer, this can easily result in networks that entail tens of millions of degrees of freedom. As a result, optimizing the weights in a deep neural network is a daunting task, requiring (1) gradient-based optimization methods and (2) regularization. Both are detailed in Section 2.3.

The task of choosing an adequate number of layers and the number of neurons in each layer represents a challenging undertaking. Montavon et al. (2012) recommend adding layers until the generalization error stops improving. Other best practices suggest adding layers until the predictive model overfits on the training data and then removing the overfitting by regularization methods. Moreover, a large dimension of neurons in each layer is generally preferred, as these seldom interfere with the generalization error. Yet, we note that optimizing these hyperparameters is still subject to active research.

For deep neural networks, the activation function is commonly set to the rectified linear unit (LeCun et al., 2015). This choice leads to sparse settings whereby a large portion of hidden units are not activated, thus having zero output. On the other hand, the recurrent network architectures (cf. Section 2.4.2) are frequently utilized with sigmoid activation functions, since these constrain the output to between 0 and 1. Thereby, so-called gates can be defined which circumvent exploding gradients during the numerical optimization.

In classification tasks, the traditional loss functions, such as L_1 and L_2 , suffer from slow learning. That is, the partial derivatives of the loss function are small as compared to a large value of the loss itself.

⁵Although some architectures currently being used consist of only two to five hidden layers, only recent innovations in the realm of deep learning render it possible to make use of them. Hence, they are also considered as being deep neural networks.

Accordingly, classification tasks generally output a discrete probability distribution over all possible classes by using a softmax activation in the output layer (Russell et al., 2010). However, for large losses, the partial derivative for traditional losses vanishes and, as a consequence, one typically prefers the cross entropy to measure the similarity between the output distribution and the target distribution (Goodfellow et al., 2017). In regression tasks, neural networks output a continuous values, for which one frequently draws upon the mean squared error or the smoothed mean absolute error.

2.3. Model estimation

2.3.1. Weight optimization

The optimization in deep learning is analogous to the general setting in predictive analytics, the loss is minimized via

$$w^* = \arg \min_{w \in \mathbb{W}} \mathcal{L}(f_{\text{DNN}}(\mathbf{x}; w), y). \quad (6)$$

While the weights in a simple perceptron can be identified through convex optimization, the optimization problem for deep neural networks is computationally challenging due to the high number of free parameters. In fact, Judd (1990) proves that it is NP-hard to optimize a neural network so that it produces the correct output for all the training samples. That study also shows that this problem remains NP-hard even if the neural network is only required to produce the correct output for two-thirds of the training examples.

The conventional solution strategy for optimizing deep neural networks involves gradient-based numerical optimization (Saad, 1998). Similar to optimizing single-layer neural networks, one computes the partial derivatives of the parameters with respect to the loss \mathcal{L} and changes the parameters in order to decrease the loss (cf. Equation (4)). However, this must be done for all layers, from the output back to the input layer. Hence, a technique called backpropagation is preferable for reasons of efficiency (Rumelhart et al., 1986). Backpropagation exploits the chain rule to reuse computations it has already calculated for the previous layer. Since these operations involve simple matrix operations from linear algebra, they can be run in parallel, turning GPUs into an effective accelerator for optimizations in deep learning.

The runtime for optimization of a deep neural network can still be very high, since a single update of the parameters requires predictions \hat{y}_i of all samples \mathbf{x}_i . As a remedy, stochastic gradient descent approximates the loss function across all samples with the loss of a smaller subset, called the “minibatch”. The size of the minibatch presents another hyperparameter, for which we typically recommend values between 32 and 256 based on our experience.

The learning rate ν in Equation (4) controls the step size during the optimization process. When this learning rate decreases at an appropriate rate, then stochastic gradient descent is guaranteed to converge to the global optimum under mild mathematical assumptions (Kiwiel, 2001). Hence, the learning rate poses another hyperparameter that is usually chosen by visual inspection of the the learning curve, which traces the loss \mathcal{L} as a function of time. Too high a learning curve coincides with oscillations, whereas a learning rate that is too low results in a slow optimization. To facilitate the choice, one may utilize an early stopping technique, which terminates the optimization when no improvement is achieved on the

validation set for a certain time period.

Even though stochastic gradient descent is popular for optimizing deep neural networks, its performance for training can be very slow when the direction of the gradients changes (similar to a second-order derivative). A technique called momentum helps solve this issue by adding a velocity vector to the gradient (Goodfellow et al., 2017). It accumulates a moving average of past gradients. As a result, momentum follows the drift of past gradients.

The above mentioned concepts have been integrated into tailored optimization routines for deep learning. This has resulted in a variety of optimizers, such as Adam (often considered a baseline), Adagrad, Adadelta and RMSProp, which are common in practice (Goodfellow et al., 2017). The question of which algorithm performs best is still to be further studied.

2.3.2. Regularization

Optimizing deep neural networks typically requires a trade-off: on the one hand, one aims at a high number of free parameters as this allows for the representation of highly non-linear relationships. On the other hand, this makes the network prone to overfitting. The following remedies (sometimes used in combination) are common in deep learning and present forms of regularization to the weights:

1. *Weight decay* adds a regularization term to the loss function, penalizing large weights in the network. With W denoting the set of all parameters, the loss function thus changes to

$$\mathcal{L}_{\text{WD}} = \sum_i \frac{\delta}{\delta w} \mathcal{L}(f(\mathbf{x}_i, w), y_i) + \frac{\lambda}{2} \|W\|_2. \quad (7)$$

Consequently, the gradient descent utilizes a new update rule given by

$$w \leftarrow w - \eta \left(\sum_i \frac{\delta}{\delta w} \mathcal{L}(f(\mathbf{x}_i, w), y_i) + \lambda \|W\|_2 \right). \quad (8)$$

As a result, the decision boundaries become smoother, thereby facilitating generalization of the network (Goodfellow et al., 2017).

2. *Dropout* discards a small but random portion of the neurons during each iteration of training (Srivastava et al., 2014). The underlying intuition is that the several neurons are likely to model the same non-linear relationship simultaneously; dropout thereby prevents neurons from co-adapting to the same features. Mathematically, this can be achieved by setting the corresponding rows in the weight matrix W to zero.
3. *Batch normalization* (not to be mistaken with the previous updating in minibatches) performs a normalization of the output of each layer before forwarding it as input to the next layer (Ioffe & Szegedy, 2015). As a result, this shrinks values closer to zero and, in practice, allows one to use higher learning rates with less care concerning the initialization of parameters.

2.4. Advanced architectures

Beyond the multi-layered perceptron discussed above, an array of alternative architectures have been proposed, often targeting specific data structures. Goodfellow et al. (2017), as well as Schmidhuber

(2015), offer comprehensive overviews, while we summarize the most widely utilized choices in the following: convolutional networks that are common in vision tasks and recurrent networks for handling sequential data (see Table 2). These can also be combined with the previous dense layers as additional building blocks.

Architecture	Data structure	Examples	Tuning parameters
DENSE NEURAL NETWORKS			
Multi-layer perceptron (MLP)	Feature vectors of fixed length	Can be simple replacement of traditional models	Activation function, number of layers and neurons
CONVOLUTIONAL NEURAL NETWORKS			
Convolutional neural network (CNN)	High-dimensional data with local dependencies	Image recognition, speech recognition	Number and width of convolution filters
RECURRENT NEURAL NETWORKS			
Long short-term memory (LSTM)	Sequential data	Text classification, time series forecasting	Number of layers and neurons
Gated recurrent unit (GRU)	Sequential data	Text mining, time series forecasting	Number of layers and neurons

Table 2: Overview of common network architectures in deep learning.

2.4.1. Convolutional neural network

The convolutional neural network (CNN) exploits spatial dependencies in the data, e.g. among neighboring pixels in an image. In other words, the idea of a CNN is to take advantage of a pyramid structure to first identify concepts at the lowest level before passing these concepts to the next layer, which, in turn, create concepts of higher level. Therefore, neurons are no longer connected with every other as a CNN goes from concepts that are identified locally to globally. Rather, neurons are now densely connected only in a small neighborhood (see Figure 4). This choice is motivated by the human visual cortex, which also experiences stimulation in a restricted region of the visual field.

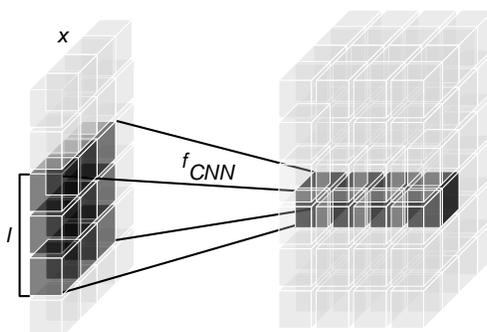


Figure 4: Example of a convolution layer that exploits the spatial structure in the input data (left). It then applies a convolution operation to a local neighborhood in order to compute the output (right). Here l denotes the kernel width.

Mathematically, the CNN applies a convolution operation (also known as kernels or filters) to the input and then passes the result to the next layer. Let \mathbf{x}_{ij} denote the element in row i and column j of a matrix $\mathbf{x} \in \mathbb{R}^{m \times m}$ that represents the values from the input layer. Then, the CNN calculates

$$f_{\text{CNN}}(\mathbf{x})_{ij} = \sum_{x=-\frac{l}{2}}^{\frac{l}{2}} \sum_{y=-\frac{l}{2}}^{\frac{l}{2}} \mathbf{x}_{i+x, j+y} k_{xy}, \quad (9)$$

when convolving \mathbf{x} with a kernel k of width l . Here one usually pads \mathbf{x} with zeros in order to avoid accessing non-existing values beyond the bounds of \mathbf{x} . The above operation can be analogously extended to higher-dimensional tensors.

2.4.2. Recurrent neural networks

Traditional machine learning is limited to input vectors $\mathbf{x} \in \mathbb{R}^m$ of a fixed dimension m . This is rarely suited for sequences, which do not fit into such a structure of fixed size, since they entail varying lengths, ranging between 1 and an arbitrary number of elements. For this reason, sequence learning requires a problem specification whereby $f_{\text{RNN}} : \mathbb{X} \rightarrow \mathbb{Y}$ can handle input from $\mathbb{X} = \{\mathbb{R}, \mathbb{R} \times \mathbb{R}, \mathbb{R}^3, \dots\}$. In other words, this formalization takes sequences $\mathbf{x}_i = [x_1^{(i)}, \dots, x_{\tau_i}^{(i)}]$ as input, but each with its own length τ_i . Prominent examples of sequential data in practical application include time series or natural language (where the running text is represented as a series of words or characters).

Recurrent neural networks are specifically developed to handle such sequential input. They iterate over input sequences \mathbf{x}_i without making assumptions regarding the length of the sequence. Hence, the same weights of the neural network are applied to process each element in the sequence; see Figure 5. Here the neural network not only processes the current element in the sequence, but also draws upon the hidden layer of the previous element in the sequence. As a result, the recurrent structure allows to pass information onwards while iterating over the sequence. This implicitly creates a “state” whereby information is stored and accumulated. It thus encodes the whole sequence in the hidden layer $h_{\tau}^{(i)}$ of the last neural network.

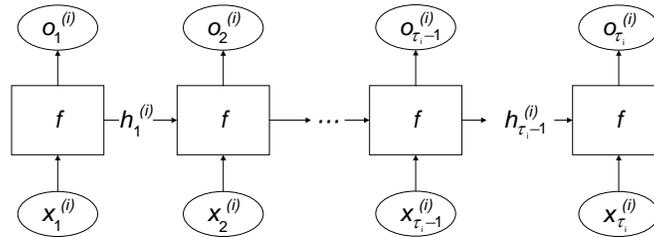


Figure 5: Schematic illustration of a recurrent neural network. Here the same network f is applied to each element $x_1^{(i)}, \dots, x_{\tau_i}^{(i)}$ of the sequence. The network utilizes not only the current element $x_k^{(i)}$ of the sequence as input, but also the internals (i. e. the hidden layer $h_{k-1}^{(i)}$) from the network belonging to the previous element in the sequence. Thereby, the knowledge of the whole sequence is accumulated in the hidden states. The hidden layers are denoted by $h_1^{(i)}, \dots, h_{\tau_i-1}^{(i)}$ and the output of the neural network by $o_1^{(i)}, \dots, o_{\tau_i}^{(i)}$.

Mathematically, the network input for element $k \in \{1, \dots, \tau_i\}$ is given by a concatenation between the current element $x_k^{(i)}$ and the previous hidden state $h_{k-1}^{(i)}$, i. e. $[x_k^{(i)}, h_{k-1}^{(i)}]$. The recurrent neural

network thus computes

$$f_{\text{RNN}}(x_1, \dots, x_\tau) = f([x_\tau^{(i)}, f([x_{\tau-1}^{(i)}, \dots, f([x_1^{(i)}]; W, b); \dots]; W, b]; W, b), \quad (10)$$

where the underlying network f can be, for instance, a simple single-layer neural network or a deep one. In practice, long sequences result in a large number of recurrent calls. This causes the optimization to become numerically highly unstable (Bengio et al., 1994). Hence, one might consider two extended network architectures that explicitly control how information is passed through:

- *Gated recurrent unit (GRU)*. The gated recurrent unit utilizes two underlying neural networks, an update gate and a reset gate, in order to explicitly determine how values in the hidden states are computed. Each of the gates is linked to an element-wise multiplication with values between 0 and 1, thus determining what ratio of each value is kept or discarded.
- *Long short-term memory (LSTM)*. In addition to the hidden state, this network structure builds upon an internal representation (called cell) which stores information (Hochreiter & Schmidhuber, 1997). It further includes multiple feedforward layers that formalize how values in cell are erased (called forget gate), how the input is written to the cell (input gate), and how the cell computes output (output gate).

The gated recurrent unit can be seen as a modification of the LSTM with fewer parameters. Yet, practical evidence suggests that both networks yield a comparable performance across a variety of tasks (Chung et al., 2014).

2.5. Embeddings

Most machine learning models and especially deep neural networks struggle with input in the form of categorical variables (Jia et al., 2014). The prime reasons are the mathematical properties of their one-hot encodings: one-hot encodings are (i) non-continuous, (ii) ignore interdependencies and (iii) their complexity scales with their cardinality. Different remedies have been suggested as follows. Zhang et al. (2016) suggest the use of a factorization machine for pre-processing categorical variables (which serves as one of our benchmarks). Cerda et al. (2018) develop similarity-based embeddings for handling noise in categorical variables (e.g. variables where texts that have spelling errors). Jia et al. (2014) present an embedding layer customized to categorical variables; however, this embedding layer is estimated via unsupervised pre-training: such estimation is not practical in most real-world OR applications where data is limited. In general, embeddings not only aid deep learning but also feature preprocessing in traditional machine learning (Jia et al., 2014; Zhu et al., 2017).

3. Methods and materials

3.1. Preprocessing

Deep neural networks entail an inherent advantage over traditional machine learning as they can handle data in its raw form without the need for manual feature engineering (LeCun et al., 2015). In contrast, the conventional approach is to first devise rules and extract specific representations from the

data. Examples are extracting edges from images instead of taking only the pixels as input, counting word triplets in natural language processing instead of processing the raw characters, or replacing individual time series observations with descriptive statistics such as minimum and maximum values. Deep learning largely circumvents the need for feature engineering and operates on the original data – pixels, characters, words or whole time series.

In a variety of applications, the only compulsory preprocessing step is to replace categorical values with a numeric representation. One usually approaches this by utilizing a one-hot encoding: the categorical value with K different entities is replaced by a K -dimensional vector where a 1 in an element indicates that the corresponding category is active and will otherwise be 0. Hence, the vector is 0 almost everywhere except for a single 1 that refers to the category.

One-hot vectors of large size lead to networks that become numerically difficult to optimize. An alternative is to replace sparse one-hot vectors with a so-called embedding, which maps them onto low-dimensional but dense representations (Goodfellow et al., 2017). Such embeddings are again modeled by a simple neural network and can even be optimized at the same time as the original deep neural network. Notably, embeddings can be constructed in such a way that they map from several neighboring input vectors onto the low-dimensional representation in order to encode additional semantics (Hirschberg & Manning, 2015). When processing natural language, one can utilize pre-computed word embeddings such as word2vec or GloVe; however, these are often not available for domain-specific applications in OR.

3.2. Proposed deep-embedded network architecture

Given the challenges surrounding categorical variables, we suggest a neural network architecture that specifically addresses the need of effectively handling categorical variables in OR practice. For this reason, our architecture presents a combination of the abovementioned embedding layers and, in addition, stacked neural layers (such as MLPs, CNNs or RNNs). Our architecture is particularly helpful in dealing with data of inhomogeneous types, i. e. consisting of both categorical and numerical input. Intuitively, our deep-embedded network architecture translates categorical variables into a dense representation, thus allowing for a numerically stable optimization of the neural network even for a large number of different categories. Mathematically, the architecture consists of two components: (i) the embedding of categorical variables and (ii) the concatenation of the embedded categorical variables with numerical variables in order to forward it to stacked neural layers. For both components, the optimization of all layers is performed simultaneously.

i) Let N_{cat} denote the number of categorical variables $\mathcal{I}_{\text{cat}}^{(i)}, i = 1, \dots, N_{\text{cat}}$ and let \mathcal{I}_{num} denote numerical variables. An embedding layer $\mathcal{E}^{(i)}$ is utilized to map one categorical variable $\mathcal{I}_{\text{cat}}^{(i)}$ to a dense vector v_i of dimension d via an embedding layer, i. e.

$$v_i = \mathcal{E}_i(\mathcal{I}_{\text{cat}}^{(i)}) \quad (11)$$

ii) Subsequently, all dense vectors $v_i, i = 1, \dots, N_{\text{cat}}$ are concatenated along with \mathcal{I}_{num} to one large

vector \mathcal{X} , i. e.

$$\mathcal{X} = [v_1, \dots, v_{N_{\text{cat}}}, \mathcal{I}_{\text{num}}]. \quad (12)$$

Then, \mathcal{X} yields input for default MLP, RNN, or CNN layers representing all input variables in dense form.

In order to counteract overfitting, the architecture is further extended by components representing dropout layers and batch normalization. In detail, batch normalization is added in all layers of the architecture up to the output layer, which utilizes batch normalization and dropout. Thereby, the non-linearity of the output layer is first applied before values are subject to batch normalization and dropout layers.

Our architecture is highly flexible: it can be utilized along MLPs, CNNs and RNNs that yield the second component of hidden and output layers. This is later shown based on different case studies involving a deep-embedded DNN and a deep-embedded LSTM. Algorithm 1 provides pseudocode of the deep-embedded network architecture utilized for our case studies.

Algorithm 1

<pre> 1: procedure DEEPEMBEDDEDDNN 2: Inputs: 3: - Vector of categorical features \mathcal{I}_{cat} 4: - Vector of numerical features \mathcal{I}_{num} 5: - Number of categorical features N_{cat} 6: - Number of numerical features N_{num} 7: - Dimension of embeddings d 8: - Number of hidden units h 9: - Number of layers l 10: - Dropout probability p 11: 12: for i in $1, \dots, N_{\text{cat}}$ do 13: $n \leftarrow$ number of different categories 14: of feature $\mathcal{I}_{\text{cat}}[i]$ 15: $v_i \leftarrow$ EmbeddingLayer(n, d)($\mathcal{I}_{\text{cat}}[i]$) 16: $v_i \leftarrow$ Dropout(v_i) 17: end for 18: $\mathcal{X} \leftarrow$ Concatenate($v_1, \dots, v_{N_{\text{cat}}}, \mathcal{I}_{\text{num}}$) 19: $\mathcal{X} \leftarrow$ FeedforwardLayer($N_{\text{num}} + N_{\text{cat}} * d, h$)($\mathcal{X}$) 20: for i in $1, \dots, l$ do 21: $\mathcal{X} \leftarrow$ FeedforwardLayer(h, h, ReLU)(\mathcal{X}) 22: $\mathcal{X} \leftarrow$ BatchNormalization(\mathcal{X}) 23: end for 24: $\mathcal{X} \leftarrow$ Dropout(\mathcal{X}) 25: $\mathcal{X} \leftarrow$ FeedforwardLayer($h, 1$)(\mathcal{X}) 26: end procedure </pre>	<pre> 1: procedure DEEPEMBEDDEDLSTM 2: Inputs: 3: - Matrix of categorical features \mathcal{I}_{cat} 4: - Number of categorical features N_{cat} 5: - Number of numerical features N_{num} 6: - Dimension of embeddings d 7: - Number of hidden units h 8: - Number of LSTM layers k 9: - Number of feedforward layers l 10: - Dropout probability p 11: 12: for i in $1, \dots, N_{\text{cat}}$ do 13: $n \leftarrow$ number of different categories 14: of feature $\mathcal{I}_{\text{cat}}[i]$ 15: $v_i \leftarrow$ EmbeddingLayer(n, d)($\mathcal{I}_{\text{cat}}[i]$) 16: $v_i \leftarrow$ Dropout(v_i) 17: end for 18: $\mathcal{X} \leftarrow$ Concatenate($(v_1, \dots, v_{N_{\text{cat}}}, \mathcal{I}_{\text{num}})$) 19: $\mathcal{X} \leftarrow$ FeedforwardLayer($N_{\text{num}} + N_{\text{cat}} * d, h$)($\mathcal{X}$) 20: for i in $1, \dots, k$ do 21: $\mathcal{X} \leftarrow$ LSTM layer(h)(\mathcal{X}) 22: end for 23: $\mathcal{X} \leftarrow$ Last hidden state of LSTM 24: for i in $1, \dots, l$ do 25: $\mathcal{X} \leftarrow$ FeedforwardLayer(h, h, ReLU)(\mathcal{X}) 26: $\mathcal{X} \leftarrow$ BatchNormalization(\mathcal{X}) 27: end for 28: $\mathcal{X} \leftarrow$ Dropout(\mathcal{X}, p) 29: $\mathcal{X} \leftarrow$ FeedforwardLayer($h, 1$)(\mathcal{X}) 30: end procedure </pre>
---	--

Our work differs from the embeddings that were used in prior literature (Jia et al., 2014; Zhu et al., 2017): (1) instead of using pre-training, our embeddings are directly integrated in the neural network architecture and are thus trained jointly. This aids a use in OR settings with scarce data. (2) In our work, we also integrate embeddings into a LSTM architecture where both LSTM and embeddings are trained

jointly. (3) Joint training introduces the risk of overfitting and, as a remedy, we refrain from using the naïve embedding layer (this only serves as a baseline). Instead, we suggest the use of a dropout-based embedding layer. The naïve embeddings later represent one of our baselines.

3.3. Estimation details

Following general conventions, we tune all parameters in traditional machine learning using 10-fold cross-validation and reduce the computational time for deep learning by selecting 10% of the training samples for validation (i. e. a random sample for case study 1 and chronological splits for the remaining time series predictions). We report our tuning ranges for each hyperparameter in the online appendix.

Given n different categories of a categorical variable, we calculate the dimension d of the embedding layer following common guidelines⁶ via $d = \sqrt[4]{n}$.

The required times for optimizing all model parameters are listed in the online appendix. As we shall see in the following, a higher optimization time for neural networks is countervailed by a higher prediction performance. We also provide loss curves for our deep-embedded networks in the online appendix.

3.4. Experimental setup

We perform a series of computational experiments in order to demonstrate the added benefit of deep learning in business analytics. For this purpose, we draw upon three case studies with public and proprietary industry data from different business areas as outlined in Table 3. As part of our benchmarks, we utilize common methods from traditional machine learning, namely, linear models (Lasso and ridge regression), a tree-based approach (i. e. the random forest) and further non-linear variants (i. e. a support vector machine and a single-layer neural network). In addition, our experiments build upon common techniques for preprocessing categorical variables from prior literature, namely a factorization machine (Zhang et al., 2016) and deep neural networks with naïve embedding layers (Jia et al., 2014). Performance is compared primarily based on the mean squared error for regression and the area under the curve for classification tasks.⁷

⁶Google Developers (2017): “Introducing TensorFlow Feature Columns”. <https://developers.googleblog.com/2017/11/introducing-tensorflow-feature-columns.html>, last accessed February 20, 2018.

⁷We refrain from using percentage errors as seen in some Kaggle competitions: there are crisp mathematical limitations (Tofallis, 2015) and these metrics do not measure costs. Instead, we evaluate all models through the lens of OR practitioners and thus report cost metrics.

Business area	Samples	Type of DNN	Public/Private	Categorical variables	Description
Operations management	273,750	GRU/LSTM	Private	Holiday flag	Predict the number of helpdesk tickets for the following hour
Inventory management	1,017,209	GRU/LSTM	Public	Store ID, state holiday, school holiday, store type, assortment	Predict the sales volume on day-ahead of pharmacy stores
Risk management	595,212	MLP	Public	14 Anonymized variables	Predict whether a customer will file a claim for an automotive insurance company

Table 3: Overview of our case studies.

4. Numerical results

4.1. Case study 1: Load forecasting of IT service requests

The second case study focuses on operations management, where the number of incoming service requests for an IT department is predicted in order to better adapt available capacities to the load. An accurate prediction of workloads is essential in order to coordinate employees in IT services (Bassamboo & Zeevi, 2009), as well as to avoid long waiting times for customers. For this purpose, we forecast the hourly number of tickets for the next work day based on historic values.

This experiment draws upon a proprietary dataset of helpdesk tickets from a Swiss organization with approximately 10,000 employees. Our dataset spans January 2011 to mid-November 2017, totaling 267,397 tickets. On average, 5.163 tickets are sent per day with a peak value of 547. In the first step, we counted the number of tickets in each hour of the study horizon. In total, this results in 60,192 samples. We augmented this dataset with an additional categorical variable that indicates holidays if there is currently a holiday. Our objective requires that we take a sequence of features x_1, \dots, x_τ as input that represent the workload at each hour. We then predict the workload of the following day $y_{\tau+1}, \dots, y_{\tau+24}$ across each hour.

Table 4 compares the prediction performance of traditional machine learning with that of deep learning. The average ticket volume per hour (i.e. 3.7 tickets/h) is used as a naïve baseline. Among the baseline models, we find the lowest mean absolute error on the test set when using a support vector machine, whereas a default deep neural network (two layers, each having 32 neurons) shows lowest mean squared error across the baseline models. The latter yields an improvement of 42.95 % compared to mean value as the predictor. Our deep-embedded architecture outperforms all baseline models including the default neural network. It yields an improvement of 3.75 % as compared to the best performing baseline, totaling in an improvement of 45.10 % points over the mean value as the baseline predictor. Statistical significance tests on the mean absolute error demonstrate that deep neural networks outperform our baselines to a statistically significant degree at the 1 % level. We also assessed the cost effectiveness in which prediction errors incur costs due to either idle workers or unsatisfied customers. Here our deep learning approach suggests savings of 29658.0 monetary units.

Model	Free parameters	Out-of-sample performance			
		MSE	MAE	Explained variance (R^2)	Prediction error costs
BASELINES					
Mean value as predictor	—	191.157	7.931	0.000	46124.6
ARMA/ARIMA	100	112.670	5.766	0.349	32454.4
Lasso	3649	115.983	5.842	0.324	32634.1
Ridge regression	3649	114.800	5.940	0.331	32663.2
Random forest	n/a	121.790	5.544	0.290	33015.0
Support vector machine	3649 [†]	129.110	5.355	0.250	32147.1
Default neural network (single-layer)	3649	112.197	5.650	0.346	33445.9
Default DNN (tuned depth)	374,902	109.040	5.785	0.365	33136.1
LSTM with factorization machine	14,596	121.782	5.458	0.322	32893.3
LSTM with embeddings	11,672	121.262	5.567	0.286	33091.3
PROPOSED NETWORK ARCHITECTURE					
Deep-embedded LSTM	11,672	104.952	5.101	0.390	29658.0

[†] We utilize a kernel approximation along with a linear support vector machine to reduce computation time.

Table 4: Numerical results for hourly prediction of incoming service requests. Here the prediction performance is compared between common baselines from traditional machine learning and sequence learning with deep neural networks. Reported are the mean squared error (MSE), the mean absolute error (MAE) on the original unscaled values, the explained variance (R^2) and the cost from prediction errors (where both an overestimated and an underestimated volume incur costs of 1.0 per ticket due to workers being idle or customers being unsatisfied). The best performance is highlighted in bold.

4.2. Case study 2: Sales forecasting

Sales forecasting facilitates the decision-making and planning within firms (Lau et al., 2018; Boone et al., 2018). Hence, this case study evaluates deep learning for sales forecasting, where the one-day-ahead sales volume for each store is predicted based on the history of previous sales. The dataset comprises 842,806 days of past sales from 1113 different pharmacies.⁸ In terms of preprocessing, we generate the following categorical variables that are embedded into a dense representation by our deep-embedded architecture: store ID, state holiday, school holiday, store type, assortment. In addition, we utilize the following numerical variables: distance to the next competitor, day of week and week of year.⁹ The average sales volume numbers to 6954.97 which presents our naïve baseline.

Table 5 lists the prediction performance of traditional machine learning, the default DNNs and our deep-embedded neural network. Among the baseline models, we find the lowest mean squared error on the test data when using a default DNN architecture (with three layers, each having 64 neurons). This model shows an improvement of 86.49% in terms of mean squared error. Our deep-embedded LSTM decreases the mean squared error by 21.73%. These improvements are also statistically significant at the 1% level as shown by a t -test. In total, our deep learning model yield an impressive improvement of up to 89.42% compared to predicting the mean value. The estimated costs from prediction errors again support the effectiveness of using deep learning. Altogether, similar to the first case study, our

⁸For reasons of reproducibility and comparability, we chose a public dataset from <https://www.kaggle.com/c/rossmann-store-sales>.

⁹These features were selected from the original dataset for the following reasons: some features are company-specific, whereas our selection should be widely applicable to other companies. Further, these features entail the largest predictive power with respect to the target variable and should thus be decisive.

deep-embedded neural network outperforms both traditional machine learning models and default neural networks.

Model	Free parameters	Out-of-sample performance			
		MSE	MAE	Explained variance (R^2)	Prediction error costs
BASELINES					
Mean value as predictor	—	9979310.0	2283.6	0.000	190565826.7
ARMA/ARIMA	50	2225451.0	1024.1	0.773	83949310.9
Lasso	512	2321334.6	1045.6	0.764	84511640.9
Ridge regression	512	2224958.3	1071.8	0.774	84212283.9
Random forest	n/a	1376444.2	805.7	0.860	64960715.8
Support vector machine	512 [†]	2453320.3	1130.6	0.751	84908542.3
Default neural network (single-layer)	512	1871942.3	968.7	0.820	71928430.8
Default DNN (tuned depth)	61,203	1348397.9	825.3	0.844	64891027.0
LSTM with factorization machine	3072	2432206.1	1185.1	0.753	95563305.5
LSTM with embeddings	10,849	1609544.4	925.3	0.837	65983043.1
PROPOSED NETWORK ARCHITECTURE					
Deep-embedded LSTM	10,849	1055438.7	713.6	0.893	57072820.5

[†] We utilize a kernel approximation along with a linear support vector machine to reduce computation time.

Table 5: Numerical results for predicting sales. The text dataset is kept the same as in previous experiments, yet we vary the size of the training set, i. e. we take a subset of our original set in order to study the sensitivity of deep learning to large-scale datasets. Here the prediction performance is compared between common baselines from traditional machine learning and sequence learning with deep neural networks. Reported are the mean squared error (MSE), the mean absolute error (MAE) on the original unscaled values, the explained variance (R^2) and the cost from prediction errors (where an overestimated volume incurs costs of 0.5 monetary units due to unsold products (incl. storage, etc.) and where an underestimated volume comes with penalty of 1.0 due to unsatisfied customers). The best performance is highlighted in bold.

4.3. Case study 3: Insurance credit scoring

Our third case study relates to credit risk analysis (Bassamboo et al., 2008), as it aims at identifying insurance credit scores. This translates into computing the likelihood of a customer filing a claim within a given time period. Based on these predictions, an insurance company can adapt their risk scoring and charge a corresponding premium. In general, predictive models found widespread application for tasks related to risk scoring (Lessmann et al., 2015).

This dataset comprises 595,212 real-world customer records with 43 numerical and 14 categorical covariates provided by an insurance company for automotives.¹⁰ The variable that is to be predicted is whether a customer would file an insurance claim in the subsequent year. Hence, we obtain a classification task with two classes: filed and not filed.

Given the inherent nature of risk modeling, the dataset is highly unbalanced, as only 3.6% of the customers filed a claim. For this reason, we assess models in terms of the AUC of the receiver operating characteristic and the Gini coefficient, both of which account for imbalanced classes in the dataset.

¹⁰For reasons of reproducibility and comparability, we experiment with a dataset from a public data science competitions (<https://www.kaggle.com/c/porto-seguro-safe-driver-prediction>). However, the participating models, as well as the test set of the competition, are confidential and we thus cannot compare our model to the deep neural network of winning contestant.

Table 6 compares the prediction performance. Among the baseline models from traditional machine learning, we find the highest Gini coefficient and the highest AUC score on the test data when using the ridge regression. This model shows an improvement of 26 % to the majority vote.

The deep-embedded DNN architecture (tuned to two hidden feedforward layers) outperforms all baseline models. The model increases the AUC score by 0.01 (i. e. 1.56 %) and the Gini coefficient by 0.02 (i. e. 7.14 %) as compared to the best-performing traditional machine learning model. Although the two-layer network is fairly shallow, deeper architectures did not result in a higher performance in this case. This can occur when the relationship between dependent variables and predicted outcome is fairly linear. Statistical tests on the AUC show that the improvement of the deep learning models is significant at the 1 % level. Altogether, our deep-embedded model can successfully improve predictive accuracy over both traditional models and out-of-the-box DNN architectures, thereby allowing for a more precise identification of insurance scores and thus premiums. This finding is additionally established when assessing the cost effectiveness, where, given the assumed numbers, deep learning results in a misclassification cost of only 19110.2 monetary units.

Model	Free parameters	Out-of-sample performance		
		Gini	AUC	Misclassification costs
BASELINES				
Majority vote (i. e. no claim)	—	0.0	0.5	21830.0
Lasso	272	0.254	0.627	19827.1
Ridge regression	272	0.260	0.630	19343.2
Random forest	n/a	0.0	0.5	21830.0
Support vector machine	272 [†]	0.0	0.5	21830.0
Default neural network (single-layer)	272	0.0	0.5	21830.0
Default DNN (tuned depth)	37,203	0.249	0.625	19931.1
DNN with factorization machine	1360	0.259	0.629	19958.1
DNN with embeddings	177,745	0.232	0.616	20981.3
PROPOSED NETWORK ARCHITECTURE				
Deep-embedded DNN	177,745	0.280	0.640	19112.0

[†] We utilize a kernel approximation along with a linear support vector machine to reduce computation time.

Table 6: Numerical results for determining insurance credit scores. Here the prediction performance is compared between common baselines from traditional machine learning and deep neural networks. Misclassification costs are computed based on an average loss of 0.1 monetary units per false positive and a lost profit of 10.0 for false negative.

The best performance is highlighted in bold.

4.4. Sensitivity to size of training set

The aforementioned case studies show that deep neural networks outperform traditional machine learning methods. This section shows results of further experiments that (1) compare the performance of deep learning and traditional machine learning for a variety of training observations and (2) give insights to the training process of deep neural networks.

To evaluate the impact of the number of training observations on the performance, we utilize a subset of the training set and train a deep neural network and a random forest with it. Figure 6 shows that the random forest outperforms the deep neural network up until 10 % of the training observations. However,

utilizing the complete dataset yields favorable results for the deep neural network. This experiment supports our statement that large datasets are needed to train deep neural networks (cp. Section 1)

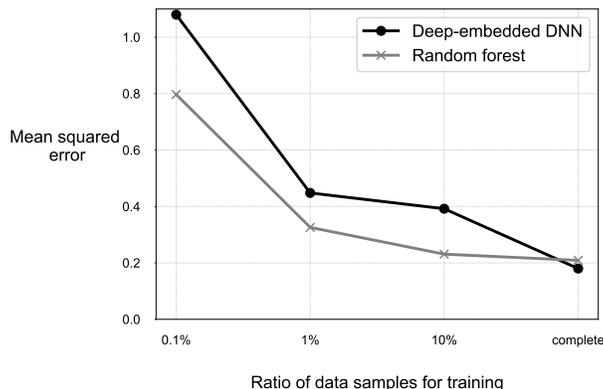


Figure 6: Experiment comparing number of training samples with the performance for sales forecasting. Deep neural networks still benefit from large amounts of data, whereas the performance increase of the random forest plateaus.

4.5. Proposed interpretation of deep neural networks

The interpretation of black-box models poses a major challenge for machine learning and can significantly reduce the barriers to technology adoption. A possible remedy is given by an ex post analysis through a visual inspection of prediction results. These same tools as for other machine learning classifiers can be applied here: for example, partial dependency plots visualize the marginal effects of one or two features having on the prediction of a machine learning model (Friedman, 2001). A partial dependency plot can show whether the relationship between the target and an input feature has a certain shape, such as linear, monotonous, or more complex. In other works, Shapley values were proposed to calculate the significance of a variable by comparing what a model predicts with and without the feature (Lundberg & Lee, 2017). In detail, the Shapley value of a feature value is the average change in the prediction that the model makes when the feature value is added as an input. Examples thereof are provided in the online appendix.

We propose an additional technique that is tailored to the use of embeddings as in our our deep-embedded networks. That is, we suggest to combine a lower-dimensional representation of the embeddings, specifically, the so-called t-SNE dimensionality reduction (van der Maaten & Hinton, 2008), together with their marginal effects. This has obvious advantages: it allows us to visualize the embedded categorical variables, despite that the input has two or three dimensions. Here the embeddings are clustered in a way that input with similar categorical variables are located in close proximity to each other, whereas dissimilar categories are distant. This is shown by the black points in Figure 7. Different from a default t-SNE plot, we additionally derive marginal effects as in a partial dependence plot, so that changes in the input dimension can be studied. This yields a natural interpretation of the inner structure behind our deep-embedded network.

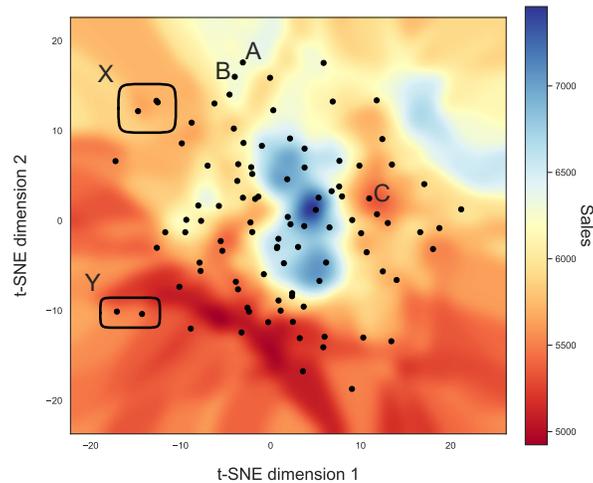


Figure 7: Proposed visual interpretation of our deep-embedded networks based on combining t-SNE dimensionality reduction with marginal effects. Black points refer to stores that were arranged based on a t-SNE dimensionality of the original categorical input (here: shown for case study 2); background color informs how embeddings are translated into an average prediction.

Figure 7 shows store embeddings for case study 2, which aids managerial decision-making follows. First, it helps management in identifying stores with similar characteristics but where some are over- or underperforming with respect to similar stores (see the red peaks). This should trigger a subsequent root cause whether the management, assortment, or promotion campaign needs adjustments. Second, the t-SNE plot supports management in location planning for new stores (here: indicated, e. g., by A, B, C). It displays the expected sales but also represents to what extent the new business is similar to the core operational area. Management might prefer store C over A or B despite lower expected sales since it is similar to other stores and, hence, prior experience on *modi operati* can be easier transferred. Third, when management wants to reduce the number of stores, this plot helps in selecting stores for disposal. Specifically, it considers trade-offs between sales and the core business area. In our example, the core business area is given by the dense inner area, whereas stores in circles X and Y are outside of it and thus present example candidates for disposal.

5. Discussion

Deep learning is essential in the context of big data and, for this purpose, its performance across different scenarios is compared and contrasted in this overview article. The empirical results of the different cases studies suggest that deep learning is a feasible and effective method, which can considerably and consistently outperform its traditional counterparts in both prediction and operational performance from the family of data-analytic models. As such, DNNs are able identify previously unknown, potentially useful, non-trivial, and interesting patterns more accurately than other popular predictive models such as random forest, artificial neural networks, and support vector machines. One of the reasons they yield

superior results originates from the strong mathematical assumptions in traditional machine learning, whereas these are relaxed by DNNs as a result of their larger parameter space. Even though various building blocks exist, actual applications benefit from customized architectures as demonstrated by the deep-embedded architecture in this research. We obtain performance improvements between 1.56% and 21.73% across all metrics by replacing the default architecture with our proposed deep-embedded network.

5.1. Managerial implications

Business analytics refers to the art and science of converting (big) data into business insights for faster and more accurate managerial decision-making (Chen et al., 2012; Baesens et al., 2016). It is positioned at the intersection of several disciplines, of which machine learning, operations research, statistics, and information systems are of particular relevance. Accordingly, the underlying objective concerns the ability of understanding and communicating insights gleaned from descriptive, predictive, and prescriptive analytics. Business analytics is arguably the most critical enabler for businesses to survive and thrive in the stiff global marketplace, where evidence-based decisions and subsequent actions are driven from data and analytical modeling. Yet, a recent survey across 3,000 executives reveals that 85 percent of them believe that predictive analytics implies a competitive advantage, though only one in 20 companies has adopted these techniques (Ransbotham et al., 2017). This becomes especially crucial in the light of deep learning, which is forecasted to deepen competition over analytics. It is thus assumed that advanced analytics, as well as big data, have widespread implications for management (see e.g. Agarwal & Dhar, 2014; George et al., 2014, 2016).

A challenge for most managers remains with regard to how they can identify valuable use cases of predictive analytics and especially deep learning. Strategically positioned at this target, the current work presents a framework to show the viability and superiority of DNNs within the business analytics domain through several case studies. Hence, the primary message of this overview article is to review the applicability of deep learning in improving decision support across core areas of businesses operations. As a direct implication, the generic approach proposed in this work can be utilized to create an automated decision support system, which in turn would increase the quality of decisions both in terms of efficiency and effectiveness. To facilitate implementation, prior research has suggested a simple rule-of-thumb: all tasks which involve only one second of thought can be replaced by predictive analytics (Ng, 2016b). This is confirmed by our experiments, since it can be safely claimed that the proposed DNN-based methodologies would bring significant financial gains for organizations across various areas of business domain as exemplified here in insurance credit scoring, IT service request predictions, and sales forecasting.

Data is the key ingredient for deep learning to become effective. In fact, a distinguishing feature of deep neural networks links to its ability to still “learn” better predictions from large-scale data as compared to traditional methods which often end up in a saturation point where larger datasets no longer improve forecasts. Hence, firms require extensive datasets and, for this topic, we refer to Corbett (2018) for an extensive discussion.

Despite the challenges surrounding the use of deep learning, this technique also entails practical benefits over traditional machine learning. In particular, feature engineering was a critical element

hitherto when deriving insights from big data (Feng & Shanthikumar, 2018). Conversely, deep neural networks circumvent the need for feature engineering. This is especially helpful in the case of sequential data, such as time series or natural language (Kratzwald et al., 2018), where the raw data can now be directly fed into the deep neural network. A similarity becomes evident between feature engineering and embeddings, yet the latter is fully data-driven and can easily be customized to domain-specific applications (Kraus & Feuerriegel, 2017).

Deep learning can achieve higher prediction accuracies than traditional machine learning, though they are still at the embryonic stage within the areas of business analytics. Therefore, practitioners should be aware of the caveat that DNNs introduce fairly complex architectures, which, in turn, necessitate a thorough understanding and careful implementation for valid and impactful results. In addition, the value of deep learning expands beyond the scope of mere business analytics within popular areas of medicine, healthcare, engineering, and etc. with a direct societal benefit could be realized. Yet conservative estimates suggest that firms at the frontier of predictive analytics need up to 1–2 years for replicating results from research (Ng, 2016b), while the actual number for the average firm is likely to be larger. Hence, practitioners can take our article as a starting point for devising a corresponding analytics strategy.

5.2. Limitations and possible remedies

Further enhancing deep learning algorithmically could also be addressed by future work. We point towards three key challenges that we consider as especially relevant for the operations research community. (1) The configuration of DNNs represents a challenging task, since it still requires extensive parameter tuning to achieve favorable results. Recently, a series of tools for automatizing the tuning process have been proposed under the umbrella term “AutoML”. This presents a promising path towards accelerating the deployment of effective DNN architectures. (2) DNNs are currently only concerned with point estimates, while their predictions commonly lack rigorous uncertainty quantification. We hope for increasing efforts with the goal of either directly modeling distributions or even develop Bayesian variants of DNNs. (3) Although ex post analysis have led to large steps in understanding the behavior of deep neural networks, accountability and interpretability are widely regarded as a weakness in deep learning and. Given the recent interest of policy-makers, these two active fields of research likely represent a pressing issue for a variety of applications in the near future. Potentially, estimating structural models via variational inferences (Hoffman et al., 2013; Kraus & Feuerriegel, 2019a), or interpreting attention mechanisms in neural networks could lead to more insights of the prediction (Kraus & Feuerriegel, 2019b).

5.3. Roadmap for future research

Deep learning has great potential to create additional value for firms, organizations, and individuals in a variety of business units and domains. Yet, its actual use in the field of operations and analytics remains scarce. Hence, it is recommended that the goal of future research in the realm of business analytics could be centered around at identifying precious use cases, as well as outlining potential value gains. This further necessitates better recommendations with regard to combinations of network architectures, training routines, parameter fine-tuning that yield favorable prediction performances in the case of DNNs.

For instance, further research efforts could eventually lead to novel techniques that customize network architectures when fusing different data sources.

Similar to the other forms of machine learning methods in predictive analytics, deep learning also merely provides predictive insights, but rarely presents actual management strategies to reach the desired outcome. As a remedy to this, future studies could focus on determining how predictions can actually be translated into effective decision-making. This is another compelling direction for future research, since the OR community is more important than ever to translate predictions into decision (Feng & Shanthikumar, 2018). Here we point towards inverse control (Bertsimas & Kallus, 2014; Bertsimas & King, 2016), Markov decision processes and bandits as promising techniques.

There is further untapped potential as a number of innovations in the domain of deep learning have not found its way into widespread adoption. First, sequence learning as presented in this paper takes sequences as input but still yields predictions in the form of a vector with fixed dimensions. So-called “seq2seq” techniques allow to make inferences where each prediction is a sequence of arbitrarily varying length (Goodfellow et al., 2017). This could be helpful when deep neural networks should compute routes or graphs as output. Second, datasets from practical applications often cannot warrant the necessary size that is needed for an effective use of deep learning. A remedy is given by transfer learning where a different yet related dataset is utilized for an inductive knowledge transfer. This can even facilitate an interesting strategy for domain customizations of predictive analytics (Kratzwald & Feuerriegel, 2019). Third, generative adversarial networks draw upon the idea of zero-sum games and train two competing neural networks, where one provides predictions in the usual fashion, while the other generates a sample input that matches a given label (Goodfellow et al., 2014). As an illustrative example, this can essentially yield a creativity mechanism (Hope et al., 2017), yet its benefit in OR applications still needs to be demonstrated.

6. Conclusion

Business analytics refers to the ability of firms and organization to collect, manage, and analyze data from a variety of sources in order to enhance the understanding of business processes, operations, and systems. As companies generate more data at ever-faster rates, the need for advances in predictive analytics becomes a pertinent issue, which can hypothetically improve decision-making processes and decision support for businesses. Consequently, competition in terms of analytics has become prevalent as even minor improvements in prediction accuracy can bolster revenues in greater folds and thus demands a better understanding of deep learning. The case studies conducted in this overview article are purposefully selected from different areas of operations research in order to validate the fact that DNNs help in improving operational performance. However, a customized network architecture is oftentimes beneficial, such as demonstrated by our deep-embedded network which attains performance improvements between 1.56 % and 21.73 % across all metrics over a default, out-of-the-box architecture.

7. Acknowledgments

This work was part-funded by the Swiss National Science Foundation (SNF), Digital Lives grant 10DL18-183149.

References

References

- Adamopoulos, P., Ghose, A., & Todri, V. (2018). The impact of user personality traits on word of mouth: Text-mining social media platforms. *Information Systems Research*, *29*, 612–640.
- Agarwal, R., & Dhar, V. (2014). Editorial – big data, data science, and analytics: The opportunity and challenge for is research. *Information Systems Research*, *25*, 443–448.
- Baesens, B., Bapna, R., Marsden, J. R., Vanthienen, J., & Zhao, J. L. (2016). Transformational issues of big data and analytics in networked business. *MIS Quarterly*, *40*, 807–818.
- Bassamboo, A., Juneja, S., & Zeevi, A. (2008). Portfolio credit risk with extremal dependence: Asymptotic analysis and efficient simulation. *Operations Research*, *56*, 593–606.
- Bassamboo, A., & Zeevi, A. (2009). On a data-driven method for staffing large call centers. *Operations Research*, *57*, 714–726.
- Bellman, R. (1972). *Dynamic Programming*. Princeton, NJ: Princeton University Press.
- Bengio, Y., Simard, P., & Frasconi, P. (1994). Learning long-term dependencies with gradient descent is difficult. *IEEE Transactions on Neural Networks*, *5*, 157–166.
- Bertsimas, D., & Kallus, N. (2014). From predictive to prescriptive analytics. *arXiv preprint arXiv:1402.5481*, .
- Bertsimas, D., & King, A. (2016). OR forum: An algorithmic approach to linear regression. *Operations Research*, *64*, 2–16.
- Bertsimas, D., & Shioda, R. (2007). Classification and regression via integer optimization. *Operations Research*, *55*, 252–271.
- Boone, T., Ganeshan, R., Hicks, R. L., & Sanders, N. R. (2018). Can google trends improve your sales forecast? *Production and Operations Management*, *27*, 1770–1774.
- Carbonneau, R., Laframboise, K., & Vahidov, R. (2008). Application of machine learning techniques for supply chain demand forecasting. *European Journal of Operational Research*, *184*, 1140–1154.
- Cerda, P., Varoquaux, G., & Kégl, B. (2018). Similarity encoding for learning with dirty categorical variables. *Machine Learning*, *107*, 1477–1494.
- Chen, Chiang, & Storey (2012). Business intelligence and analytics: From big data to big impact. *MIS Quarterly*, *36*, 1165–1188.
- Choromanska, A., Henaff, M., Mathieu, M., Arous, G., & LeCun, Y. (2015). The loss surfaces of multilayer networks. In *International Conference on Artificial Intelligence and Statistics (AISTATS '15)*.
- Chung, J., Gulcehre, C., Cho, K., & Bengio, Y. (2014). Empirical evaluation of gated recurrent neural networks on sequence modeling. In *Advances in Neural Information Processing Systems (NIPS '14)*.
- Corbett, C. J. (2018). How sustainable is big data? *Production and Operations Management*, *27*, 1685–1695.

- Cybenko, G. (1989). Approximation by superpositions of a sigmoidal function. *Mathematics of Control, Signals, and Systems*, 2, 303–314.
- Davenport, T. H., & Harris, J. G. (2007). *Competing on analytics: The new science of winning*. Boston, MA: Harvard Business School.
- Delen, D., Oztekin, A., & Tomak, L. (2012). An analytic approach to better understanding and management of coronary surgeries. *Decision Support Systems*, 52, 698–705.
- Feng, Q., & Shanthikumar, J. G. (2018). How research in production and operations management may evolve in the era of big data. *Production and Operations Management*, 27, 1670–1684.
- Fischer, T., & Krauss, C. (2018). Deep learning with long short-term memory networks for financial market predictions. *European Journal of Operational Research*, 270, 654–669.
- Friedman, J. H. (2001). Greedy function approximation: A gradient boosting machine. *Annals of Statistics*, (pp. 1189–1232).
- George, G., Haas, M. R., & Pentland, A. (2014). Big data and management. *Academy of Management Journal*, 57, 321–326.
- George, G., Osinga, E. C., Lavie, D., & Scott, B. A. (2016). Big data and data science methods for management research. *Academy of Management Journal*, 59, 1493–1507.
- Goodfellow, I., Bengio, Y., & Courville, A. (2017). *Deep Learning*. Cambridge, MA: MIT Press.
- Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., & Bengio, Y. (2014). Generative adversarial nets. In *Advances in Neural Information Processing Systems (NIPS '14)*.
- He, K., Zhang, X., Ren, S., & Sun, J. (2016). Deep residual learning for image recognition. In *Conference on Computer Vision and Pattern Recognition (CVPR '16)*.
- Henke, N., Bughin, J., Chui, M., Manyika, J., Saleh, T., Wiseman, B., & Sethupathy, G. (2016). The age of analytics: Competing in a data-driven world. *McKinsey Global Institute*, .
- Hinton, G. E., & Salakhutdinov, R. R. (2006). Reducing the dimensionality of data with neural networks. *Science*, 313, 504–507.
- Hirschberg, J., & Manning, C. D. (2015). Advances in natural language processing. *Science*, 349, 261–266.
- Hochreiter, S., & Schmidhuber, J. (1997). Long short-term memory. *Neural Computation*, 9, 1735–1780.
- Hoffman, M. D., Blei, D. M., Wang, C., & Paisley, J. (2013). Stochastic variational inference. *The Journal of Machine Learning Research*, 14, 1303–1347.
- Hope, T., Chan, J., Kittur, A., & Shahaf, D. (2017). Accelerating innovation through analogy mining. In *International Conference on Knowledge Discovery and Data Mining (KDD '17)*.
- Hu, Q., Chakhar, S., Siraj, S., & Labib, A. (2017). Spare parts classification in industrial manufacturing using the dominance-based rough set approach. *European Journal of Operational Research*, 262, 1136–1163.
- Ioffe, S., & Szegedy, C. (2015). Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *International Conference on Machine Learning (ICML '15)*.

- Jia, Y., Shelhamer, E., Donahue, J., Karayev, S., Long, J., Girshick, R., Guadarrama, S., & Darrell, T. (2014). Caffe: Convolutional architecture for fast feature embedding. In *ACM International Conference on Multimedia* (pp. 675–678).
- Judd, J. S. (1990). *Neural Network Design And The Complexity Of Learning*. Cambridge, MA: MIT Press.
- Kiwiel, K. C. (2001). Convergence and efficiency of subgradient methods for quasiconvex minimization. *Mathematical Programming*, *90*, 1–25.
- Kratzwald, B., & Feuerriegel, S. (2019). Putting question-answering systems into practice. *ACM Transactions on Management Information Systems*, *9*.
- Kratzwald, B., Ilić, S., Kraus, M., Feuerriegel, S., & Prendinger, H. (2018). Deep learning for affective computing: Text-based emotion recognition in decision support. *Decision Support Systems*, *115*, 24–35.
- Kraus, M., & Feuerriegel, S. (2017). Decision support from financial disclosures with deep neural networks and transfer learning. *Decision Support Systems*, *104*, 38–48.
- Kraus, M., & Feuerriegel, S. (2019a). Forecasting remaining useful life: Interpretable deep learning approach via variational Bayesian inferences. *Decision Support Systems*, .
- Kraus, M., & Feuerriegel, S. (2019b). Sentiment analysis based on rhetorical structure theory: Learning deep neural networks from discourse trees. *Expert Systems with Applications*, *118*, 65–79.
- Krauss, C., Do, X. A., & Huck, N. (2017). Deep neural networks, gradient-boosted trees, random forests: Statistical arbitrage on the S&P 500. *European Journal of Operational Research*, *259*, 689–702.
- Lau, R. Y. K., Zhang, W., & Xu, W. (2018). Parallel aspect-oriented sentiment analysis for sales forecasting with big data. *Production and Operations Management*, *27*, 1775–1794.
- LeCun, Y., Bengio, Y., & Hinton, G. (2015). Deep learning. *Nature*, *521*, 436–444.
- LeCun, Y., Bottou, L., Orr, G. B., & Müller, K.-R. (1998). Efficient BackProp. In *Neural Networks* (pp. 9–50).
- Lee, H. L. (2018). Big data and the innovation cycle. *Production and Operations Management*, *27*, 1642–1646.
- Lessmann, S., Baesens, B., Seow, H.-V., & Thomas, L. C. (2015). Benchmarking state-of-the-art classification algorithms for credit scoring: An update of research. *European Journal of Operational Research*, *247*, 124–136.
- Li, W., Chen, H., & Nunamaker, J. F. (2017). Identifying and profiling key sellers in cyber carding community: Azsecure text mining system. *Journal of Management Information Systems*, *33*, 1059–1086.
- Lim, E.-P., Chen, H., & Chen, G. (2013). Business intelligence and analytics: Research directions. *ACM Transactions on Management Information Systems*, *3*.
- Lundberg, S. M., & Lee, S.-I. (2017). A unified approach to interpreting model predictions. In *Advances in Neural Information Processing Systems (NIPS '17)*.
- Mahamad, A. K., Saon, S., & Hiyama, T. (2010). Predicting remaining useful life of rotating machinery

- based artificial neural network. *Computers & Mathematics with Applications*, 60, 1078–1087.
- Mazhar, M., Kara, S., & Kaebernick, H. (2007). Remaining life estimation of used components in consumer products: Life cycle data analysis by Weibull and artificial neural networks. *Journal of Operations Management*, 25, 1184–1193.
- Misiunas, N., Oztekin, A., Chen, Y., & Chandra, K. (2016). DEANN: A healthcare analytic methodology of data envelopment analysis and artificial neural networks for the prediction of organ recipient functional status. *Omega*, 58, 46–54.
- Montavon, G., Orr, G., & Müller, K.-R. (2012). *Neural Networks: Tricks Of The Trade*. (2nd ed.). Heidelberg: Springer.
- Mortenson, M. J., Doherty, N. F., & Robinson, S. (2015). Operational research from taylorism to terabytes: A research agenda for the analytics age. *European Journal of Operational Research*, 241, 583–595.
- Ng, A. (2016a). *Machine Learning Yearning: Technical Strategy for AI Engineers, In the Era of Deep Learning, Draft Version 0.5*.
- Ng, A. (2016b). What artificial intelligence can and can't do right now. *Harvard Business Review*, 9.
- Oztekin, A., Al-Ebbini, L., Sevcli, Z., & Delen, D. (2018). A decision analytic approach to predicting quality of life for lung transplant recipients: A hybrid genetic algorithms-based methodology. *European Journal of Operational Research*, 266, 639–651.
- Oztekin, A., Kizilaslan, R., Freund, S., & Iseri, A. (2016). A data analytic approach to forecasting daily stock returns in an emerging market. *European Journal of Operational Research*, 253, 697–710.
- Probst, M., Rothlauf, F., & Grahl, J. (2017). Scalability of using restricted boltzmann machines for combinatorial optimization. *European Journal of Operational Research*, 256, 368–383.
- Ransbotham, S., Kiron, D., Gerbert, P., & Reeves, M. (2017). Reshaping business with artificial intelligence: Closing the gap between ambition and action. *MIT Sloan Management Review*, 59.
- Ranyard, J. C., Fildes, R., & Hu, T.-I. (2015). Reassessing the scope of OR practice: The influences of problem structuring methods and the analytics movement. *European Journal of Operational Research*, 245.
- Rumelhart, D. E., Hinton, G. E., & Williams, R. J. (1986). Learning representations by back-propagating errors. *Nature*, 323, 533–536.
- Russell, S. J., Norvig, P., & Davis, E. (2010). *Artificial Intelligence: A Modern Approach*. (3rd ed.). Upper Saddle River, NJ: Prentice Hall.
- Saad, D. (1998). *Online Learning in Neural Networks*. Cambridge, MA: Cambridge University press.
- Schmidhuber, J. (2015). Deep learning in neural networks: An overview. *Neural networks*, 61, 85–117.
- Scholz, M., Pfeiffer, J., & Rothlauf, F. (2017). Using PageRank for non-personalized default rankings in dynamic markets. *European Journal of Operational Research*, 260, 388–401.
- Silver, D., Schrittwieser, J., Simonyan, K., Antonoglou, I., Huang, A., Guez, A., Hubert, T., Baker, L., Lai, M., Bolton, A., Chen, Y., Lillicrap, T., Hui, F., Sifre, L., van den Driessche, G., Graepel, T., & Hassabis, D. (2017). Mastering the game of Go without human knowledge. *Nature*, 550, 354–359.

- Srivastava, N., Hinton, G. E., Krizhevsky, A., Sutskever, I., & Salakhutdinov, R. (2014). Dropout: A simple way to prevent neural networks from overfitting. *Journal of Machine Learning Research*, *15*, 1929–1958.
- Stachenfeld, K. L., Botvinick, M. M., & Gershman, S. J. (2017). The hippocampus as a predictive map. *Nature Neuroscience*, *20*, 1643–1653.
- Sun, Y., Ma, L., & Morris, J. (2009). A practical approach for reliability prediction of pipeline systems. *European Journal of Operational Research*, *198*, 210–214.
- Tofallis, C. (2015). A better measure of relative prediction accuracy for model selection and model estimation. *Journal of the Operational Research Society*, *66*, 1352–1362.
- van der Maaten, L., & Hinton, G. (2008). Visualizing data using t-SNE. *Journal of Machine Learning Research*, *9*, 2579–2605.
- Venkatesh, K., Ravi, V., Prinzie, A., & van den Poel, D. (2014). Cash demand forecasting in ATMs by clustering and neural networks. *European Journal of Operational Research*, *232*, 383–392.
- Zhang, W., Du, T., & Wang, J. (2016). Deep learning over multi-field categorical data. In *European Conference on Information Retrieval* (pp. 45–57).
- Zhu, J., Shan, Y., Mao, J. C., Yu, D., Rahmanian, H., & Zhang, Y. (2017). Deep embedding forest: Forest-based serving with deep embedding features. In *International Conference on Knowledge Discovery and Data Mining (KDD '17)* (pp. 1703–1711).

Appendix A. Runtime analysis

Table A.1 reports the runtime for optimizing the parameters inside the machine learning models.

Model	Optimization time (min:sec)		
	Load forecasting	Sales forecasting	Credit scoring
BASELINES			
ARMA/ARIMA	0:00	0:07	–
Lasso	0:00	0:35	0:07
Ridge regression	0:00	0:08	0:04
Random forest	1:37	358:10	129:30
Support vector machine	0:12	17:33	7:16
Default neural network (single-layer)	0:07	0:23	0:46
Default DNN (tuned depth)	1:47	2:58	1:25
PROPOSED NETWORKS ARCHITECTURE			
Deep-embedded DNN/LSTM	4:17	7:29	3:24

Table A.1: Time (in minutes:seconds) for optimizing the machine learning models for our case studies.

Appendix B. Ex post interpretation

Figure B.1 and Figure B.2 illustrate a partial dependence plot and Shapley values in order to interpret model behavior. The former shows increasing sales at around the 20th calendar week and, in particular, at the end of the year. Figure B.2 shows that the assortment and specific store types ($= 0$) are associated with a decreasing effect on the predicted sales for the upper sample, whereas the middle of the week (day of week equals three) and school holidays lead to an increase in the predicted value.

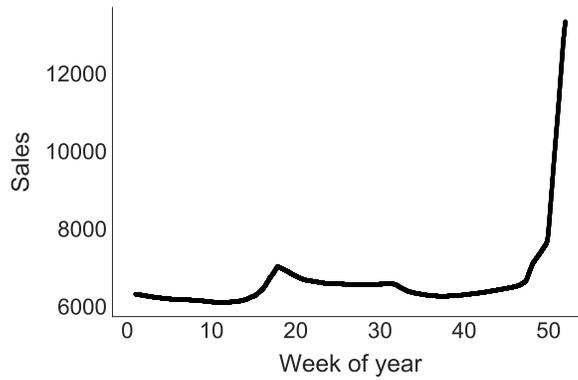


Figure B.1: Partial dependence plot that shows the marginal effect of the (scaled) week-of-year variable on the prediction outcome.

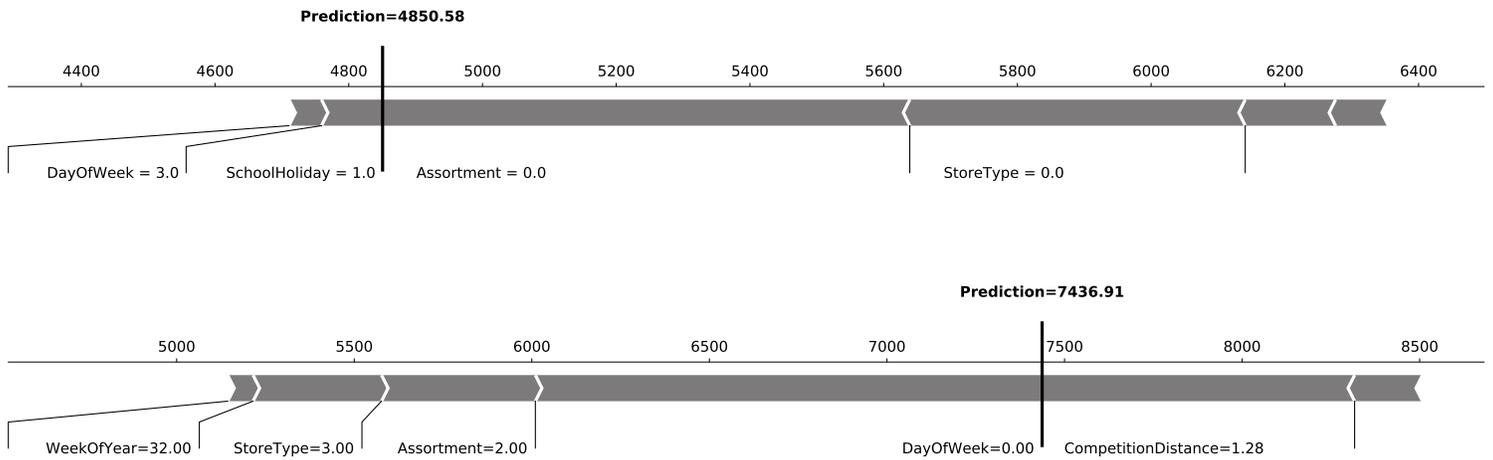


Figure B.2: Shapley values show the effect of the input features on the prediction outcome (4850.58 and 7436.91).

Appendix C. Hyperparameter tuning

Table C.2 reports the tuning parameters used in our grid search.

Model	Tuning parameter	Tuning range
Lasso	Regularization strength α	$10^{-3}, \dots, 10^{+3}$
Ridge regression	Regularization strength α	$10^{-3}, \dots, 10^{+3}$
ARMA/ARIMA [†]	Auto-regressive lags p	2, 5, 10, 50, 100, 200
	Difference iterations d	0, 1, 2, 3
	Moving-average terms q	0, 1, 2, 3
Random forest	Number of trees	200, 500
	Maximum depth of trees	2, 5, 10
	Number of randomly-sampled variables	1, 3, 5, 10
Support vector machine	Kernel function	Linear, radial [‡]
	Cost	$10^{-3}, \dots, 10^{+3}$
Single-layer neural network	Learning rate	0.001, 0.005, 0.01, 0.05
	Batch size	32, 64, 256
Default neural network	Learning rate	0.001, 0.005, 0.01, 0.05
	Batch size	32, 64, 256
	Number of neurons in hidden layers	Task-dependent
Deep neural network	Number of hidden layers	1, 2, 3
	Learning rate	0.001, 0.005, 0.01, 0.05
	Dropout rate	0.0, 0.25, 0.5, 0.75
	Number of neurons in hidden layers	Task-dependent
	Batch size	32, 64, 256

[†] The ARMA and ARIMA model is only evaluated for time series prediction tasks.

[‡] We utilize a kernel approximation to reduce computation time.

Table C.2: Grid search for hyperparameter tuning.

Appendix D. Loss curves

Figure D.3 depicts our optimization processes for the LSTM and the GRU network for the load forecasting and the sales forecasting tasks.

(a) Training and validation graphs of load forecasting (b) Training and validation graphs of sales forecasting

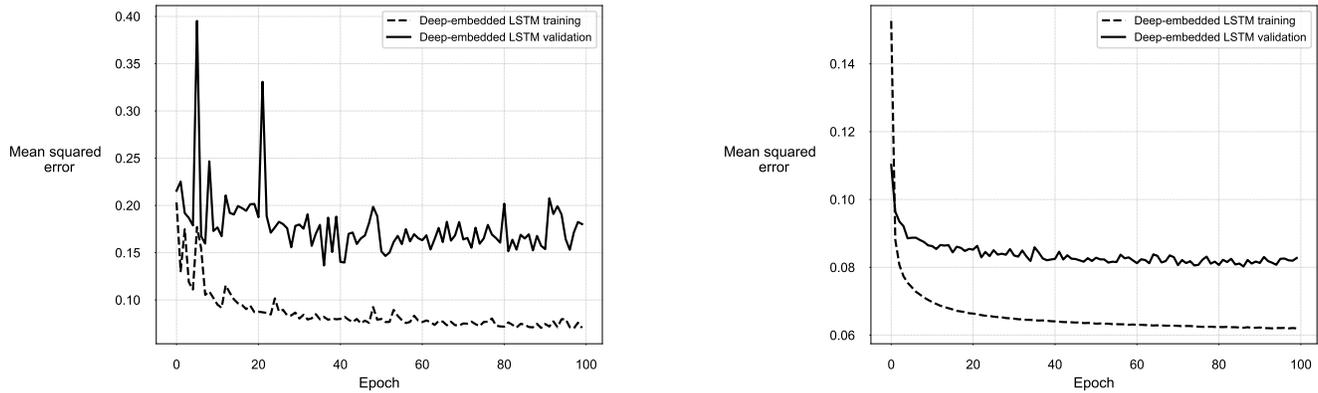


Figure D.3: Training and validation loss by epoch for our best performing deep neural networks on the load forecasting and sales forecasting studies. One epoch depicts a complete pass of all training samples through the optimization process. On the left, we see that the gated recurrent unit overfits after around 40 epochs. Therefore, additional regularization is needed to stabilize the optimization. On the right, both long short-term memory network and gated recurrent unit perform equally and converge nicely.