



Published in final edited form as:

Data Min Knowl Discov. 2021 January ; 35(1): 46–87. doi:10.1007/s10618-020-00722-8.

A Survey of Deep Network Techniques All Classifiers Can Adopt

Alireza Ghods, Diane J. Cook

School of Electrical Engineering and Computer Science, Washington State University, Pullman, WA, 99164

Abstract

Deep neural networks (DNNs) have introduced novel and useful tools to the machine learning community. Other types of classifiers can potentially make use of these tools as well to improve their performance and generality. This paper reviews the current state of the art for deep learning classifier technologies that are being used outside of deep neural networks. Non-neural network classifiers can employ many components found in DNN architectures. In this paper, we review the feature learning, optimization, and regularization methods that form a core of deep network technologies. We then survey non-neural network learning algorithms that make innovative use of these methods to improve classification performance. Because many opportunities and challenges still exist, we discuss directions that can be pursued to expand the area of deep learning for a variety of classification algorithms.

Keywords

Deep Learning; Deep Neural Networks; Optimization; Regularization

1 Introduction

The objective of supervised learning algorithms is to identify an optimal mapping between input features and output values based on a given training dataset. A supervised learning method that is attracting substantial research and industry attention is DNN. DNNs have a profound effect on our daily lives; they are found in search engines (Guo et al. 2017), self-driving cars (Ndikumana and Hong 2019), health care systems (Esteva et al. 2019), and consumer devices such as smart-phones and cameras (Gjoreski et al. 2020; Yang et al. 2020). Convolutional Neural Networks (CNN) have become the standard for processing images (Feng et al. 2019), whereas Recurrent Neural Networks (RNN) dominate the processing of sequential data such as text and voice (Smagulova and James 2019). DNNs allow machines to automatically discover the representations needed for the detection or classification of raw input (LeCun et al. 2015). Additionally, the neural network community developed unsupervised algorithms to help with the learning of unlabeled data. These unsupervised methods have found their way to real-world applications, such as creating generative adversarial networks (GANs) that design clothes (Singh et al. 2020). The term

alireza.ghods@wsu.edu .

Conflict of interest

The authors declare that they have no conflict of interest.

deep has been used to distinguish these networks from shallow networks with only one hidden layer; in contrast, DNNs have multiple hidden layers. The two terms *deep learning* and *deep neural networks* have been used synonymously. However, we observe that deep learning itself conveys a broader meaning, which can also shape the field of machine learning outside the realm of neural network algorithms.

The remarkable recent DNN advances were made possible by the availability of massive amounts of computational power and labeled data. However, these advances do not overcome all of the difficulties associated with DNNs. For example, there are many real-world scenarios, such as analyzing power distribution data (Tang et al. 2018), for which large annotated datasets do not exist due to the complexity and expense of collecting data. While applications like clinical interpretations of medical diagnoses require that the learned model be understandable, most DNNs resist interpretation due to their complexity (Caruana et al. 2015). DNNs can be insensitive to noisy training data (Nguyen et al. 2015; Zhang et al. 2017; Krueger et al. 2017), and they also require appropriate parameter initialization to converge (Sutskever et al. 2013; Mishkin and Matas 2016).

Despite these shortcomings, DNNs have reported higher predictive accuracy than other supervised learning methods for many datasets, given enough supervised data and computational resources. Deep models offer structural advantages that may improve the quality of learning in complex datasets as empirically shown by Bengio (2009). Recently, researchers have designed hybrid methods which combine unique DNN techniques with other classifiers to address some of these identified problems or to boost other classifiers. This survey paper investigates these methods, reviewing classifiers that have adapted DNN techniques to alternative classifiers.

1.1 Survey Objectives and Outline

While DNN research is growing rapidly, this paper aims to draw a broader picture of deep learning methods. Although some studies provide evidence that DNN models offer greater generalization than classic machine learning algorithms for complex data (Szegedy et al. 2015; Wu et al. 2016; Józefowicz et al. 2016; Graves et al. 2013; Ji et al. 2013), there is no “silver bullet” approach to concept learning (Wolpert and Macready 1997). Numerous studies comparing DNNs and other supervised learning algorithms (King et al. 1995; Lim et al. 2000; Caruana and Niculescu-Mizil 2006; Caruana et al. 2008; Baumann et al. 2019) observe that the choice of algorithm depends on the data - no ideal algorithm exists which generalizes optimally on all types of data. Recognizing the unique and important role other classifiers thus play, we aim to investigate how non-neural network machine learning algorithms can benefit from the advances in deep neural networks. Many deep learning survey papers have been published that provide a primer on the topic (Pouyanfar et al. 2019) or highlight diverse applications such as object detection (Shickel et al. 2018), medical record analysis (Han et al. 2018), activity recognition (Wang et al. 2019b), and natural language processing (Hatcher and Yu 2018). In this survey, we do not focus solely on deep neural network models but rather on how deep learning can inspire a broader range of classifiers. We concentrate on research breakthroughs that transform non-neural network classifiers into deep learners. Further, we review deep network techniques such as stochastic

gradient descent that can be used more broadly, and we discuss ways in which non-neural network models can benefit from network-inspired deep learning innovations.

The literature provides evidence that non-neural network models may offer improved generalizability over deep networks, depending on the amount and type of data that is available. By surveying methods for transforming non-neural network classifiers into deep learners, these approaches can become stronger learners. To provide evidence of the need for continued research on this topic, we also implement a collection of shallow and deep learners surveyed in this paper, both network and non-neural network classifiers, to compare their performance. Figure 1 highlights the deep learning components that we discuss in this survey. This graph also summarizes the deep classifiers that we survey and the relationships that we highlight between techniques.

2 Brief Overview of Deep Neural Networks

2.1 The Origin

In 1985, Rosenblatt introduced the Perceptron (Rosenblatt 1958), an online binary classifier which flows input through a weight vector to an output layer. Perceptron learning uses a form of gradient descent to adjust the weights between the input and output layers to optimize a loss function (Widrow and Hoff 1960). A few years later, Minsky proved that a single-layer Perceptron is unable to learn nonlinear functions, including the XOR function (Minsky and Papert 1987). Multilayer perceptrons (MLPs, see Table 3 for a complete list of abbreviations) addressed the nonlinearity problem by adding layers of hidden units to the networks and applying alternative differentiable activation functions, such as sigmoid, to each node. Stochastic gradient descent was then applied to MLPs to determine the weights between layers that minimize function approximation errors (Rumelhart et al. 1985). However, the lack of computational power caused DNN research to stagnate for decades, and other classifiers rose in popularity. In 2006, a renaissance began in DNN research, spurred by the introduction of Deep Belief Networks (DBNs) (Hinton et al. 2006).

2.2 Deep Neural Network Architectures

Due to the increasing popularity of deep learning, many DNN architectures have been introduced with variations such as Neural Turing Machines (Graves et al. 2014) and Capsule Neural Networks (Sabour et al. 2017). In this paper, we summarize the general form of DNNs together with architectural components that not only appear in DNNs but can be incorporated into other models. We start by reviewing popular types of DNNs that have been introduced and that play complementary learning roles.

2.3 Supervised Learning

2.3.1 Multilayer Perceptron—A multilayer perceptron (MLP) is one of the essential bases of many deep learning algorithms. The goal of a MLP is to map input X to class y by learning a function $y = f(X, \theta)$, where θ represents the best possible function approximation. For example, in Figure 2 the MLP maps input X to y using function $f(x) = f^{(3)}(f^{(2)}(f^{(1)}(x)))$, where $f^{(1)}$ is the first hidden-layer, $f^{(2)}$ is the hidden-second layer, and $f^{(3)}$ represents the third, output layer. This chain structure is a common component of many DNN architectures.

The network depth is equal to the length of the chain, and the width of each layer represents the number of nodes in that layer (Goodfellow et al. 2016).

In networks such as the MLP, the connections are not cyclic and thus belong to a class of DNNs called *feedforward networks*. Feedforward networks move information in only one direction, from the input to the output layer. Figure 2 depicts a particular type of feedforward network which is a fully-connected multilayer perceptron because each node at one layer is connected to all of the nodes at the next layer. Special cases of feedforward networks and MLPs have drawn considerable recent attention, which we describe next.

2.3.2 Deep Convolutional Neural Network—A convolutional neural network (CNN) (LeCun et al. 1989) is a specialized class of feedforward DNNs for processing data that can be discretely presented. Examples of data that can benefit from CNNs include time series data that can be presented as samples of discrete regular time intervals and image data presented as samples of 2-D pixels at discrete locations. Most CNNs involve three stages: a convolution operation; an activation function, such as the rectified linear activation (ReLU) function (Krizhevsky et al. 2012); and a pooling function, such as max pooling (Zhou and Chellappa 1988). A convolution operation is a weighted average or smooth estimation of a windowed input. One of the strengths of the convolution operation is that the connections between nodes in a network become sparser by learning a small kernel for unimportant features. Another benefit of convolution is parameter sharing. A CNN makes an assumption that a kernel learned for one input position can be used at every position, in contrast to an MLP, which deploys a separate element of a weight matrix for each connection. Applying the convolution operator frequently improves the network's learning ability.

A pooling function replaces the output of specific nearby nodes by their statistical summary. For example, the max-pooling function returns the maximum of a rectangular neighborhood. The motivation behind adding a pooling layer is that statistically down-sampling the number of features makes the representation approximately invariant to small translations of the input by maintaining the essential features. The final output of the learner is generated via a Fully-Connected (FC) layer that appears after the convolutional and max-pooling layers (see Figure 3 for an illustration of the process).

2.3.3 Recurrent Neural Network—A recurrent Neural Network (RNN) is a sequential model that can capture the relationship between items in a sequence. Unlike traditional neural networks, wherein all inputs are independent, RNNs contain artificial neurons with one or more feedback loops. Feedback loops are recurrent cycles over time or sequence, as shown in Figure 4. An established RNN problem is exploding or vanishing gradients. For a long data sequence, the gradient could become increasingly smaller or increasingly larger, which halts the learning. To address this issue, Hochreiter and Schmidhuber (1997) introduced a long short-term memory (LSTM) model and Cho et al. (2014) proposed a gated recurrent unit (GRU) model. Both of these networks allow the gradient to flow unchanged in the network, thus preventing exploding or vanishing gradients.

2.3.4 Siamese Neural Network—There are settings where the number of training samples is limited, such as in facial recognition scenarios where only one image is available

per person. When there is a limited number of examples for each class, DNNs struggle with generalizing the model. One strategy for addressing this problem is to learn a similarity function. This function computes the degree of difference between two samples instead of learning each class. As an example, let x_1 represents one facial image and x_2 represents a second. If $d(x_1, x_2) \leq \tau$, we can conclude that the images are of the same person while $d(x_1, x_2) > \tau$ implies that they are different people. Siamese Neural Networks (SNN) (Taigman et al. 2014) build on this idea by encoding examples x_i and x_j on two separate DNNs with shared parameters. The SNN learns a function d using encoded features, as shown in Figure 5. The network then outputs $y > 0$ for similar objects (i.e., when d is less than a threshold value) and $y < 0$ otherwise. Thus, SNNs can be used for similarity learning by learning a distance function over objects. In addition to their value for supervised learning from limited samples, SNNs are also beneficial for unsupervised learning tasks (Riad et al. 2018; Alaverdyan et al. 2020).

2.4 Unsupervised Learning

2.4.1 Generative Adversarial Network—Until this point in the survey, we have focused on deep learning for its power in classifying data points. However, researchers have exploited deep learning for other uses as well, such as generating synthetic data that shares characteristics with known real data.

One way to create synthetic data is to learn a generative model. A generative model learns the parameters that govern a distribution based on observation of real data points from that distribution. The learned model can then be used to create arbitrary amounts of synthetic data that emulate real data observations. Recently, researchers have found a way to exploit multiplayer games for the purpose of improving generative machine learning algorithms. In the adversarial training scenario, two agents compete against each other, as inspired by Samuel (1959) who designed a computer program to play checkers against itself. Goodfellow et al. (2014) put this idea to use when developing Generative Adversarial Networks (GANs), in which a DNN (generator) tries to generate synthetic data that is so similar to real data that it fools its opponent DNN (discriminator), whose job is to distinguish real from fake data (see Figure 6 for an illustration). The goal of GANs is to simultaneously improve the ability of the generator to produce realistic data and of the discriminator to distinguish synthetic from real data. GANs have found successful application in diverse tasks, including translating text to images (Reed et al. 2016), discovering drugs (Kadurin et al. 2017), and transforming sketches to images (Chen and Hays 2018; Park et al. 2019).

2.4.2 Autoencoder—Yet another purpose for deep neural networks is to provide data compression and dimensionality reduction. An Autoencoder (AE) is a DNN that accomplishes this goal by creating an output layer that resembles the input layer, using a reduced set of terms represented by the middle layers (Goodfellow et al. 2016). Architecturally, an AE combines two networks. The first network, called the encoder, learns a new representation of input x with fewer features $h = f(x)$; the second part, called the decoder, maps h onto a reconstruction of the input space $\hat{y} = g(h)$, as shown in Figure 7. The goal of an AE is not simply to recreate the input features. Instead, an

AE learns an approximation of the input features to identify useful properties of the data. AEs are vital tools for dimensionality reduction (Hinton and Salakhutdinov 2006), feature learning (Vincent et al. 2008), image colorization (Zhang et al. 2016), higher-resolution data generation (Huang et al. 2018), and latent space clustering (Yeh et al. 2017). Additionally, other versions of AEs such as variational autoencoders (VAEs) (Kingma and Welling 2014) can be used as generative models.

2.5 Optimization for Training Deep Neural Networks

In the previous section, we described common DNN architecture components. In this section, we offer a brief overview of optimization approaches for training DNNs. Learning methods may optimize a function $f(x)$ (e.g., minimize a loss function) by modifying model parameters (e.g., changing DNN weights). However, as Bengio (2013) point out, DNN optimization during training may be further complicated by local minima and ill-conditioning (see Figure 8 for an illustration of an ill-condition).

The most common type of optimization strategy employed by DNNs is gradient descent. This intuitive approach computes the error derivative with respect to a higher-level layer of the network to learn the weights of connections between layers, which reduces the network's objective function. Input x is fed forward through a network to predict \hat{y} . A cost function $\mathcal{J}(\theta)$ measures the error of the network at the output layer. Gradient descent then directs the cost value to flow backward through the network by computing the gradient of the objective function $\nabla_{\theta}\mathcal{J}(\theta)$. This process is sometimes alternatively referred to as backpropagation because the training error propagates backward through the network from output to input layers. Many variations of gradient descent have been tested for DNN optimization, such as stochastic gradient descent, mini-batch gradient descent, momentum (Sutskever et al. 2013), Ada-Grad (Duchi et al. 2011), and Adam (Kingma and Ba 2015).

Deep network optimization is an active area of research. Along with gradient descent, many other algorithms such as derivative-free optimization (Rios and Sahinidis 2013) and feedback-alignment (Nøklund 2016) have appeared. However, none of these algorithms are as popular as the gradient descent algorithms.

2.6 Regularization

Regularization was an optimization staple for decades prior to the development of DNNs. The rationale behind adding a regularizer to a classifier is to avoid the overfitting problem, where the classifier fits the training set too closely instead of generalizing to the entire data space. Goodfellow et al. (2016) defined regularization as “any modification to a learning algorithm that is intended to reduce its generalization error but not its training error”. While regularization methods such as bagging have been popular for neural networks and other classifiers, recently, the DNN community has developed novel regularization methods that are unique to deep neural networks. In some cases, backpropagation training of fully-connected DNNs results in poorer performance than shallow structures because the deeper structure is prone to being trapped in local minima and overfitting the training data (Zhang et al. 2017). To improve the generalizability of DNNs, regularization methods have

thus been adopted during training. Here we review the intuition behind the most frequent regularization methods that are currently found in DNNs.

2.6.1 Parameter Norm Penalty—A conventional method for avoiding overfitting is to penalize large weights by adding a p-norm penalty function to the optimization function of the form $f(x)+p\text{-norm}(x)$, where the p-norm p for weights w is denoted as $\|w\|_p = (\sum_i |w_i|^p)^{\frac{1}{p}}$. Popular p-norms are the L_1 and L_2 norms which have been used by other classifiers such as logistic regression and SVMs prior to the introduction of DNNs. L_1 adds a regularization term $\mathcal{Q}(\theta) = \|w\|_1$ to the objective function for weights w , while L_2 adds a regularization term $\mathcal{Q}(\theta) = \|w\|_2$. The difference between the L_1 and L_2 norm penalty functions is that L_1 penalizes features more heavily by setting the corresponding edge weights to zero compared to L_2 . Therefore, a classifier with the L_1 norm penalty tends to prefer a sparse model. The L_2 norm penalty is more common than the L_1 norm penalty. However, it is often advised to use the L_1 norm penalty when the amount of training data is small and the number of features is large to avoid noisy and less-important features. Because of its sparsity property, the L_1 penalty function is a key component of LASSO feature selection (Tibshirani 1996).

2.6.2 Dropout—A powerful method to reduce generalization error is to create an ensemble of classifiers. Multiple models are trained separately, then as an ensemble they output a combination of the models' predictions on test points. Some examples of ensemble methods included bagging (Breiman 1996), which trains k models on k different folds of random samples with replacement and boosting (Freund 1995), which applies a similar process to weighted data. A variety of DNNs use boosting to achieve lower generalization error (Hinton et al. 2006; Moghimi et al. 2016; Eickholt and Cheng 2013).

Dropout (Srivastava et al. 2014) is a popular regularization method for DNNs, which can be viewed as a computationally-inexpensive application of bagging to deep networks. A common way to apply dropout to a DNN is to deactivate a randomly-selected 50% of the hidden nodes and a randomly-selected 20% of the input nodes for each mini-batch of data. The difference between bagging and dropout is that in bagging, the models are independent of each other, while in dropout, each model inherits a subset of parameters from the parent deep network.

2.6.3 Data Augmentation—DNNs can generalize better when they have more training data; however, the amount of available data is often limited. One way to circumvent this limitation is to generate artificial data from the same distribution as the training set. Data augmentation has been particularly effective when used in the context of classification. The goal of data augmentation is to generate new training samples from the original training set (X, y) by transforming the X inputs. Data augmentation may include generating noisy data to improve robustness (denoising) or creating additional training data for the purpose of regularization (synthetic data generation). Dataset augmentation has been adopted for a variety of tasks such as image recognition (Perez and Wang 2017; Cubuk et al. 2018), speech recognition (Jaitly and Hinton 2013), and activity recognition (Ohashi et al. 2017). Additionally, GANs (Bowles et al. 2018; Antoniou et al. 2017) and AEs (Jorge et al. 2018;

Liu et al. 2018), described in Sections 2.4.1 and 2.4.2, can be employed to generate such new examples.

Injecting noise into a copy of the input is another data augmentation method. Although DNNs are not consistently robust to noise (Tang and Eliasmith 2010), Poole et al. (2014) show that DNNs can benefit from carefully-tuned noise.

3 Deep Learning Architectures Outside of Deep Neural Networks

Recent research has introduced numerous enhancements to the basic neural network architecture that enhance network classification power, particularly for deep networks. In this section, we survey non-neural network classifiers that also make use of these advances.

3.1 Supervised Learning

3.1.1 Feedforward Learning—A DNN involves multiple layers of operations that are performed sequentially. The idea of creating a sequence of operations, each of which manipulates the data before passing them to the next operator, may be used to improve many types of classifiers. One way to construct a model with a deep feedforward architecture is to use stacked generalization (Wolpert 1992; Ting and Witten 1999). Stacked generalization classifiers are comprised of multiple layers of classifiers stacked on top of each other, as found in DNNs. In stacked generalization classifiers, one layer generates the next layer's input by concatenating its own input to its output. Stacked generalization classifiers typically only implement forward propagation, in contrast to DNNs, which propagate information both forward and backward through the model.

In general, learning methods that employ stacked generalization can be categorized into two strategies. In the first stacked generalization strategy, the new feature space for the current layer comes from the concatenation of the predicted output of the previous layer with the original feature vector. Here, layers refer not to layers of neural network operations but instead refer to sequences of other types of operations. Examples of this strategy include Deep Forest (DF) (Zhou and Feng 2017) and the Deep Transfer Additive Kernel Least Square SVM (DTA-LS-SVM) (Wang et al. 2019a). At any given layer, for each instance x , DF extends x 's previous feature vector to include the previous layer's predicted class value for the instance. The prediction represents a distribution over class values, averaged over all trees in the forest. Furthermore, Zhou and Feng (2017) introduce a method called Multi-Grained Scanning for improving the accuracy of DFs. Inspired by CNNs and RNNs where spatial relationships between features are critical, Multi-Grained Scanning splits a D -dimensional feature vector into smaller segments by moving a window over the features. For example, given 400 features and a window size of 100, the original features convert to 301 features of length 100, $\{ \langle 1-100 \rangle, \langle 2-101 \rangle, \dots, \langle 301-400 \rangle \}$, where the new instances have the same labels as the original instances. The new samples, described by a subset of the original features, might have incorrectly-associated labels. At first glance, it seems these noisy data could hurt the generalization. But as Breiman (2000) illustrates, perturbing a percentage of the training labels can actually help generalization.

Furthermore, Ho (1995) demonstrates that feature sub-sampling can enhance the generalization capability for RFs. Zhou and Feng (2017) tested three different window sizes ($D/4$, $D/8$, and $D/16$), where data from each different window size fits a different level of a DF model. Then the newly-learned representation from these three layers is fed to a multilayer DF. If the transformed features are too long, Zhou and Feng (2017) apply feature sub-sampling. Multi-Grained Scanning can improve the performance of a DF model for continuous data, as Zhou and Feng (2017) report that accuracy increased by 1.24% on the MNIST (LeCun 1998) dataset. An alternative method, DTA-LS-SVM, applies an Additive Kernel Least Squares SVM (AK-LS-SVM) (Cawley 2006; Yang and Wu 2012) at each layer and concatenates the original feature vector x with the prediction of the previous level to feed to the next layer. In addition, DTA-LS-SVM incorporates a parameter-transfer approach between the source (previous-layer learner) and target (next-layer learner) to enhance the classification capability of the higher level.

In the second stacked generalization strategy, the current layer's new feature space comes from the concatenation of predictions from all previous layers with the original input feature vector. Examples of this strategy include the Deep SVM (D-SVM) (Abdullah et al. 2009) and the Random Recursive SVM (R2-SVM) (Vinyals et al. 2012). The D-SVM contains multiple layers of SVMs, where the first layer is trained in the normal fashion. Following this step, each successive layer employs the kernel activation from the previous layer with the desired labels. The R2-SVM is a multilayer SVM model which at each layer transforms the data based on the sigmoid of a projection of all previous layers' outputs. For the data (X, Y) where $X \in R^D$ and $Y \in R^C$, the random projection matrix is $W \in R^{D \times C}$, where each element is sampled from $N(0, 1)$. The input data for the next layer is:

$$X_{l+1} = \sigma\left(d + \beta W_{l+1} [o_1^T, o_2^T, \dots, o_l^T]^T\right), \quad (1)$$

where β is a weight parameter that controls the degree with which a data sample in X_{l+1} moves from the previous layer, $\sigma(\cdot)$ is the sigmoid function, W_{l+1} is the concatenation of l random projection matrices $[W_{l+1,1}, W_{l+1,2}, \dots, W_{l+1,l}]$, one for each previous layer, and o is the output of each layer. Addition of a sigmoid function to the recursive model prevents deterioration to a trivial linear model in a similar fashion as MLPs. The purpose of the random projection is to push data from different classes in different directions.

It is important to note here that stacked generalization can be found in DNNs as well as non-neural network classifiers. Examples of DNNs with stacked generalization include Deep Stacking Networks (Deng et al. 2012; Hutchinson et al. 2013) and Convex Stacking Architectures (Yu and Deng 2011; Deng et al. 2012). This is clearly one enhancement that benefits all types of classifier strategies. However, there is no evidence that stack generalization could add nonlinearity to the model.

DNN classifiers learn a new representation of data at each layer with a goal that the newly-learned representation maximally separates the classes. Unsupervised DNNs often share this goal. As an example, Deep PCA's model (Liong et al. 2013) is made of two layers that each learn a new data representation by applying a Zero Components Analysis (ZCA) whitening filter (Krizhevsky and Hinton 2009) followed by a principal components analysis

(PCA) (Shlens 2014). The final data representation is derived from concatenating the output of the two layers. The motivation behind applying a ZCA whitening filter is to force the model to focus on higher-order correlations. One motivation for combining output from the first and second layers could be to preserve the learned representation from the first layer and to prevent feature loss after applying PCA at each layer. Experiments demonstrate that Deep PCA exhibits superior performance for face recognition tasks compared to standard PCA and a two-layer PCA without a whitening filter. However, as empirically confirmed by Damianou and Lawrence (2013), stacking PCAs does not necessarily result in an improved representation of the data because Deep PCA is unable to learn a nonlinear representation of data at each layer. Damianou and Lawrence (2013) fed a Gaussian to a Deep PCA and observed that the model learned just a lower rank of the input Gaussian at each layer.

As pointed out earlier in this survey, the invention of the deep belief net (DBN) (Hinton et al. 2006) drew the attention of researchers to developing deep models. A DBN can be viewed as a stacked restricted Boltzmann machine (RBM), where each layer is trained separately and alternates functionality between hidden and input units. In this model, features learned at hidden layers then represent inputs to the next layer. An RBM is a generative model that contains a single hidden layer. Unlike the Boltzmann machine, hidden units in the restricted model are not connected to each other and contain undirected, symmetrical connections from a layer of visible units (inputs). All of the units in each layer of an RBM are updated in parallel by inputting the current state of the unit to the other layer. This updating process repeats until the system is sampling from an equilibrium distribution. The RBM learning rule is shown in Equation 2.

$$\frac{\partial \log P(v)}{\partial W_{ij}} \approx \langle v_i h_j \rangle_{data} - \langle v_i h_j \rangle_{reconstruction} \quad (2)$$

In this equation, W_{ij} represents the weight vector between a visible unit v_j and a hidden unit h_j , and $\langle . \rangle$ is the average value over all training samples. Since the introduction of DBNs, many other different variations of Deep RBMs have been proposed, such as temporal RBMs (Sutskever and Hinton 2007), gated RBMs (Memisevic and Hinton 2007), and cardinality RBMs (Swersky et al. 2012).

Another novel form of a deep belief net is a deep Gaussian process (DGP) model (Damianou and Lawrence 2013). DGP is a deep directed graph where multiple layers of Gaussian processes map the original features to a series of latent spaces. DGPs offer a more general form of Gaussian Processes (GPs) (Rasmussen 2003) where a one-layer DGP consists of a single GP, f . In a multilayer DGP, each GP, f_p , maps data from one latent space to the next. As shown in Equation 3, each data point Y is generated from the corresponding function f_j with ϵ Gaussian noise applied to data X_j that is obtained from a previous layer.

$$Y = f_l(X_l) + \epsilon_l, \epsilon_l \sim \mathcal{N}(0, \sigma_l^2 I) \quad (3)$$

Figure 9 illustrates a DGP expressed as a series of Gaussian processes mapping data from one latent space to the next. Functions f_j are drawn from a Gaussian process, i.e. $f(x) \sim \mathcal{GP}(0, k(x, x'))$. In this setting, the covariance function k defines the properties of the

mapping function. DGP can be utilized for both supervised and unsupervised learning. In the supervised setting, the top hidden layer is observed, whereas in the unsupervised setting, the top hidden layer is set to a unit Gaussian as a fairly uninformative prior. DGP is a powerful non-parametric model, but it has only been tested on small datasets. Also, we note that researchers have developed deep Gaussian process models with alternative architectures such as recurrent Gaussian processes (Mattos et al. 2016), convolutional Gaussian processes (van der Wilk et al. 2017) and variational auto-encoded deep Gaussian processes (Dai et al. 2016). There exists a vast amount of literature on this topic that provides additional insights on deep Gaussian processes (Duvenaud et al. 2014; Damianou 2015; Dunlop et al. 2018).

As we discussed, non-neural network classifiers have been designed that contain multiple layers of operations, similar to a DNN. We observe that a common strategy for creating a deep non-neural network model is to add the prediction of the previous layer or layers to the original input feature. Likewise, novel methods can be applied to learn a new representation of data at each layer. We discuss these methods next.

3.1.2 Siamese Model—As discussed in Section 2.3.4, an SNN represents a powerful method for similarity learning. However, one problem with SNNs is overfitting when there is a small number of training examples. The Siamese Deep Forest (SDF) (Utkin and Ryabinin 2018) is a method based on DF which offers an alternative to a standard SNN. The SDF, unlike the SNN, uses only one DF. The first step in training an SDF is to modify the training examples. The training set consists of the concatenation of each pair of samples in the original set. If the sample points x_i and x_j are semantically similar, the corresponding class label is set to zero; otherwise, the class label is set to one. The difference between the SDF and the DF in training is that the Siamese Deep Forest concatenates the original feature vector with a weighted sum of the tree class probabilities. Training of SDF is similar to DF; the primary difference is that SDF learns the class probability weights w for each forest separately at each layer. Learning the weights for each forest can be accomplished by minimizing the function in Equation 4.

$$\min_w J_q(w) = \min_w \sum_{i,j} l(x_i, x_j, y_{ij}, w) + \lambda R(w) \quad (4)$$

Here, w represents a concatenation of vectors w^k , $k = 1, \dots, M$, q is the SDF layer, $R(w)$ is a regularization term, and λ is a hyper-parameter to control regularization. Detailed instructions on minimizing Equation 4 are found in the literature (Utkin and Ryabinin 2018). The results of SDF experiments indicate that the SDF can achieve better classification accuracy than DF for small datasets. In general, all non-neural network models that learn data representations can take advantage of the Siamese architecture like SDF.

3.2 Unsupervised Learning

3.2.1 Generative Adversarial Model—A common element found in GANs is the inclusion of an FC layer in the discriminator. One issue with the FC layer is that it cannot deal with the ill-condition in which local minima are not surrounded by spherical wells, as shown in Figure 8. The Generative Adversarial Forest (GAF) (Zuo et al. 2018) replaces the FC layer of the discriminator with a deep neural decision forest (DNDF), which is discussed

in Section 4. GAF and DNDF are distinguished based on how leaf node values are learned. Instead of learning leaf node values iteratively, as DNDF does, GAF learns them in parallel across the ensemble members. The strong discriminatory power of the decision forest is the reason the authors recommend this method in lieu of the fully-connected discriminator layer.

In this previous work, the discriminator is replaced by an unconventional model. We hypothesize that replacing the discriminator with other classifiers such as Random Forest, SVM, or K-nearest neighbor based on the data could result in a diverse GAN strategy, each of which may offer benefits for alternative learning problems.

3.2.2 Autoencoder—As we discussed in Section 2.4.2, AEs offer strategies for dimensionality reduction and data reconstruction from compressed information. The autoencoding methodology can be found in neural networks, non-neural networks, and hybrid methods. As an example, the multilayer SVM (ML-SVM) autoencoder is a variation of ML-SVM with the same number of output nodes as input features and a single hidden layer that consists of fewer nodes than the input features. ML-SVM is a model with the same structure as an MLP. The distinction here is that the network contains SVM models as its nodes. A review of ML-SVM is discussed in Section 4. The outputs of hidden nodes are fed as input to each SVM output node c as follows:

$$g_c(f(X | \theta)) = \sum_{i=1}^l (\alpha_i^{c*} - \alpha_i^c) K_o(f(x_i | \theta), f(x | \theta)) + b_c, \quad (5)$$

where α_i^{c*} and α_i^c are the support vector coefficients, K_o is the kernel function, and b_c is their bias. The error backpropagates through the network to update the parameters.

Another exciting emerging research area is the combination of Kalman filters with deep networks. A Kalman filter is a well-known algorithm that estimates the optimal state of a system from a series of noisy observations. The classical Kalman filter (Kalman 1960) is a linear dynamical system and therefore is unable to model complex phenomena. For this reason, researchers developed nonlinear versions of Kalman filters. In a seminal contribution, Krishnan et al. (2015) introduced a model that combines a variational autoencoder with Kalman filters for counterfactual inference of patient information. In a standard autoencoder, the model learns a latent space that represents the original data minus extraneous information or “signal noise”. In contrast, a variational autoencoder (VAE) (Kingma and Welling 2014) adds a constraint to the encoder that learns a Gaussian distribution of the original input data. Therefore, a VAE is able to generate a latent vector by sampling from the learned Gaussian distribution. Deep Kalman filters (DKF) learn a generative model from observed sequences $\mathbf{x} = (x_1, \dots, x_T)$ and actions $\mathbf{u} = (u_1, \dots, u_{T-1})$, with a corresponding latent space $\mathbf{z} = (z_1, \dots, z_T)$, as follows:

$$\begin{aligned} z_1 &\sim \mathcal{N}(\mu_0, \Sigma_0) \\ z_t &\sim \mathcal{N}(G_\alpha(z_{t-1}, u_{t-1}, \Delta_t), S_\beta(z_{t-1}, y_{t-1}, \Delta_t)) \\ x_t &\sim \Pi(F_k(z_t)), \end{aligned} \quad (6)$$

where $\mu_0 = 0$ and $\Sigma_0 = I_{d_t}$ represents the difference between times t and $t - 1$, and Π represents a distribution (e.g., Bernoulli for binary data) over observation x_t . The functions $G_\alpha, S_\beta, F_\kappa$ are parameterized by a neural net. As a result, the autoencoder will learn $\theta = \{\alpha, \beta, \kappa\}$ parameters. Additionally, Shashua and Mannor (2017) introduced deep Q-learning with Kalman filters and Lu et al. (2018) presented a deep Kalman filter model for video compression.

As we highlighted in this section, non-neural network methods have been designed that are inspired by AEs. Although ML-SVM mimics the architecture of AEs, its computational cost prevents the algorithm from being a practical choice. DKF takes advantage of the VAE idea by learning a Kalman Filter in its middle layer. Additionally, Feng and Zhou (2018) introduced an encoder forest, a model inspired by the DNN autoencoder. Because the encoder forest is not a deep model, we do not include the details of this algorithm in our survey.

4 Deep Learning Optimization Outside of Deep Neural Networks

As discussed in Section 2.5, gradient descent has been a prominent optimization algorithm for DNNs; however, it has been underutilized by non-neural network classifiers. Some notable exceptions are found in the literature. We discuss these here.

A resourceful method for constructing a deep model is to start with a DNN architecture and then replace nodes with non-neural network classifiers. As an example, the multilayer SVM (ML-SVM) (Wiering and Schomaker 2014) replaces nodes in an MLP with standard SVMs. ML-SVM is a multiclass classifier which contains SVMs within the network. At the output layer, the ML-SVM contains the same number of SVMs as the number of classes learned at the perceptron output layer. Each SVM at the ML-SVM output layer is trained in a one-versus-all fashion for one of the classes. When observing a new data point, ML-SVM outputs the class label corresponding to the SVM that generates the highest confidence. At each hidden layer, SVMs are associated with each node that learns latent variables. These variables are then fed to the output layer. At hidden layer $f(X|\theta)$ where X is the training set and θ denotes the trainable parameters of the SVM, ML-SVM maps the hidden layer features to an output value as follows:

$$g(f(X|\theta)) = \sum_{i=1}^l y_i^c \alpha_i^c K_o(f(x_i|\theta), f(X|\theta)) + b_c, \quad (7)$$

where g is the output layer function, $y_i^c \in \{-1, 1\}$ for each class c , K_o is the kernel function for the output layer, α_i^c are the support vector coefficients for SVM nodes of the output layer, and b_c is their bias. The goal of ML-SVM is to learn the maximum support vector coefficient of each SVM at the output layer with respect to the objective function $J_c(\cdot)$, as shown in Equation 8.

$$\min_{w^c, b, \xi, \theta} J_c = \frac{1}{2} \|w^c\|^2 + C \sum_i \xi_i \quad (8)$$

Here, w^c represents the set of weights for class c , C represents a trade-off between margin width and misclassification risk and ξ_i are slack variables. ML-SVM applies gradient ascent to adapt its support vector coefficient towards a local maximum of $J_C(\cdot)$. The support vector coefficient is defined as zero for values less than zero and is assigned to C for values larger than C . The data is backpropagated through the network in a way that is similar to traditional MLPs, by calculating the gradient of the objective function.

The SVMs in the hidden layer are identical. Given the same inputs, they would thus generate the same outputs. The hidden layers train on a perturbed version of the training set to eliminate producing similar outputs before training the combined ML-SVM model to diversify the SVMs. The outputs of hidden layer nodes are constrained to generate values in the range $[-1, 1]$. Despite the effort of ML-SVMs to learn a multi-layer data representation, this approach is currently not practical because adding a new node incurs a dramatic computational expense for large datasets.

Kontschieder et al. (2015) further incorporate gradient descent into a Random Forest (RF), which is a popular classification method. One of the drawbacks of an RF is that it does not traditionally learn new internal representations like DNNs. The Deep Network Decision Forest (DNDF) (Kontschieder et al. 2015) integrates a DNN into each decision tree within the forest to reduce the uncertainty at each decision node. In DNDF, the result of a decision node $d_n(x, \Theta)$ corresponds to the output of a DNN $f_n(x, \Theta)$, where x is an input and Θ is the parameter of a decision node. DNDF must have differentiable decision trees to apply gradient descent to the process of updating decision nodes. In a standard decision tree, the result of a decision node $d_n(x, \Theta)$ is deterministic. DNDF replaces the traditional decision node with a sigmoid function $d_n(x, \Theta) = \sigma(f_n(x; \Theta))$ to create a stochastic decision node. The probability of reaching a leaf node l is calculated as the product of all decision node outputs from the root to the leaf l , which is expressed as μ_l in this context. The set of leaf nodes \mathcal{L} learns the class distribution π , and the class with the highest probability is the prediction of the tree. The aim of DNDF is to minimize its empirical risk with respect to the decision node parameter Θ and the class distribution π of \mathcal{L} under the log-loss function for a given data set.

The optimization of the empirical risk is a two-step process which is executed iteratively. The first step is to optimize the class distribution of leaf nodes $\pi_{\mathcal{L}}$ while fixing the decision node parameters and the corresponding DNN. At the start of optimization (iteration 0), class distribution π^0 is set to a uniform distribution across all leaves. DNDF then iteratively updates the class distribution across the leaf nodes as follows for iteration $t + 1$:

$$\pi_{ly}^{(t+1)} = \frac{1}{Z_l^{(t)}} \sum_{(x, y') \in T} \frac{\mathbb{1}_{y=y'} \pi_{ly}^{(t)} \mu_l(x | \Theta)}{\mathbb{P}_T[y | x, \Theta, \pi^{(t)}]}, \quad (9)$$

where $Z_l^{(t)}$ is a normalization factor ensuring that $\sum_y \pi_{ly}^{t+1} = 1$, $\mathbb{1}_q$ is the indicator function on the argument q , and \mathbb{P}_T is the prediction of the tree.

The second step is to optimize decision node parameters Θ while fixing the class distribution $\pi_{\mathcal{L}}$. DNDF employs gradient descent to minimize log-loss with respect to Θ as follows:

$$\frac{\partial L}{\partial \Theta}(\Theta, \pi; x, y) = \sum_{n \in \mathcal{N}} \frac{\partial L(\Theta, \pi; x, y)}{\partial f_n(x; \Theta)} \frac{\partial f_n(x; \Theta)}{\partial \Theta}. \quad (10)$$

The second term in Equation 10 is the gradient of the DNN. Because this is commonly known, we only discuss calculating the gradient of the differentiable decision tree. Here, the gradient of the differentiable decision tree is given by:

$$\frac{\partial L(\Theta, \pi; x, y)}{\partial f_n(x; \Theta)} = d_n(x; \Theta) A_{n_l} - \bar{d}_n(x; \Theta) A_{n_r}, \quad (11)$$

where d_n is the probability of transitioning to the left child, $\bar{d}_n = 1 - d_n$ is the probability of transitioning to the right child calculated by a forward pass through the DNN, and n_l and n_r indicate the left and right children of node n . To calculate the term A in Equation 11, DNDF performs one forward pass and one backward pass through the differentiable decision tree. Upon completing the forward pass, a value A_l can be initially computed for each leaf node as follows:

$$A_l = \frac{\pi_{l_y} \mu_l}{\sum_l \pi_{l_y} \mu_l}. \quad (12)$$

Next, the values of A_l for each leaf node are used to compute the values of A_m for each internal node m . To do this, a backward pass is made through the decision tree, during which the values are calculated as $A_m = A_{n_l} + A_{n_r}$, where n_l and n_r represent the left and the right children of node m , respectively.

Each layer of a standard DNN produces the output o_i at layer i . As mentioned earlier, the goal of the DNN is to learn a mapping function $F_i: o_{i-1} \rightarrow o_i$ that minimizes the empirical loss at the last layer of DNN on a training set. Because each F_i is differentiable, a DNN updates its parameters efficiently by applying gradient descent to reduce the empirical loss.

Adopting a different methodology, Frosst and Hinton (2017) distill a neural network into a soft decision tree. This model benefits from both neural network-based representation learning and decision tree-based concept explainability. In the Frosst soft decision tree (FSDT), each tree's inner node learns a filter w_i and a bias b_i , and leaf nodes l learn a distribution of classes. Like the hidden units of a neural network, each inner node of the tree determines the probability of input x at node i as follows:

$$p_i(x) = \sigma(\beta(xw_i + b_i)) \quad (13)$$

where σ represents the sigmoid function and β represents an inverse temperature whose function is to avoid soft decisions in the tree. Filter activation routes the sample x to the left branch for values of p_i less than 0.5, and to the right branch otherwise. The probability distribution Q^l for each leaf node l represents the learned parameter ϕ^l at that leaf over the possible k output classes:

$$Q_k^l = \frac{\exp(\phi_k^l)}{\sum_{k'} \exp(\phi_{k'}^l)}. \quad (14)$$

The predictive distribution over classes is calculated by traversing the greatest-probability path. To train this soft decision tree, Frosst and Hinton (2017) calculate a loss function L that minimizes the cross entropy between each leaf, weighted by input vector x path probability and target distribution T , as follows:

$$L(x) = -\log \left(\sum_{l \in \text{Leaf Nodes}} p^l(x) \sum_k T_k \log Q_k^l \right) \quad (15)$$

where $P^l(x)$ is the probability of reaching leaf node l given input x . Frosst and Hinton (2017) also introduce a regularization term to avoid internal nodes routing all data points on one particular path and encourage them to equally route data along the left and right branches. The penalty function calculates a sum over all internal nodes from the root to node i , as follows:

$$C = -\lambda \sum_{i \in \text{Inner Nodes}} 0.5 \log(\alpha_i) + 0.5 \log(1 - \alpha_i) \quad (16)$$

where λ is a hyper-parameter set prior to training to determine the effect of the penalty. The cross entropy α for a node i is the sum of the path probability $P^i(x)$ from the root to node i multiplied by the probability of that node p_i divided by the path probability, as follows:

$$\alpha_i = \frac{\sum_x P^i(x) p_i(x)}{\sum_x P^i(x)}. \quad (17)$$

Because the probability distribution is not uniform across nodes in the penultimate level, this penalty function could actually hurt the generalization. The authors address this problem by decaying the strength of penalty function λ exponentially with the depth d of the node to 2^d . Another challenge is that in any given batch of data, as the data descends the tree, the number of samples decreases exponentially. Therefore, the estimated probability loses accuracy further down the tree. Frosst and Hinton (2017) recommend addressing this problem by decaying a running average of the actual probabilities with a time window that is exponentially proportional to the depth of the nodes (Frosst and Hinton 2017). Although the authors report that the accuracy of this model was less than the deep neural network, the model offers an advantage of concept interpretability.

Both DNDF and the soft decision tree fix the depth of the learned tree to a predefined value. In contrast, Tanno et al. (2019) introduced the Adaptive Neural Tree (ANT), which can grow to any arbitrary depth. The ANT architecture is similar to a decision tree, but at each internal node and edge, ANT learns a new data representation. For example, an ANT may contain one or more convolution layers followed by a fully-connected layer at each inner node, one or more convolution layers followed by an activation function such as *ReLU* or *tanh* at each edge, and a linear classifier at each leaf node.

Training an ANT requires two phases: growth and refinement. In the growth phase, starting from the root in breadth-first order, one of the nodes is selected. The learner then evaluates three choices: 1) split the node and add a sub-tree, 2) deepen edge transformation by adding another layer of convolution, or 3) keep the current model. The model optimizes the parameters of the newly-added components by minimizing log likelihood via gradient descent while fixing the parameters of the previous portion of the tree. Eventually, the model selects the choice that yields the lowest log likelihood. This process repeats until the model converges. In the refinement phase, the model performs gradient descent on the final architecture. The purpose of the refinement phase is to correct suboptimal decisions that may have occurred during the growth phase. The authors evaluate their method on several standard testbeds, and the results indicate that ANT is competitive with many deep networks and non-neural network learners for these tasks.

In contrast to the soft decision trees, Carreira-Perpiñán and Tavallali (2018) introduce Tree Alternation Optimization (TAO), which learns a tree with linear decision nodes. Traditional decision tree algorithms such as CART (Breiman et al. 1984) and C4.5 (Salzberg 1994) create a decision tree from scratch in a way that optimizes a proxy measure such as impurity. In contrast, TAO modifies an existing tree in a way that minimizes classification error. This modification is performed incrementally, in a way that reflects the incremental adjustment of weights in a neural network. Specifically, given tree T , TAO minimizes a loss function representing the classification error resulting from all leaf nodes Θ in the tree:

$$\mathcal{L}(\Theta) = \sum_{n=1}^N L(y_n, T(x_n; \Theta)). \quad (18)$$

One advantage of TAO is that it not only learns axis-aligned trees, it can also learn oblique trees where a linear or nonlinear combination of attributes split the nodes.

Yang et al. (2018) took a different approach. They created a decision tree using a neural network. The Deep Neural Decision Tree (DNDT) employs a soft binning function to learn the split rules of the tree. DNDT constructs a one-layer neural network with softmax as its activation function. The objective function of this network is:

$$\text{softmax}\left(\frac{wx + b}{\tau}\right). \quad (19)$$

Here, for a continuous variable x , we want to bin it to $n+1$, $w = [1, 2, \dots, n+1]$ is an untrainable constant, b is a learnable bin or the cutting rule in the tree, and τ is a temperature variable. After training this model, the decision tree is constructed via the Kronecker product \otimes . Given an input $x \in R^D$ with D features, the tree rule to reach a leaf node is:

$$z = f_1(x_1) \otimes f_2(x_2) \otimes \dots \otimes f_D(x_D) \quad (20)$$

Here, z is an almost-one-hot encoded vector that indicates the index of a leaf node. One of the shortcomings of this method is that it cannot handle a high-dimensional dataset because the cost of calculating the Kronecker product becomes prohibitive. To overcome

this problem, authors learn a classifier forest by training each tree on a random subset of features.

In some cases, the mapping function is not differentiable. Feng et al. (2018) propose a new learning paradigm for training a multilayer Gradient Boosting decision tree (mGBDT) (Feng et al. 2018) where F_i is not differentiable. Gradient boosting decision tree (GBDT) is an iterative method which learns an ensemble of regression predictors. In GBDT, a decision tree first learns a model on a training set, then it computes the corresponding error residual for each training sample. A new tree learns a model on the error residuals, and by combining these two trees, GBDT is able to learn a more complex model. The algorithm follows this procedure iteratively until it meets a prespecified number of trees for training.

Since gradient descent is not applicable to mGBDT, Feng et al. (2018) obtain a “pseudo-inverse” mapping. In this mapping, G_i^t represents the pseudo-inverse of F_i^{t-1} at iteration t , such that $G_i^t(F_i^{t-1}(o_i - 1)) \sim o_i - 1$. After performing backward propagation and calculating G_i^t , forward propagation is performed by fitting a pseudo-label z_{i-1}^t from G_i^t to F_i^{t-1} . The last layer F_m computes z_m^t based on the true labels at iteration t , where $i \in \{2 \dots m\}$. After this step, pseudo-labels for previous layers are computed via pseudo-inverse mapping. To initialize mGBDT at iteration $t=0$, each intermediate (hidden) layer outputs Gaussian noise and F_i^0 represent depth-constrained trees that will later be refined. Feng et al. (2018) thus create a method that is inspired by gradient descent yet is applicable in situations where true gradient descent cannot be effectively applied.

In this section, we examine methods that apply gradient descent to non-neural network models. As we observed, one way of utilizing gradient descent is to replace the hidden units in a network with a differentiable algorithm like SVM. Another common theme we recognized was to transform deterministic decision-tree nodes into stochastic versions that offer greater representational power. Alternatively, trees or other ruled-based models can be built using neural networks.

5 Deep Learning Regularization Outside of Deep Neural Networks

We have discussed some of the common regularization methods used by DNNs in Section 2.6. Now we focus on how these methods have been applied to non-neural network classifiers in the literature. It is worth mentioning that while most models introduced in this section are not deep models, we investigate how non-neural network models can improve their performance by applying regularization methods typically associated with the deep operations found in DNNs.

5.1 Parameter Norm Penalty

Problems arise when a model is learned from data that contain a large number of redundant features. For example, selecting relevant genes associated with different types of cancer is challenging because of a large number of redundancies may exist in the gene’s long string of features. There are two common ways to eliminate redundant features: the first way is to perform feature selection and then train a classifier from the selected features; the second

way is to simultaneously perform feature selection and classification. As we discussed in Section 2.6.1, DNNs apply a L_1 or L_2 penalty function to penalize large weights. In this section, we investigate how the traditional DNN idea of penalizing features can be applied to non-neural network classifiers to simultaneously select high-ranked features and perform classification.

Standard SVMs employ the L_2 norm penalty to penalize weights in a manner similar to DNNs. However, the Newton Linear Programming SVM (NLP-SVM) (Fung and Mangasarian 2004) replaces the L_2 norm penalty with the L_1 norm penalty. This has the effect of setting small hyperparameter coefficients to zero, thus enabling NLP-SVM to select important features automatically. A different way to penalize non-important features in SVMs is to employ a Smoothly Clipped Absolute Deviation (SCAD) (Zhang et al. 2006) function. The L_1 penalty function can be biased because it imposes a larger penalty on large coefficients; in contrast, SCAD can give a nearly unbiased estimation of large coefficients. SCAD learns a non-convex penalty function as shown in Equation 21.

$$p_\lambda(|w|) = \begin{cases} \lambda|w| & \text{if } |w| \leq \lambda \\ -\frac{(|w|^2 - 2a\lambda|w| + \lambda^2)}{2(a-1)} & \text{if } \lambda < |w| \leq a\lambda \\ \frac{(a+1)\lambda^2}{2} & \text{if } |w| > a\lambda \end{cases} \quad (21)$$

SCAD equates with L_1 penalty function until $|w| = \lambda$, then smoothly transitions to a quadratic function until $|w| = a\lambda$, after which it remains a constant for all $|w| > a\lambda$. As shown by Fan and Li (2001), SCAD has better theoretical properties than the L_1 function.

One limitation of decision tree classifiers is that the number of training instances that can be selected at each branch in the tree decreases with the tree depth. This downward sampling may cause less relevant or redundant features to be selected near the bottom of the tree. To address this issue, Deng and Runger (2012) proposed to penalize features that were never selected in the process of making a tree. In a Regularized Random Forest (RRF) (Deng and Runger 2012), the information gain for a feature j is specified as follows:

$$Gain(j) = \begin{cases} \lambda \cdot Gain(j) & j \notin F \\ Gain(f_i) & j \in F \end{cases} \quad (22)$$

where F is the set of features used earlier in the path, $f_i \in F$, and $\lambda \in [0, 1]$ is the penalty coefficient. RRF avoids including a new feature j , except when the value of $Gain(j)$ is greater than $\max_i(Gain(f_i))$.

To improve RRF, Guided RRF (GRRF) (Deng and Runger 2013) assigns a different penalty coefficient λ_j to each feature instead of assigning the same penalty coefficient to all features. GRRF employs the importance score from a pre-trained RF on the training set to refine the selection of features at a given node. The importance score of feature j in an RF with T trees is the mean of gain for features in the RF. The important scores evaluate the contribution of

features for predicting classes. The GRRF uses the normalized importance score to control the degree of regularization of the penalty coefficient as follows:

$$\lambda_j = (1 - \gamma)\lambda_0 + \gamma Imp'_j, \quad (23)$$

where $\lambda_0 \in (0, 1]$ is the base penalty coefficient and $\gamma \in [0, 1]$ controls the weight of the normalized importance score. The GRRF and RRF are computationally inexpensive methods that are able to select stronger features and avoid redundant features.

5.2 Dropout

As detailed in Section 2.6.2, dropout is a method that prevents DNNs from overfitting by randomly dropping nodes during the training. Dropout can be added to other machine learning algorithms through two methods: by dropping features or by dropping models in the case of ensemble methods. Dropout has also been employed by dropping input features during training (Wang and Manning 2012, 2013). Here we look at techniques that have been investigated for dropping input features, particularly in non-neural network classifiers.

Rashmi and Gilad-Bachrach (2015) applied dropout to Multiple Additive Regression Trees (MART) (Friedman 2001, 2002). MART is a regression tree ensemble that iteratively refines its model by continually adding trees that fit the loss function derivatives from the previous version of the ensemble. Because trees added at later iterations may only impact a small fraction of the training set and thus over-specialize, researchers previously used shrinkage to exclude a random subset of leaf nodes during each tree-adding step. More recently, Rashmi and Gilad-Bachrach (2015) integrated the deep-learning idea of dropout into MART. Here, a subset of the trees is temporarily dropped. A new tree then is created based on the loss function for the on-dropped trees. This new tree is combined with the previously-dropped trees into a new ensemble. This method, Dropout Multiple Additive Regression Trees (DART) (Rashmi and Gilad-Bachrach 2015), weights the votes for the new and re-integrated trees to have the same effect on the final model output as the original set of trees. Other researchers have experimented with permanently removing a strategic subset of the dropped trees as well (Lucchese et al. 2017).

5.3 Early Stopping

The core concept of early stopping is to terminate DNN training once performance on the validation set is not improving. One potential advantage of Deep Forest (Zhou and Feng 2017) over DNNs is that DF can determine the depth of a model automatically. In DF, if the model performance does not increase on the validation set after adding a new layer, the learning terminates. Unlike DNNs, DF may avoid the tendency to overfit as more layers are added. Thus, while early stopping does not necessarily enjoy the primary outcome of preventing such overfitting, it can provide additional benefits such as shortening the validation cycle when searching for the optimal tree depth.

5.4 Data Augmentation

As discussed in Section 2.6.3, data augmentation is a powerful method for improving DNN generalization. However, little research has investigated the effects of data augmentation

methods on non-neural network classifiers. As demonstrated by Wong et al. (2016), the SVM classifier does not always benefit from data augmentation, in contrast to DNNs. However, Xu (2013) ran several data augmentation experiments on synthetic datasets and observed that data augmentation did enhance the performance of random forest classifiers. Offering explanations for the circumstances in which such augmentation is beneficial is a needed area for future research.

6 Hybrid Models

Hybrid models can be defined as a combination of two or more classes of models. There are many ways to construct hybrid models, such as DNDF (Kontschieder et al. 2015), which integrates a deep network into a decision forest, as explained in Section 4. In this section, we discuss other examples of hybrid models.

6.1 Neural Network and Decision Trees

Neural decision trees can be categorized into two groups: (1) decision trees with linear decision nodes (Carreira-Perpiñán and Tavallali 2018), and (2) soft decision trees with differentiable decision nodes (Kontschieder et al. 2015; Ioannou et al. 2016; Frosst and Hinton 2017; Tanno et al. 2019). One motivation for combining aspects of multiple models is to find a balance between classification accuracy and computational cost. Energy consumption by mobile devices and cloud servers is an increasing concern for responsive applications and green computing. Decision forests are computationally inexpensive models because of the conditional property of decision trees. Conversely, while CNNs are less efficient, they can achieve higher accuracy because of their representation-learning capabilities. Ioannou et al. (2016) introduced the Conditional Neural Network (CondNN) to reduce computation in a CNN model by introducing a routing method similar to that found in decision trees. In CondNN, each node in layer l is connected to a subset of nodes from the previous layer, $l-1$. Given a fully trained network, for every two consecutive layers, a matrix $A_{(l-1,l)}$ stores the activation values of these two layers. By rearranging elements of $A_{(l-1,l)}$ based on highly-active pairs for each class in the diagonal and zeroing out off-diagonal elements, the CondNN develops explicit routes $A_{(l,l-1)}^{route}$ through nodes in the network. CondNN incurs profoundly lower computation cost compared to other DNNs at test time; whereas, CondNN's accuracy remains similar to larger models. We note that DNN size can also be reduced by employing Bayesian optimization, as investigated by Blundell et al. (2015) and by Fortunato et al. (2017). These earlier efforts provide evidence that Bayesian neural networks can decrease network size even more than CondNNs while maintaining a similar level of accuracy.

Another motivation is to make the DNNs more interpretable. Zhao et al. (2018) replace the last layer of a deep network with a visual hierarchical tree to learn a better solution for image classification problems. A visual hierarchical tree with L levels organizes N object classes based on their visual similarities in its nodes. Deeper in the tree, groups become more separated wherein each leaf node should contain instances of one class. The class similarity between the class c_i and c_j is defined as follows:

$$S_{i,j} = S(c_i, c_j) = \exp\left(-\frac{d(x_i, x_j)}{\sigma}\right). \quad (24)$$

Here, $d(x_i, x_j)$ represents the distance between the deep representation of instances of classes c_i and c_j , and σ is automatically determined by a self-tuning technique. After calculating matrix S , hierarchical clustering is employed to learn a visual hierarchical tree.

In a traditional visual hierarchical tree, some objects might be assigned to incorrect groups. A level-wise mixture model (LMM) (Zhao et al. 2018) aims to improve this visual hierarchical tree by learning a new representation of data via a DNN, then updating the tree during training. For a given tree, matrix Ψ_{y_i, t_i} denotes the probability of objects with label y belonging to group t in the tree. First, LMM updates the DNN parameters and the visual hierarchical tree as is done with a traditional DNN. The only difference is a calculation of two gradients, one based on the parameters of the DNN and the other based on the parameters of the tree. Second, LMM updates the matrix Ψ_{y_i, t_i} for each training sample separately and then updates the parameters of the DNN and the tree. To update the Ψ , the posterior probability of the assigning group t_j for the object x_j is calculated based on the number of samples having the same label y as the label of x_j in a group t . For a given test image, LMM learns a new representation of the image based on the DNN and then obtains a prediction by traversing the tree. One of the advantages of an LMM is that, over time, by learning a better representation of data via DNN, the algorithm can update the visual hierarchical tree.

6.2 Neural Networks and K-nearest Neighbors

Another direction for blending a deep network with a non-neural network classifier is to improve the non-neural network model by learning a better representation of data via a deep network. Zoran et al. (2017) introduce a differentiable boundary tree (DBT) to integrate a DNN into the boundary tree (Mathy et al. 2015) to learn a better representation of data. The newly-learned data representation leads to a simpler boundary tree because the classes are well separated. The boundary tree is an online algorithm in which each node in the tree corresponds to a sample in the training set. The first sample together with its label are established as the tree root. Given a new query sample z , the sample traverses through the tree from the root to find the closest node n based on some distance function like the Euclidean distance function. If the label of the nearest node in the tree is different from the query sample, a new node containing the query z is added as a child of the closest node n in the tree; otherwise, the query node z is discarded. Therefore, each edge in the tree marks the boundary between two classes and each node tends to be close to these boundaries.

Transitions between nodes in a standard boundary tree are deterministic. DBT combines a SoftMax cost function with a boundary tree, resulting in stochastic transitions. Let x be a training sample and c be the one-hot encoding label of that sample. Given the current node x_j in the tree and a query node z , the transition probability from node x_j to node x_j , where $x_j \in \{child(x_j), x_j\}$ is the SoftMax of the negative distance between x_j and z . This is shown in Equation 25.

$$\begin{aligned}
 p(x_i \rightarrow x_j | z) &= \text{SoftMax}_{i, j \in \text{child}(i)}(-d(x_j, z)) \\
 &= \frac{\exp(-d(x_j, z))}{\sum_{j' \in \{i, j \in \text{child}(i)\}} \exp(-d(x_{j'}, z))}
 \end{aligned} \tag{25}$$

The probability of traversing a particular path in the boundary tree, given a query node z , is the product of the probability of each transition along the path from the root to the final node x_{final}^* in the tree. The final class log probability of DBT is computed by summing the probabilities of all transitions to the parent of x_{final}^* together with the probabilities of the final node and its siblings. The set $sibling(x_j)$ consists of all nodes sharing the same parent with node x_j and the node x_j itself. As discussed earlier, a DNN $f_\theta(x)$ transforms the inputs to learn a better representation. The final class log probabilities for the query node z are calculated as follows:

$$\begin{aligned}
 \log p(c | f_\theta(z)) &= \sum_{x_i \rightarrow x_j \in \text{path}^\dagger | f_\theta(z)} \log p(f_\theta(x_i) \rightarrow f_\theta(x_j) | f_\theta(z)) \\
 + \log &\sum_{x_k \in \text{sibling}(x_{final}^*)} p(\text{parent}(f_\theta(x_k)) \rightarrow f_\theta(x_k) | f_\theta(z))c(x_k).
 \end{aligned} \tag{26}$$

In Equation 26, path^\dagger denotes path^* (the path to the final node x_{final}^*) without the last transition, and $\text{sibling}(x)$ represents node x and all other nodes sharing the same parent with node x . The gradient descent algorithm can be applied to Equation 26 by plugging in a loss function to learn parameter θ of the DNN. However, gradient descent cannot be applied easily to DBT because of the node and edge manipulations in the graph. To address this issue, DBT transforms a small subset of training examples via a DNN and builds a boundary tree based on the transformed examples. Next, DBT transforms a query node z via the same DNN and calculates the log probability of a class according to Equation 26. The DNN employs gradient descent to update its parameters by propagating the gradient of log loss probability. DBT discards this boundary tree and iteratively builds a new boundary tree as described until a convergence criteria is met. In the described method, the authors set a specific threshold for the loss value to terminate the training. DBT is able to achieve greater accuracy with a simpler tree than the original boundary tree, as shown by the authors on the MNIST dataset (LeCun 1998). One of the biggest advantages of DBT is its interpretability. However, DBT is computationally an expensive method because a new computation graph needs to be built, which makes batching inefficient. Another limitation is that the algorithm needs to switch between building the tree and updating the tree. Therefore, scaling to large datasets is fairly prohibitive.

Often, k nearest neighbor (kNN) models are disregarded because of their computational cost and need for a large training set. In the traditional k-nearest neighbor algorithm (kNN), the posterior probability is estimated by the class distributions provided by points that are the closest neighbors to the point in question. In a special case of kNN, the 1-nearest neighbor (1NN) classifies the new point based on the nearest training (labeled prototype). To improve the kNN model, a variation uses a prototype learning model to generate prototypes that replace the original training set (Liu and Nakagawa 2001). In recent years, many neural

prototype learning (NPL) models have been developed. The NPL models can be categorized in two ways: (1) learned prototypes are points in the feature space that represent each class (Snell et al. 2017; Mettes et al. 2019), (2) learned prototypes are very close to the training set examples and a set of prototypes represents the training set (Li et al. 2018; Chen et al. 2019).

The first type of NPL learns a vector representing the mean of all of the points in a given class through an encoder network. The second type of NPL employs an encoder to learn a fixed-length feature vector z of size m . Next, a predefined number of prototypes n of size m utilizes z to learn meaningful prototypes. In general, the goal of these models is to minimize the sum of the misclassification loss plus two regularizers. The first regularizer pushes the prototype vectors to be meaningful by minimizing the average squared distance between the prototypes and the encoded vector. The second regularizer helps with clustering the training examples around prototypes by minimizing the average squared distance between the encoded vector and prototypes.

6.3 Neural Networks and SVMs

Yet another way of building a hybrid model is to learn a new representation of data with a DNN, then hand the resulting feature vectors off to other classifiers to learn a model. Tang (2013) explored replacing the last layer of DNNs with a linear SVM for classification tasks. The activation values of the penultimate layer are fed as input to an SVM with a L_2 regularizer. The weights of the lower layer are learned through momentum gradient descent by differentiating the SVM objective function with respect to activation of the penultimate layer. The author's experiments on the MNIST (LeCun 1998) and CIFAR-10 (Krizhevsky et al. 2010) datasets demonstrate that replacing a CNN's SoftMax output layer with SVM yields a lower test error. Tang (2013) postulate that the performance gain is due to the superior regularization effect of the SVM loss function.

It is worth mentioning that in their experiment on MNIST (LeCun 1998), Tang (2013) first used PCA to reduce the features and then fed the reduced feature vectors as input to their model. Also, Niu and Suen (2012) replaced the last layer of a CNN with an SVM, which similarly resulted in lowering test error of the model compared to a CNN on the MNIST dataset. Similar to these methods, Bellili et al. (2001), Azevedo and Zanchettin (2011), Nagi et al. (2012), and Zareapoor et al. (2018) replace the last layer of a DNN with an SVM. In these cases, their results from multiple datasets reveal that employing a SVM as the last layer of a neural network can improve the generalization of the network.

6.4 Neural Networks and Statistical Models

In some cases, two different data views are available. As an example, one view might contain video and another sound. Canonical correlation analysis (CCA) (Hotelling 1992) and kernel canonical correlation analysis (KCCA) (Hardoon et al. 2004) find basis vectors that maximize the correlations between the projections of the two views onto the basis vectors. Nonlinear representations learned by KCCA can achieve a higher correlation than linear representations learned by CCA. Despite the advantages of KCCA, the kernel function faces some drawbacks. Specifically, the representation is bound to the fixed kernel. Furthermore, the training time, as well as the time to compute the new data representation,

scales poorly with the size of the training set because of the non-parametric nature of kernel models.

Andrew et al. (2013) proposed to apply deep networks to learn a nonlinear data representation instead of employing a kernel function. Their resulting deep canonical correlation analysis (DCCA) consists of two separate deep networks for learning a new representation for each view. The new representation learned by the final layer of networks H_1 and H_2 is fed to CCA. To compute the objective gradient of DCCA, the gradient of the output of the correlation objective with respect to the new representation can be calculated as follows:

$$\frac{\partial_{corr}(H_1, H_2)}{\partial H_1} \quad (27)$$

After this computation, backpropagation is applied to find the gradient with respect to all parameters. The details of calculating the gradient in Equation 27 are provided by the authors (Andrew et al. 2013).

While researchers have also created LSTM methods that employ tree structures (Tai et al. 2015; Alvarez-Melis and Jaakkola 2017), these methods utilize the data structure to improve a network model rather than employing tree-based learning algorithms. Similarly, other approaches integrate non-neural network classifiers into a network structure. Cimino and Dell'Orletta (2016) and Agarap (2018) introduce hybrid models. These two methods apply LSTM and GRU, respectively, to learn a network representation. Unlike traditional DNNs, the last layer employs an SVM for classification. The work surveyed in this section provides evidence that deep neural nets are capable methods for learning high-level features. These features, in turn, can be used to improve the modeling capability for many types of supervised classifiers. In this survey, we aim to provide a thorough review of non-neural network models that utilize the unique features of deep network models. Table 1 provides a summary of such non-neural network models, organized based on four aspects of deep networks: model architecture, optimization, regularization, and hybrid model fusing. A known advantage of traditional deep networks compared with non-neural network models has been the ability to learn a better representation of input features. Inspired by various deep network architectures, deep learning of non-neural network classifiers has resulted in methods to also learn new feature representations. Another area where non-neural network classifiers have benefited from recent deep network research is applying backpropagation optimization to improve generalization. This table summarizes published efforts to apply regularization techniques that improve neural network generalization. The last category of models combines deep network classifiers and non-neural network classifiers to increase overall performance.

7 Experiments

In this paper, we survey a wide variety of models and methods. Our goal is to demonstrate that diverse types of models can benefit from deep learning techniques. To highlight this point, we empirically compare the performance of many techniques described in this

survey, as shown in Table 2. This comparison includes deep and shallow networks as well as non-neural network learning algorithms. Because of the variety of classifiers that are surveyed, we organize the comparison based on the learned model structure. We compare the performance of the models utilizing three datasets: (1) MNIST, (2) CIFAR-10, and (3) UCI Human Activity Recognition (HAR) (Anguita et al. 2013). MNIST instances contain 28×28 pixel grayscale images of handwritten digits and their labels. The MNIST labels are drawn from 10 object classes, with a total of 6000 training samples and 1000 testing samples. CIFAR-10 is also a well-known dataset containing 10 object classes with approximately 5000 examples per class, where each sample is a 32×32 pixel RGB image. The HAR data were collected from 30 participants performing six scripted activities (walking, walking upstairs, walking downstairs, sitting, standing, and laying) while wearing smartphones. The dataset contains 561 features extracted from sensors, including an accelerometer and a gyroscope. The training set contains 7352 samples from 70% of the volunteers and the testing set contains 2947 samples from the remaining 30% of volunteers.

We report the test error and model parameters provided by the authors for the mentioned datasets. If the performance of a model was not available for any of these datasets, we ran that experiment with the authors' code when available. We employ default values for parameters that are not specified in the original papers. In the event that the authors did not provide their code, we did not report any results. These omissions prevent the report of erroneous performances that result from implementation differences.

For a fair comparison, we divide Table 2 into different sections based on the type of models. First, we investigate the performance of the models that have a tree structure. The popularity of both neural networks and decision trees (DT) gives rise to a type of model that combines positive aspects of both models. We observe that models that integrated neural networks into their architecture, such as DNDF and ANT, outperform RF. Whereas stacking RF on DF displayed performance improvement for the MNIST and CIFAR-10 datasets, the multi-layer XGBoost model, mGBDT, did not perform well, even on a small dataset such as HAR. Additionally, we could not run mGBDT on MNIST data because of the computational cost, as mentioned in the original paper (Feng et al. 2018). In the case of models such as ANT, TAO, and SFDT, their performance does not necessarily exceed the other approaches, because these authors try to balance classification accuracy with model interpretability. Both DART and RRF make use of regularizers frequently used by neural networks. The results in Table 2 indicate that DART achieves consistently-strong performance. Although RRF did not perform well on MNIST, it did outperform other methods on HAR. We can conclude that there is no specific model that performs consistently well on all types of data.

Second, we observe that models DBT and NPL, which combine nearest neighbor strategies with neural networks, do yield strong classification performance on MNIST data while retaining the interpretability of kNN methods. Lastly, we study the models that utilize SVM as part of their structure. This type of model focuses on improving accuracy. We observe that a strategy like stacking SVMs (R2SVM) can improve performance over standard SVMs. Table 2 shows that the average error for R2SVM is 20.3 on the HAR dataset, while the average error for R2SVM is 21. The strategy of swapping the traditional Softmax deep network final layer with an SVM further improves accuracy to 11.9.

The results from our experiments reveal that both network classifiers and non-neural network classifiers benefit from deep learning. The methods surveyed in this paper and evaluated in these experiments demonstrate that non-neural network machine learning models do improve performance by incorporating DNN components into their algorithms. Whereas models without feature learning such as RF usually do not perform well on unstructured data such as images, we observe that adding deep learning to these models drastically improves their performance, as shown in Table 2. Additionally, non-deep models may achieve improved performance on structured data by adding regularizers, as shown in Table 2. The methods surveyed in this paper demonstrate that deep learning components can be added to any type of machine learning model and are not specific to DNNs. The incorporation of deep learning strategies is a promising direction for all types of classifiers, both network and non-neural network methods.

8 Conclusions and Directions for Ongoing Research

DNNs have emerged as a powerful force in the machine learning field for the past few years. This survey paper reviews the latest attempts to incorporate methods that are traditionally found in DNNs into other learning algorithms. DNNs work well when there is a large body of training data and available computational power. DNNs have consistently yielded strong results for a variety of datasets and competitions, such as winning the Large Scale Visual Recognition Challenge (Russakovsky et al. 2015) and achieving strong results for energy demand prediction (Paterakis et al. 2017), identifying gender of a text author (Sboev et al. 2018), stroke prediction (Hung et al. 2017), network intrusion detection (Yin et al. 2017), speech emotion recognition (Fayek et al. 2017), and taxi destination prediction (de Brébisson et al. 2015). Since there are many applications which lack large amounts of training data or for which the interpretability of a learned model is important, there is a need to integrate the benefits of DNNs with other classifier algorithms. Other classifiers have demonstrated improved performance on some types of data; therefore, the field can benefit from examining ways of combining deep learning elements between the network and non-neural network methods.

Although some work to date provides evidence that DNN techniques can be used effectively by other classifiers, there are still many challenges that researchers need to address, both to improve DNNs and to extend deep learning to other types of classifiers. Based on our survey of existing work, some related areas where supervised learners can benefit from unique DNN methods are outlined below.

The most characteristic feature of DNNs is a deep architecture and its ability to learn a new representation of data. A variety of stacked generalization methods have been developed to allow other machine learning methods to utilize deep architectures as well. These methods incorporate multiple classification steps in which the input of the next layer represents the concatenation of the output of the previous layer and the original feature vector, as discussed in Section 3.1.1. Future work can explore the many other possibilities that exist for refining the input features to each layer to better separate instances of each class at each layer.

Previous studies provide evidence that DNNs are effective data generators (Radford et al. 2016; Hoffman et al. 2018), while in some cases, non-neural network classifiers may actually be the better discriminators. Future research can consider using a DNN as a generator and an alternative classifier as a discriminator in generative adversarial models. Incorporating this type of model diversity could improve the robustness of the models.

Gradient descent can be applied to any differentiable algorithm. We observed that Kotschieder et al. (2015), Frosst and Hinton (2017), Tanno et al. (2019), and Zoran et al. (2017) all applied gradient descent to two different tree-based algorithms by making them differentiable. In the future, additional classifiers can be altered to be differentiable. Applying gradient descent to other algorithms could be an effective way to adjust the probability distribution of parameters.

Another area which is vital to investigate is the application of regularization methods that are customized for non-neural network classifiers. As discussed in Section 5, the non-neural network classifiers can benefit from the regularization methods that are unique to DNNs. However, there exist many different ways that these regularization methods can be adapted by non-neural network classifiers to improve model generalization.

An important area of research is interpretable models. There exist applications such as credit score, insurance risk, health status because of their sensitivity, models need to be interpretable. Further research needs to exploit the use of DNNs in interpretable models such as DNDT (Yang et al. 2018).

As we discussed in this survey, an emerging area of research is to combine the complementary benefits of statistical models with neural networks. Statistical models offer mathematical formalisms as well as possible explanatory power. This combination may provide a more effective model than either approach used in isolation.

There are cases in which the amount of ground truth-labeled data is limited, but a large body of labeled data from the same or similar distribution is available. One possible area of ongoing exploration is to couple the use of DNNs for learning from unlabeled data with the use of other classifier strategies for learning from labeled data. The simple model learned from labeled data can be exploited to further tune and improve learned representation patterns in the DNN.

We observe that currently, there is a general interest among the machine learning community to transfer new deep network developments to other classifiers. While a substantial effort has been made to incorporate deep learning ideas into the general machine learning field, continuing this work may spark the creation of new learning paradigms. However, the benefit between network-based learners and non-neural network learners can be bi-directional. Because a tremendous variety of classifiers has shown superior performance for a wide range of applications, future research can focus not only on how DNN techniques can improve non-neural network classifiers but on how DNNs can incorporate and benefit from non-neural network learning ideas as well.

Acknowledgements

The authors would like to thank Tharindu Adikari, Chris Choy, Ji Feng, Yani Ioannou, Stanislaw Jastrzebski and Marco A. Wiering for their valuable assistance in providing code and additional implementation details of the algorithms that were evaluated in this paper. We would also like to thank Samaneh Aminikhanghahi and Tinghui Wang for their feedback and guidance on the methods described in this survey. This material is based upon work supported by the National Science Foundation under Grant No. 1543656.

References

- Abdullah A, Veltkamp RC, Wiering MA (2009) An ensemble of deep support vector machines for image categorization. In: International Conference of Soft Computing and Pattern Recognition, SoCPaR, pp 301–306
- Agarap AF (2018) A neural network architecture combining gated recurrent unit (GRU) and support vector machine (SVM) for intrusion detection in network traffic data. In: Proceedings of the 10th International Conference on Machine Learning and Computing, ICMLC, ACM, pp 26–30
- Alaverdyan Z, Jung J, Bouet R, Lartizien C (2020) Regularized siamese neural network for unsupervised outlier detection on brain multiparametric magnetic resonance imaging: Application to epilepsy lesion screening. *Medical Image Analysis* 60:101618 [PubMed: 31841950]
- Alvarez-Melis D, Jaakkola TS (2017) Tree-structured decoding with doubly-recurrent neural networks. In: International Conference on Learning Representations, ICLR
- Andrew G, Arora R, Bilmes JA, Livescu K (2013) Deep canonical correlation analysis. In: Proceedings of the 30th International Conference on Machine Learning, ICML, vol 28, pp 1247–1255
- Anguita D, Ghio A, Oneto L, Parra X, Reyes-Ortiz JL (2013) A public domain dataset for human activity recognition using smartphones. In: 21st European Symposium on Artificial Neural Networks, ESANN
- Antoniou A, Storkey AJ, Edwards H (2017) Data augmentation generative adversarial networks. arXiv preprint: 1711.04340
- Azevedo WW, Zanchettin C (2011) A MLP-SVM hybrid model for cursive handwriting recognition. In: The International Joint Conference on Neural Networks, IJCNN, IEEE, pp 843–850
- Baumann P, Hochbaum DS, Yang YT (2019) A comparative study of the leading machine learning techniques and two new optimization algorithms. *European journal of operational research* 272(3):1041–1057
- Bellili A, Gilloux M, Gallinari P (2001) An hybrid MLP-SVM handwritten digit recognizer. In: International Conference on Document Analysis and Recognition, ICDAR, IEEE Computer Society, pp 28–33
- Bengio Y (2009) Learning deep architectures for AI. *Foundations and Trends in Machine Learning* 2(1):1–127
- Bengio Y (2013) Deep learning of representations: Looking forward. In: International Conference of Statistical Language and Speech Processing, SLSLP, vol 7978, pp 1–37
- Blundell C, Cornebise J, Kavukcuoglu K, Wierstra D (2015) Weight uncertainty in neural networks. arXiv preprint: 1505.05424
- Bowles C, Chen L, Guerrero R, Bentley P, Gunn RN, Hammers A, Dickie DA, del C Valdés Hernández M, Wardlaw JM, Rueckert D (2018) GAN augmentation: Augmenting training data using generative adversarial networks. arXiv preprint: 1810.10863
- de Brébisson A, Simon É, Auvolet A, Vincent P, Bengio Y (2015) Artificial neural networks applied to taxi destination prediction. In: Proceedings of the ECML/PKDD, vol 1526
- Breiman L (1996) Bagging predictors. *Machine learning* 24(2):123–140
- Breiman L (2000) Randomizing outputs to increase prediction accuracy. *Machine Learning* 40(3):229–242
- Breiman L, Friedman JH, Olshen RA, Stone CJ (1984) *Classification and Regression Trees*. Wadsworth

- Carreira-Perpiñán MÁ, Tavallali P (2018) Alternating optimization of decision trees, with application to learning sparse oblique trees. In: Advances in Neural Information Processing Systems, NeurIPS, pp 1219–1229
- Caruana R, Niculescu-Mizil A (2006) An empirical comparison of supervised learning algorithms. In: Proceedings of the 23rd international conference on Machine learning, ICML, ACM, vol 148, pp 161–168
- Caruana R, Karampatziakis N, Yessenalina A (2008) An empirical evaluation of supervised learning in high dimensions. In: Proceedings of the 23rd international conference on Machine learning, ICML, ACM, ACM International Conference Proceeding Series, vol 307, pp 96–103
- Caruana R, Lou Y, Gehrke J, Koch P, Sturm M, Elhadad N (2015) Intelligible models for healthcare: Predicting pneumonia risk and hospital 30-day readmission. In: Proceedings of the 21th ACM SIGKDD international conference on knowledge discovery and data mining, ACM, pp 1721–1730
- Cawley GC (2006) Leave-one-out cross-validation based model selection criteria for weighted ls-svms. In: Proceedings of the International Joint Conference on Neural Networks, IJCNN, pp 1661–1668
- Chen C, Li O, Tao D, Barnett A, Rudin C, Su J (2019) This looks like that: Deep learning for interpretable image recognition. In: Advances in Neural Information Processing Systems, NeurIPS, pp 8928–8939
- Chen W, Hays J (2018) Sketchygan: Towards diverse and realistic sketch to image synthesis. In: IEEE Conference on Computer Vision and Pattern Recognition, CVPR, pp 9416–9425
- Cho K, van Merriënboer B, Gülçehre Ç, Bahdanau D, Bougares F, Schwenk H, Bengio Y (2014) Learning phrase representations using RNN encoder-decoder for statistical machine translation. In: Proceedings of the Conference on Empirical Methods in Natural Language Processing, EMNLP, ACL, pp 1724–1734
- Cimino A, Dell’Orletta F (2016) Tandem LSTM-SVM approach for sentiment analysis. In: Proceedings of Third Italian Conference on Computational Linguistics, CLiC-it 2016, & Fifth Evaluation Campaign of Natural Language Processing and Speech Tools for Italian. Final Workshop, EVALITA, vol 1749
- Cubuk ED, Zoph B, Mané D, Vasudevan V, Le QV (2018) Autoaugment: Learning augmentation policies from data. arXiv preprint: 1805.09501
- Dai Z, Damianou AC, González J, Lawrence ND (2016) Variational auto-encoded deep gaussian processes. In: International Conference on Learning Representations, ICLR
- Damianou A (2015) Deep gaussian processes and variational propagation of uncertainty. PhD thesis, University of Sheffield
- Damianou AC, Lawrence ND (2013) Deep gaussian processes. In: Proceedings of the International Conference on Artificial Intelligence and Statistics, AISTATS, vol 31, pp 207–215
- Deng H, Runger GC (2012) Feature selection via regularized trees. In: The International Joint Conference on Neural Networks, IJCNN, IEEE, pp 1–8
- Deng H, Runger GC (2013) Gene selection with guided regularized random forest. Pattern Recognition 46(12):3483–3489
- Deng L, Yu D, Platt JC (2012) Scalable stacking and learning for building deep architectures. In: IEEE International Conference on Acoustics, Speech and Signal Processing, ICASSP, pp 2133–2136
- Duchi JC, Hazan E, Singer Y (2011) Adaptive subgradient methods for online learning and stochastic optimization. Journal of Machine Learning Research 12:2121–2159
- Dunlop MM, Girolami MA, Stuart AM, Teckenrump AL (2018) How deep are deep gaussian processes? Journal of Machine Learning Research 19:54:1–54:46
- Duvenaud D, Rippel O, Adams RP, Ghahramani Z (2014) Avoiding pathologies in very deep networks. In: Proceedings of the International Conference on Artificial Intelligence and Statistics, AISTATS, vol 33, pp 202–210
- Eickholt J, Cheng J (2013) Dndisorder: predicting protein disorder using boosting and deep networks. BMC Bioinform 14:88
- Esteva A, Robicquet A, Ramsundar B, Kuleshov V, DePristo M, Chou K, Cui C, Corrado G, Thrun S, Dean J (2019) A guide to deep learning in healthcare. Nature medicine 25(1):24–29
- Fan J, Li R (2001) Variable selection via nonconcave penalized likelihood and its oracle properties. Journal of the American statistical Association 96(456):1348–1360

- Fayek HM, Lech M, Cavedon L (2017) Evaluating deep learning architectures for speech emotion recognition. *Neural Networks* 92:60–68 [PubMed: 28396068]
- Feng J, Zhou Z (2018) Autoencoder by forest. In: *Proceedings of the AAAI Conference on Artificial Intelligence*, pp 2967–2973
- Feng J, Yu Y, Zhou Z (2018) Multi-layered gradient boosting decision trees. In: *Advances in Neural Information Processing Systems, NeurIPS*, pp 3555–3565
- Feng X, Jiang Y, Yang X, Du M, Li X (2019) Computer vision algorithms and hardware implementations: A survey. *Integration* 69:309–320
- Fortunato M, Blundell C, Vinyals O (2017) Bayesian recurrent neural networks. arXiv preprint: 1704.02798
- Freund Y (1995) Boosting a weak learning algorithm by majority. *Information and computation* 121(2):256–285
- Friedman JH (2001) Greedy function approximation: a gradient boosting machine. *Annals of statistics* pp 1189–1232
- Friedman JH (2002) Stochastic gradient boosting. *Computational Statistics & Data Analysis* 38(4):367–378
- Frosst N, Hinton GE (2017) Distilling a neural network into a soft decision tree. In: *Proceedings of the First International Workshop on Comprehensibility and Explanation in AI and ML*, vol 2071
- Fung G, Mangasarian OL (2004) A feature selection newton method for support vector machine classification. *Computational Optimization and Applications* 28(2):185–202
- Gjoreski M, Janko V, Slapnicar G, Mlakar M, Resçiç N, Bizjak J, Drobnic V, Marinko M, Mlakar N, Lustrek M, Gams M (2020) Classical and deep learning methods for recognizing human activities and modes of transportation with smartphone sensors. *Information Fusion* 62:47–62
- Goodfellow IJ, Pouget-Abadie J, Mirza M, Xu B, Warde-Farley D, Ozair S, Courville AC, Bengio Y (2014) Generative adversarial nets. In: *Advances in Neural Information Processing Systems, NeurIPS*, pp 2672–2680
- Goodfellow IJ, Bengio Y, Courville AC (2016) *Deep Learning*. Adaptive computation and machine learning, MIT Press
- Graves A, Mohamed A, Hinton GE (2013) Speech recognition with deep recurrent neural networks. In: *IEEE International Conference on Acoustics, Speech and Signal Processing, ICASSP, IEEE*, pp 6645–6649
- Graves A, Wayne G, Danihelka I (2014) Neural Turing machines. arXiv preprint: 1410.5401
- Guo C, Gao J, Wang YY, Deng L, He X (2017) Context-sensitive search using a deep learning model. US Patent 9,535,960
- Han J, Zhang D, Cheng G, Liu N, Xu D (2018) Advanced deep-learning techniques for salient and category-specific object detection: A survey. *IEEE Signal Processing Magazine* 35(1):84–100
- Hardoon DR, Szedmák S, Shawe-Taylor J (2004) Canonical correlation analysis: An overview with application to learning methods. *Neural Computation* 16(12):2639–2664 [PubMed: 15516276]
- Hatcher WG, Yu W (2018) A survey of deep learning: Platforms, applications and emerging research trends. *IEEE Access* 6:24411–24432
- Hinton GE, Salakhutdinov RR (2006) Reducing the dimensionality of data with neural networks. *Science* 313(5786):504–507 [PubMed: 16873662]
- Hinton GE, Osindero S, Teh YW (2006) A fast learning algorithm for deep belief nets. *Neural Computation* 18(7):1527–1554 [PubMed: 16764513]
- Ho TK (1995) Random decision forests. In: *International Conference on Document Analysis and Recognition, ICDAR*, pp 278–282
- Hochreiter S, Schmidhuber J (1997) Long short-term memory. *Neural Computation* 9(8):1735–1780 [PubMed: 9377276]
- Hoffman J, Tzeng E, Park T, Zhu J, Isola P, Saenko K, Efros AA, Darrell T (2018) Cycada: Cycle-consistent adversarial domain adaptation. In: *Proceedings of the 35th International Conference on Machine Learning, ICML*, vol 80, pp 1994–2003
- Hotelling H (1992) Relations between two sets of variates. In: *Breakthroughs in statistics*, Springer, pp 162–190

- Huang D, Huang W, Yuan Z, Lin Y, Zhang J, Zheng L (2018) Image superresolution algorithm based on an improved sparse autoencoder. *Information* 9(1):11
- Hung C, Chen W, Lai P, Lin C, Lee C (2017) Comparing deep neural network and other machine learning algorithms for stroke prediction in a large-scale population-based electronic medical claims database. In: 39th Annual International Conference of the IEEE Engineering in Medicine and Biology Society (EMBC), IEEE, pp 3110–3113
- Hutchinson B, Deng L, Yu D (2013) Tensor deep stacking networks. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 35(8):1944–1957 [PubMed: 23267198]
- Ioannou Y, Robertson DP, Zikic D, Kotschieder P, Shotton J, Brown M, Criminisi A (2016) Decision forests, convolutional networks and the models in-between. *arXiv preprint: 1603.01250*
- Jaitly N, Hinton GE (2013) Vocal tract length perturbation (VTLP) improves speech recognition. In: *Proceedings ICML Workshop on Deep Learning for Audio, Speech and Language*, vol 117
- Ji S, Xu W, Yang M, Yu K (2013) 3D convolutional neural networks for human action recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 35(1):221–231 [PubMed: 22392705]
- Jorge J, Vieco J, Paredes R, Sánchez J, Benedí J (2018) Empirical evaluation of variational autoencoders for data augmentation. In: *Proceedings of the International Joint Conference on Computer Vision, Imaging and Computer Graphics Theory and Applications, VISIGRAPP*, pp 96–104
- Józefowicz R, Vinyals O, Schuster M, Shazeer N, Wu Y (2016) Exploring the limits of language modeling. *arXiv preprint: 1602.02410*
- Kadurin A, Aliper A, Kazennov A, Mamoshina P, Vanhaelen Q, Khrabrov K, Zhavoronkov A (2017) The cornucopia of meaningful leads: applying deep adversarial autoencoders for new molecule development in oncology. *Oncotarget* 8(7):10883 [PubMed: 28029644]
- Kalman RE (1960) A new approach to linear filtering and prediction problems. *Journal of basic Engineering* 82(1):35–45
- King RD, Feng C, Sutherland A (1995) Statlog: comparison of classification algorithms on large real-world problems. *Applied Artificial Intelligence an International Journal* 9(3):289–333
- Kingma DP, Ba J (2015) Adam: A method for stochastic optimization. In: *International Conference on Learning Representations, ICLR*
- Kingma DP, Welling M (2014) Auto-encoding variational bayes. In: *International Conference on Learning Representations, ICLR*
- Kotschieder P, Fiterau M, Criminisi A, Bulò SR (2015) Deep neural decision forests. In: *IEEE International Conference on Computer Vision, ICCV*, pp 1467–1475
- Krishnan RG, Shalit U, Sontag DA (2015) Deep kalman filters. *arXiv preprint: 1511.05121*
- Krizhevsky A, Hinton G (2009) Learning multiple layers of features from tiny images. *Tech. rep, Citeseer*
- Krizhevsky A, Nair V, Hinton G (2010) Cifar-10 (canadian institute for advanced research) URL <http://www.cs.toronto.edu/~kriz/cifar.html>
- Krizhevsky A, Sutskever I, Hinton GE (2012) Imagenet classification with deep convolutional neural networks. In: *Advances in Neural Information Processing Systems, NeurIPS*, pp 1097–1105
- Krueger D, Ballas N, Jastrzebski S, Arpit D, Kanwal MS, Maharaj T, Bengio E, Fischer A, Courville AC (2017) Deep nets don't learn via memorization. In: *International Conference on Learning Representations, ICLR*
- LeCun Y (1998) The mnist database of handwritten digits URL <http://yann.lecun.com/exdb/mnist/>
- LeCun Y, Bengio Y, Hinton G (2015) Deep learning. *Nature* 521(7553):436 [PubMed: 26017442]
- LeCun Y, et al. (1989) Generalization and network design strategies. *Connectionism in perspective* pp 143–155
- Li O, Liu H, Chen C, Rudin C (2018) Deep learning for case-based reasoning through prototypes: A neural network that explains its predictions. In: *Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence*, pp 3530–3537
- Lim T, Loh W, Shih Y (2000) A comparison of prediction accuracy, complexity, and training time of thirty-three old and new classification algorithms. *Machine Learning* 40(3):203–228

- Liong VE, Lu J, Wang G (2013) Face recognition using deep PCA. In: International Conference on Information, Communications & Signal Processing, ICICS, pp 1–5
- Liu C, Nakagawa M (2001) Evaluation of prototype learning algorithms for nearest-neighbor classifier in application to handwritten character recognition. *Pattern Recognition* 34(3):601–615
- Liu X, Zou Y, Kong L, Diao Z, Yan J, Wang J, Li S, Jia P, You J (2018) Data augmentation via latent space interpolation for image classification. In: International Conference on Pattern Recognition, ICPR, pp 728–733
- Lu G, Ouyang W, Xu D, Zhang X, Gao Z, Sun M (2018) Deep kalman filtering network for video compression artifact reduction. In: Proceedings of the European Conference Computer Vision, ECCV, vol 11218, pp 591–608
- Lucchese C, Nardini FM, Orlando S, Perego R, Trani S (2017) X-DART: blending dropout and pruning for efficient learning to rank. In: Proceedings of the 40th International ACM SIGIR Conference on Research and Development in Information Retrieval, ACM, pp 1077–1080
- Mathy C, Derbinsky N, Bento J, Rosenthal J, Yedidia JS (2015) The boundary forest algorithm for online supervised and unsupervised learning. In: Proceedings of the Twenty-Ninth AAAI Conference on Artificial Intelligence, pp 2864–2870
- Mattos CLC, Dai Z, Damianou AC, Forth J, Barreto GA, Lawrence ND (2016) Recurrent gaussian processes. In: International Conference on Learning Representations, ICLR
- Memisevic R, Hinton GE (2007) Unsupervised learning of image transformations. In: IEEE Computer Society Conference on Computer Vision and Pattern Recognition, CVPR
- Mettes P, van der Pol E, Snoek C (2019) Hyperspherical prototype networks. In: Advances in Neural Information Processing Systems, NeurIPS, pp 1485–1495
- Minsky M, Papert S (1987) *Perceptrons - an introduction to computational geometry*. MIT Press
- Mishkin D, Matas J (2016) All you need is a good init. In: International Conference on Learning Representations, ICLR
- Moghimi M, Belongie SJ, Saberian MJ, Yang J, Vasconcelos N, Li L (2016) Boosted convolutional neural networks. In: Proceedings of the British Machine Vision Conference, BMVC
- Nagi J, Caro GAD, Giusti A, Nagi F, Gambardella LM (2012) Convolutional neural support vector machines: Hybrid visual pattern classifiers for multi-robot systems. In: International Conference on Machine Learning and Applications, ICMLA, IEEE, pp 27–32
- Ndikumana A, Hong CS (2019) Self-driving car meets multi-access edge computing for deep learning-based caching. In: International Conference on Information Networking, ICOIN, IEEE, pp 49–54
- Nguyen AM, Yosinski J, Clune J (2015) Deep neural networks are easily fooled: High confidence predictions for unrecognizable images. In: IEEE Conference on Computer Vision and Pattern Recognition, CVPR, IEEE Computer Society, pp 427–436
- Niu X, Suen CY (2012) A novel hybrid CNN-SVM classifier for recognizing handwritten digits. *Pattern Recognition* 45(4):1318–1325
- Nøklund A (2016) Direct feedback alignment provides learning in deep neural networks. In: Advances in Neural Information Processing Systems, NeurIPS, pp 1037–1045
- Ohashi H, Al-Nasser M, Ahmed S, Akiyama T, Sato T, Nguyen P, Nakamura K, Dengel A (2017) Augmenting wearable sensor data with physical constraint for dnn-based human-action recognition. In: ICML 2017 Times Series Workshop, pp 6–11
- Park T, Liu M, Wang T, Zhu J (2019) Semantic image synthesis with spatially-adaptive normalization. In: IEEE Conference on Computer Vision and Pattern Recognition, CVPR, pp 2337–2346
- Paterakis NG, Mocanu E, Gibescu M, Stappers B, van Alst W (2017) Deep learning versus traditional machine learning methods for aggregated energy demand prediction. In: IEEE PES Innovative Smart Grid Technologies Conference Europe, ISGT, IEEE, pp 1–6
- Perez L, Wang J (2017) The effectiveness of data augmentation in image classification using deep learning. arXiv preprint: 1712.04621
- Poole B, Sohl-Dickstein J, Ganguli S (2014) Analyzing noise in autoencoders and deep networks. arXiv preprint: 1406.1831

- Pouyanfar S, Sadiq S, Yan Y, Tian H, Tao Y, Reyes MEP, Shyu M, Chen S, Iyengar SS (2019) A survey on deep learning: Algorithms, techniques, and applications. *ACM Computing Surveys* 51(5):92:1–92:36
- Radford A, Metz L, Chintala S (2016) Unsupervised representation learning with deep convolutional generative adversarial networks. In: *International Conference on Learning Representations, ICLR*
- Rashmi KV, Gilad-Bachrach R (2015) DART: dropouts meet multiple additive regression trees. In: *Proceedings of the Eighteenth International Conference on Artificial Intelligence and Statistics, AISTATS*, vol 38
- Rasmussen CE (2003) Gaussian processes in machine learning. In: *Summer School on Machine Learning*, Springer, vol 3176, pp 63–71
- Reed SE, Akata Z, Yan X, Logeswaran L, Schiele B, Lee H (2016) Generative adversarial text to image synthesis. In: *Proceedings of the International Conference on Machine Learning, ICML, JMLR Workshop and Conference Proceedings*, vol 48, pp 1060–1069
- Riad R, Dancette C, Karadayi J, Zeghidour N, Schatz T, Dupoux E (2018) Sampling strategies in siamese networks for unsupervised speech representation learning. In: *Conference of the International Speech Communication Association, ISCA*, pp 2658–2662
- Rios LM, Sahinidis NV (2013) Derivative-free optimization: a review of algorithms and comparison of software implementations. *Journal of Global Optimization* 56(3):1247–1293
- Rosenblatt F (1958) The perceptron: a probabilistic model for information storage and organization in the brain. *Psychological review* 65(6):386 [PubMed: 13602029]
- Rumelhart DE, Hinton GE, Williams RJ (1985) Learning internal representations by error propagation. Tech. rep. California Univ San Diego La Jolla Inst for Cognitive Science
- Russakovsky O, Deng J, Su H, Krause J, Satheesh S, Ma S, Huang Z, Karpathy A, Khosla A, Bernstein M, Berg AC, Fei-Fei L (2015) ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision, IJCV* 115(3):211–252
- Sabour S, Frosst N, Hinton GE (2017) Dynamic routing between capsules. In: *Advances in Neural Information Processing Systems, NeurIPS*, pp 3856–3866
- Salzberg S (1994) Book review: C4.5: programs for machine learning by j. ross quinlan. morgan kaufmann publishers, inc., 1993. *Machine Learning* 16(3):235–240
- Samuel AL (1959) Some studies in machine learning using the game of checkers. *IBM Journal of research and development* 3(3):210–229
- Sboev A, Moloshnikov I, Gudovskikh D, Selivanov A, Rybka R, Litvinova T (2018) Deep learning neural nets versus traditional machine learning in gender identification of authors of rusprofiling texts. *Procedia Computer Science* 123:424–431
- Shashua SD, Mannor S (2017) Deep robust kalman filter. arXiv preprint: 1703.02310
- Shickel B, Tighe P, Bihorac A, Rashidi P (2018) Deep EHR: A survey of recent advances in deep learning techniques for electronic health record (EHR) analysis. *IEEE Journal of Biomedical And Health Informatics* 22(5):1589–1604 [PubMed: 29989977]
- Shlens J (2014) A tutorial on principal component analysis. arXiv preprint: 1404.1100
- Singh M, Bajpai U, Prasath S, et al. (2020) Generation of fashionable clothes using generative adversarial networks: A preliminary feasibility study. *International Journal of Clothing Science and Technology* 32(2):177–187
- Smagulova K, James AP (2019) A survey on lstm memristive neural network architectures and applications. *The European Physical Journal Special Topics* 228(10):2313–2324
- Snell J, Swersky K, Zemel RS (2017) Prototypical networks for few-shot learning. In: *Advances in Neural Information Processing Systems, NeurIPS*, pp 4077–4087
- Srivastava N, Hinton GE, Krizhevsky A, Sutskever I, Salakhutdinov R (2014) Dropout: a simple way to prevent neural networks from overfitting. *Journal of Machine Learning Research* 15(1):1929–1958
- Sutskever I, Hinton GE (2007) Learning multilevel distributed representations for high-dimensional sequences. In: *Proceedings of the International Conference on Artificial Intelligence and Statistics, AISTATS*, vol 2, pp 548–555

- Sutskever I, Martens J, Dahl GE, Hinton GE (2013) On the importance of initialization and momentum in deep learning. In: Proceedings of the International Conference on Machine Learning, ICML, vol 28, pp 1139–1147
- Swersky K, Tarlow D, Sutskever I, Salakhutdinov R, Zemel RS, Adams RP (2012) Cardinality restricted boltzmann machines. In: Advances in Neural Information Processing Systems, NeurIPS, pp 3302–3310
- Szegedy C, Liu W, Jia Y, Sermanet P, Reed SE, Anguelov D, Erhan D, Vanhoucke V, Rabinovich A (2015) Going deeper with convolutions. In: IEEE Conference on Computer Vision and Pattern Recognition, CVPR, IEEE Computer Society, pp 1–9
- Tai KS, Socher R, Manning CD (2015) Improved semantic representations from tree-structured long short-term memory networks. In: Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing of the Asian Federation of Natural Language Processing, ACL, The Association for Computer Linguistics, pp 1556–1566
- Taigman Y, Yang M, Ranzato M, Wolf L (2014) Deepface: Closing the gap to human-level performance in face verification. In: IEEE Conference on Computer Vision and Pattern Recognition, CVPR, IEEE Computer Society, pp 1701–1708
- Tang C, Wang Y, Xu J, Sun Y, Zhang B (2018) Efficient scenario generation of multiple renewable power plants considering spatial and temporal correlations. *Applied Energy* 221:348–357
- Tang Y (2013) Deep learning using support vector machines. arXiv preprint: 1306.0239
- Tang Y, Eliasmith C (2010) Deep networks for robust visual recognition. In: Proceedings of the International Conference on Machine Learning, ICML, pp 1055–1062
- Tanno R, Arulkumaran K, Alexander DC, Criminisi A, Nori AV (2019) Adaptive neural trees. In: Proceedings of the International Conference on Machine Learning, ICML, vol 97, pp 6166–6175
- Tibshirani R (1996) Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society Series B Methodological* pp 267–288
- Ting KM, Witten IH (1999) Issues in stacked generalization. *Journal of Artificial Intelligence Research* 10:271–289
- Utkin LV, Ryabinin MA (2018) A siamese deep forest. *Knowledge-Based Systems* 139:13–22
- Vincent P, Larochelle H, Bengio Y, Manzagol P (2008) Extracting and composing robust features with denoising autoencoders. In: Proceedings of the International Conference of Machine Learning, ICML, vol 307, pp 1096–1103
- Vinyals O, Jia Y, Deng L, Darrell T (2012) Learning with recursive perceptual representations. In: Advances in Neural Information Processing Systems, NeurIPS, pp 2834–2842
- Wang G, Zhang G, Choi K, Lu J (2019a) Deep additive least squares support vector machines for classification with model transfer. *IEEE Transactions on Systems, Man, and Cybernetics: Systems* 49(7):1527–1540
- Wang J, Chen Y, Hao S, Peng X, Hu L (2019b) Deep learning for sensor-based activity recognition: A survey. *Pattern Recognition Letters* 119:3–11
- Wang SI, Manning CD (2012) Baselines and bigrams: Simple, good sentiment and topic classification. In: Proceedings of the 50th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers), pp 90–94
- Wang SI, Manning CD (2013) Fast dropout training. In: Proceedings of the International Conference on Machine Learning, ICML, vol 28, pp 118–126
- Widrow B, Hoff ME (1960) Adaptive switching circuits. Tech. rep, Stanford Univ CA Stanford Electronics Labs
- Wiering MA, Schomaker LR (2014) Multi-layer support vector machines. Regularization, optimization, kernels, and support vector machines p 457
- van der Wilk M, Rasmussen CE, Hensman J (2017) Convolutional gaussian processes. In: Advances in Neural Information Processing Systems, NeurIPS, pp 2849–2858
- Wolpert DH (1992) Stacked generalization. *Neural Networks* 5(2):241–259
- Wolpert DH, Macready WG (1997) No free lunch theorems for optimization. *IEEE Transactions on Evolutionary Computation* 1(1):67–82

- Wong SC, Gatt A, Stamatescu V, McDonnell MD (2016) Understanding data augmentation for classification: When to warp? In: International Conference on Digital Image Computing: Techniques and Applications, DICTA, IEEE, pp 1–6
- Wu Y, Schuster M, Chen Z, Le QV, Norouzi M, Macherey W, Krikun M, Cao Y, Gao Q, Macherey K, Klingner J, Shah A, Johnson M, Liu X, Kaiser L, Gouws S, Kato Y, Kudo T, Kazawa H, Stevens K, Kurian G, Patil N, Wang W, Young C, Smith J, Riesa J, Rudnick A, Vinyals O, Corrado G, Hughes M, Dean J (2016) Google’s neural machine translation system: Bridging the gap between human and machine translation. arXiv preprint: 1609.08144
- Xu R (2013) Improvements to random forest methodology. PhD thesis, Iowa State University
- Yang F, Poostchi M, Yu H, Zhou Z, Silamut K, Yu J, Maude RJ, Jäger S, Antani SK (2020) Deep learning for smartphone-based malaria parasite detection in thick blood smears. *IEEE Journal of Biomedical and Health Informatics* 24(5):1427–1438 [PubMed: 31545747]
- Yang H, Wu J (2012) Practical large scale classification with additive kernels. In: Proceedings of the Asian Conference on Machine Learning, ACML, vol 25, pp 523–538
- Yang Y, Morillo IG, Hospedales TM (2018) Deep neural decision trees. arXiv preprint: 1806.06988
- Yeh C, Wu W, Ko W, Wang YF (2017) Learning deep latent space for multilabel classification. In: Proceedings of AAAI Conference on Artificial Intelligence, pp 2838–2844
- Yin C, Zhu Y, Fei J, He X (2017) A deep learning approach for intrusion detection using recurrent neural networks. *IEEE Access* 5:21954–21961
- Yiu D, Deng L (2011) Deep convex net: A scalable architecture for speech pattern classification. In: Annual Conference of the International Speech Communication Association, INTERSPEECH, pp 2285–2288
- Zareapoor M, Shamsolmoali P, Jain DK, Wang H, Yang J (2018) Kernelized support vector machine with deep learning: An efficient approach for extreme multiclass dataset. *Pattern Recognition Letters* 115:4–13
- Zhang C, Bengio S, Hardt M, Recht B, Vinyals O (2017) Understanding deep learning requires rethinking generalization. In: International Conference on Learning Representations, ICLR
- Zhang HH, Ahn J, Lin X, Park C (2006) Gene selection using support vector machines with non-convex penalty. *Bioinformatics* 22(1):88–95 [PubMed: 16249260]
- Zhang R, Isola P, Efros AA (2016) Colorful image colorization. In: European Conference on Computer Vision, ECCV, vol 9907, pp 649–666
- Zhao T, Zhang B, He M, Zhang W, Zhou N, Yu J, Fan J (2018) Embedding visual hierarchy with deep networks for large-scale visual recognition. *IEEE Transactions on Image Processing* 27(10):4740–4755
- Zhou Y, Chellappa R (1988) Computation of optical flow using a neural network. In: Proceedings of International Conference on Neural Networks, ICNN’88, IEEE, pp 71–78
- Zhou Z, Feng J (2017) Deep forest: Towards an alternative to deep neural networks. In: Proceedings of the International Joint Conference on Artificial Intelligence (IJCAI), pp 3553–3559
- Zoran D, Lakshminarayanan B, Blundell C (2017) Learning deep nearest neighbor representations using differentiable boundary trees. arXiv preprint: 1702.08833
- Zuo Y, Avraham G, Drummond T (2018) Generative adversarial forests for better conditioned adversarial learning. arXiv preprint: 1805.05185

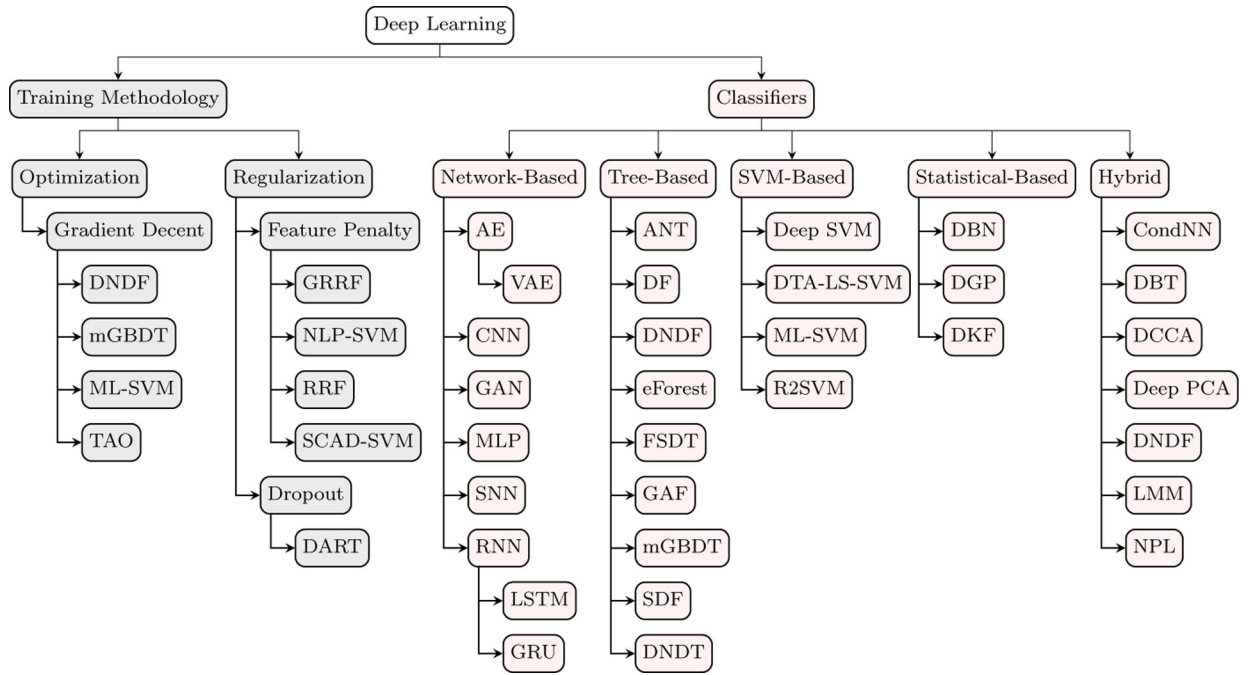


Fig. 1. Content map of the methods covered in this survey.

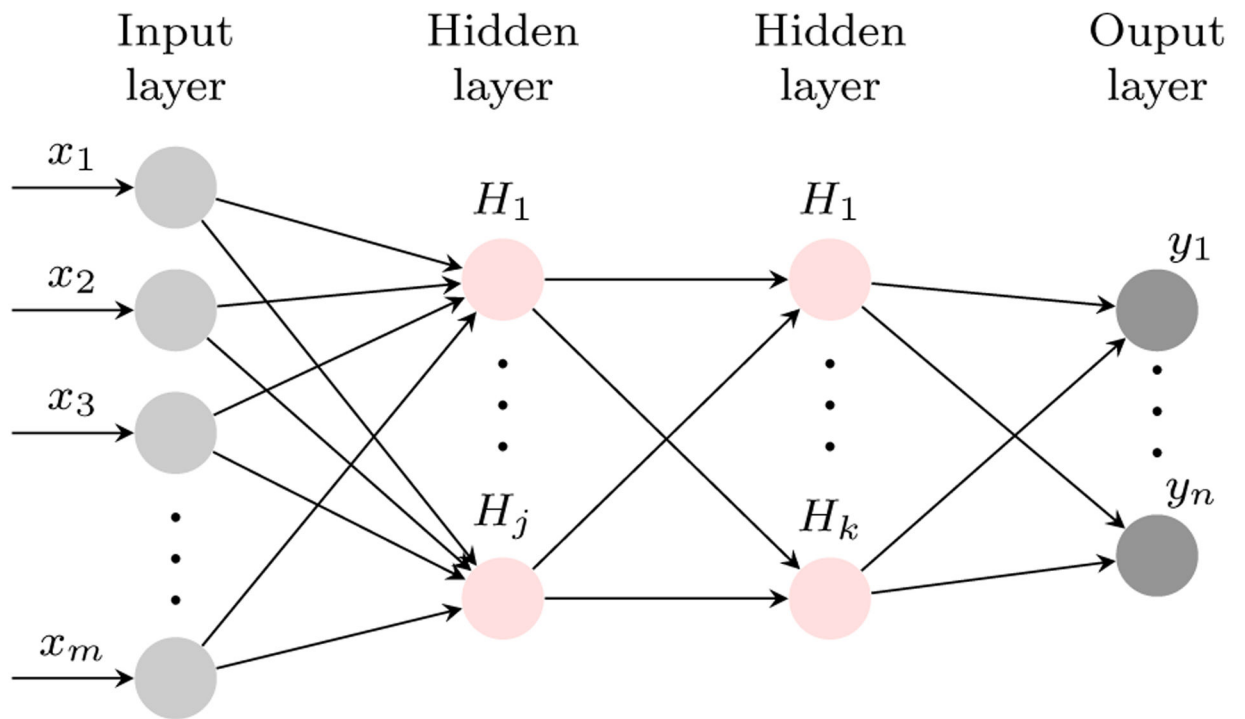


Fig. 2.
An illustration of a three-layered MLP with j nodes at the first hidden layer and k at the second layer.

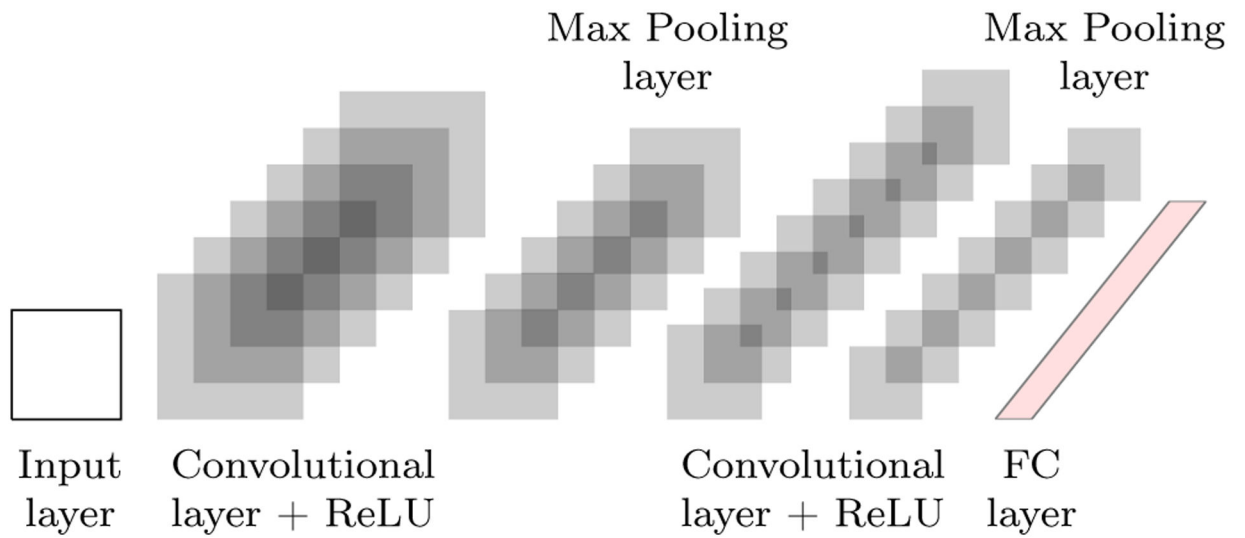


Fig. 3.

An illustration of a three-layered CNN made of six convolution filters followed by six max pooling filters at the first layer, and eight convolution filters followed by seven max pooling filters at the second layer. The last layer is a fully-connected layer (FC).

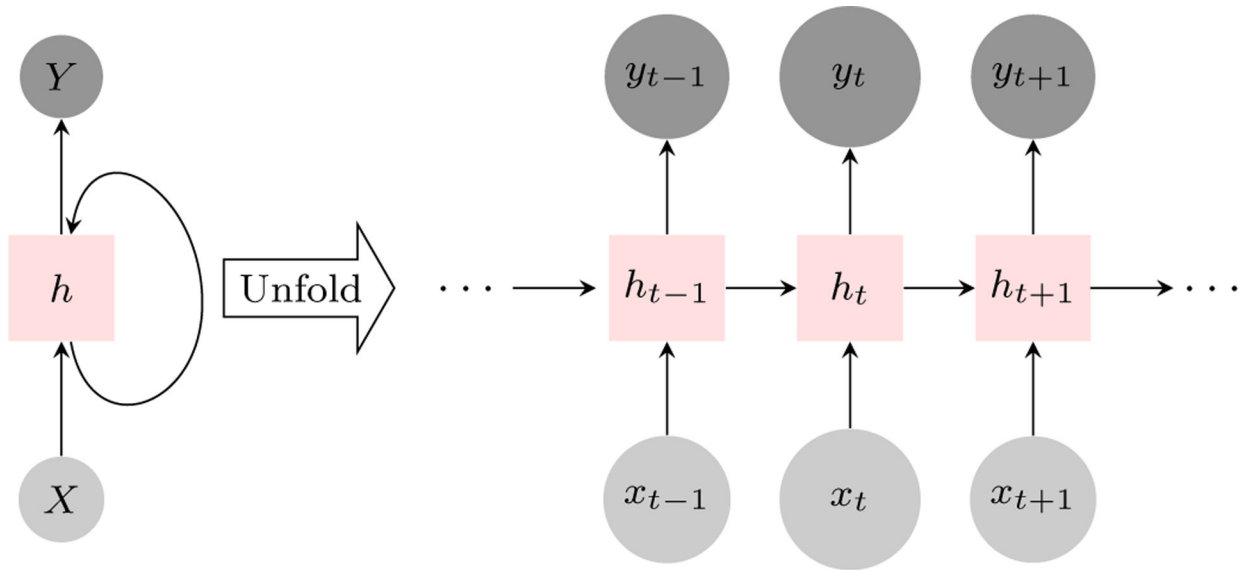


Fig. 4.
An illustration of a simple RNN and its unfolded structure through time t .

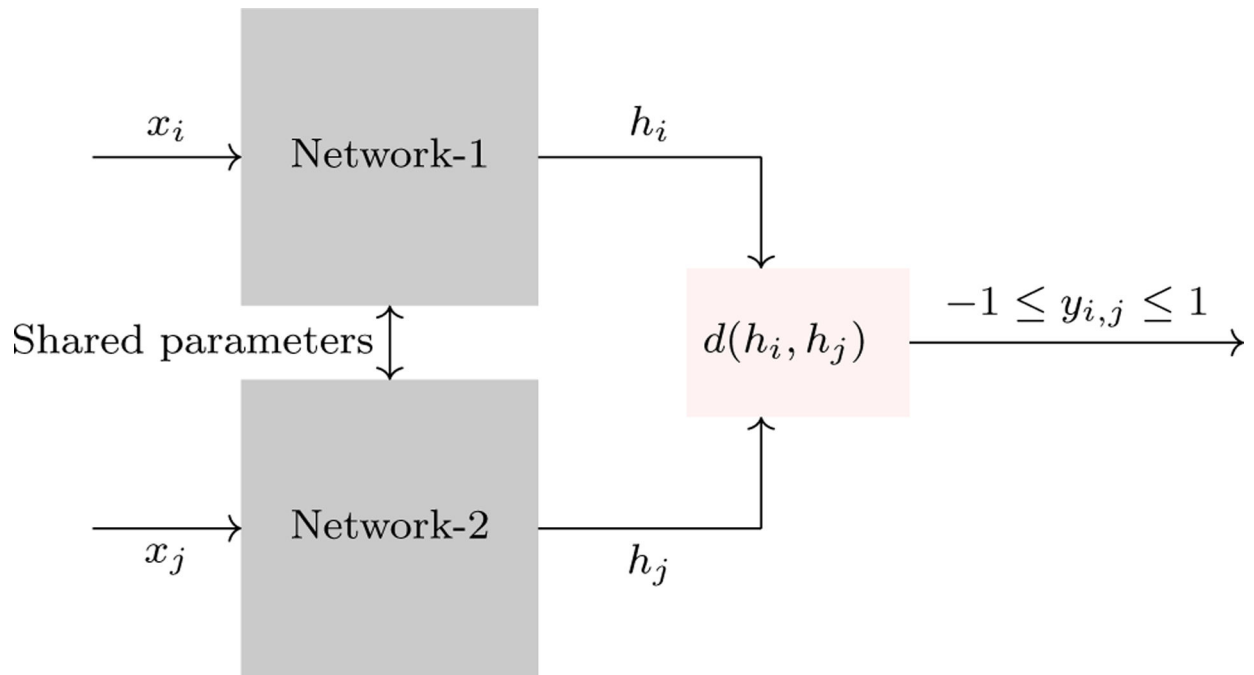


Fig. 5.

An illustration of an SNN. In this figure, x_i and x_j are two data vectors corresponding to a pair of instances from the training set. Both networks share the same weights and map the input to a new representation. By comparing the outputs of the networks using a distance measure such as Euclidean, we can determine the compatibility between instances x_i and x_j .

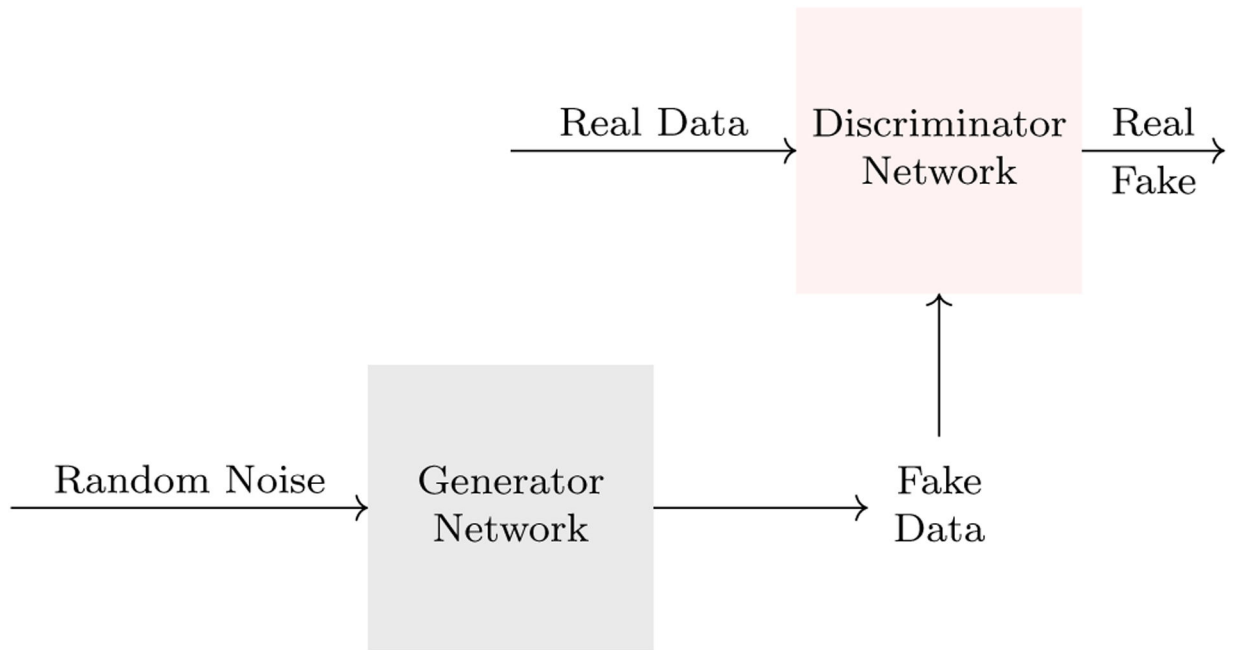


Fig. 6. An illustration of a GAN. The goal of the discriminator network is to distinguish real data from fake data, and the goal of the generator network is to use the feedback from the discriminator to generate data that the discriminator cannot distinguish from real.

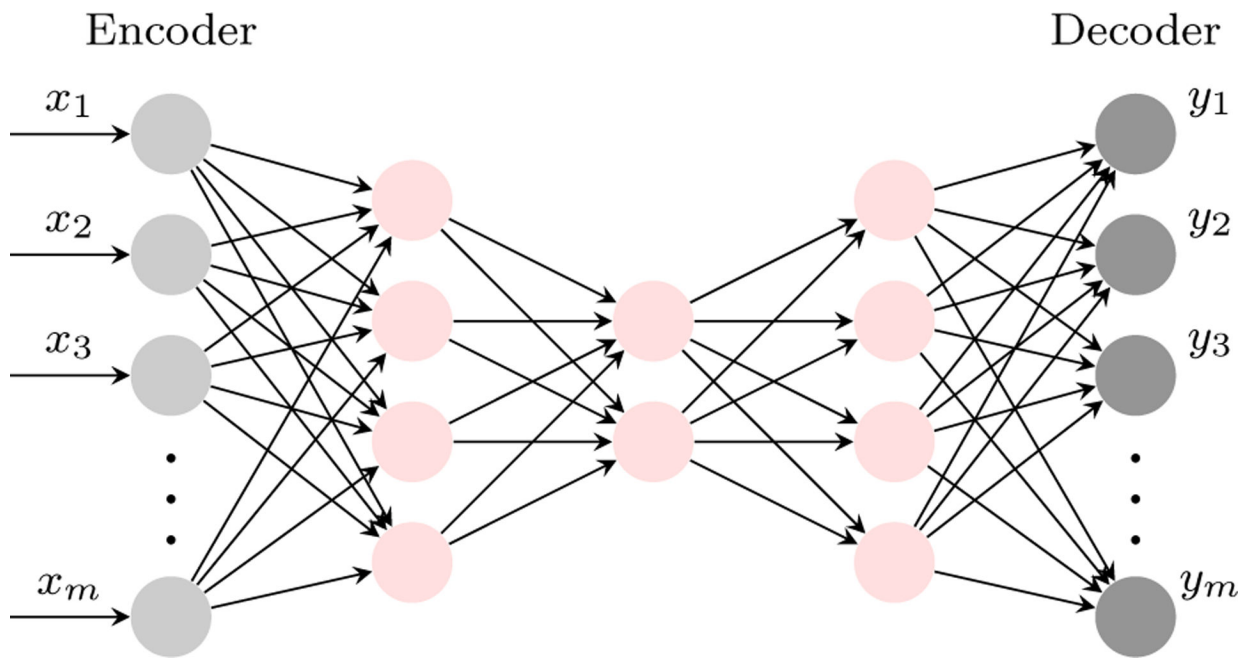


Fig. 7.

An illustration of an AE. The first part of the network, called the encoder, compresses input into a latent-space by learning the function $h = f(x)$. The second part, called the decoder, reconstructs the input from the latent-space representation by learning the function $\hat{y} = g(h)$.

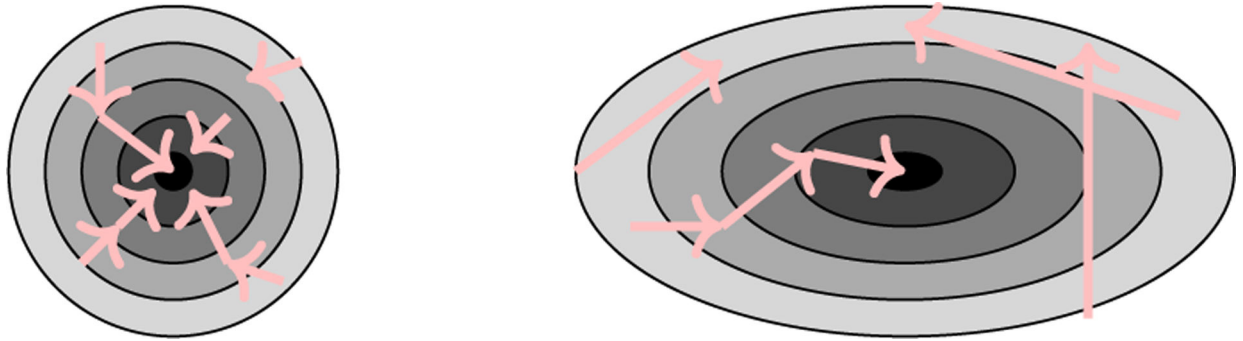


Fig. 8.

The loss surface on the left depicts a well-conditioned model where local minima can be reached from all directions. The loss surface on the right depicts an ill-conditioned model where there are several ways to overshoot or never reach the minima.

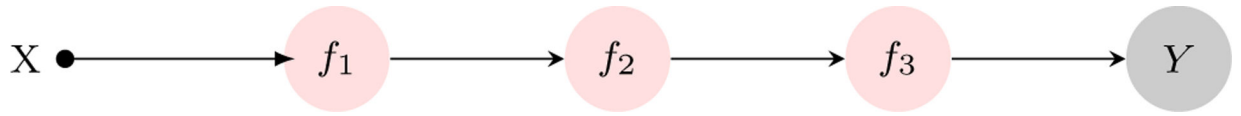


Fig. 9.
A deep Gaussian process with two hidden layers.

Table 1

Summary of classifiers which integrate deep network components into non-neural network classifiers.

	Methods	Classifiers
Architecture	Feedforward	ANT (Tanno et al. 2019), DNDT (Yang et al. 2018), DBN (Hinton et al. 2006), Deep PCA (Liong et al. 2013), DF (Zhou and Feng 2017), DPG (Damianou 2015), R2-SVM (Vinyals et al. 2012), D-SVM (Abdullah et al. 2009), DTA-LS-SVM (Wang et al. 2019a), SFDT (Frosst and Hinton 2017)
	Autoencoder	DKF (Krishnan et al. 2015), eForest (Feng and Zhou 2018), ML-SVM (Wiering and Schomaker 2014)
	Siamese Model	SDF (Utkin and Ryabinin 2018)
	Generative Adversarial Model	GAF (Zuo et al. 2018)
Optimization	Gradient Decent	DNDF (Kontschieder et al. 2015), mGBDT (Feng et al. 2018), ML-SVM (Wiering and Schomaker 2014), TAO (Carreira-Perpiñán and Tavallali 2018)
Regularization	Parameter Norm Penalty	NLP-SVM (Fung and Mangasarian 2004), GRRF (Deng and Runger 2013), RRF (Deng and Runger 2012), SCAD-SVM (Zhang et al. 2006)
	Dropout	DART (Rashmi and Gilad-Bachrach 2015)
Hybrid Model		CondCNN Ioannou et al. (2016), DBT (Zoran et al. 2017), DCCA (Andrew et al. 2013), DNDF (Kontschieder et al. 2015), LMM (Zhao et al. 2018), <i>DNN+SVM</i> : (Tang 2013) (Niu and Suen 2012) (Zareapoor et al. 2018) (Nagi et al. 2012) (Bellili et al. 2001) (Azevedo and Zanchettin 2011) <i>NPL</i> : (Snell et al. 2017)(Mettes et al. 2019)(Li et al. 2018) (Chen et al. 2019)

Table 2

Comparison of model performance on MNIST, CIFAR-10, and HAR datasets. Depth indicates the depth of a decision tree. Ens. is the ensemble size, which in the case of multi-layer models is the number of trees at each layer. The number of parameters refers to the number of weights. Time reflects the elapsed training time in seconds. (*) reported result reflects the code was ran on a machine with 31GB memory, Intel Core i7–8700K CPU, and Nvidia RTX2070 GPU.

	Model	Structure	Depth	Ens.	Parms.	Time	Avg. Err	
MNIST	RF (Breiman 2000)*	Ensemble of DTs	NA	200	NA	53.78	2.95	
	ANT (Tanno et al. 2019)	Soft DT with integrated neural network	NA	1	100,596	NA	0.64	
	ANT (Tanno et al. 2019)	Soft DT with integrated neural network	NA	8	850,775	NA	0.29	
	DF (Zhou and Feng 2017)	Stacked Forest	NA	1000	NA	NA	0.74	
	Tree	SFDT (Frosst and Hinton 2017)	Soft DT based on neural network	NA	1	NA	NA	5.55
	DNDF (Kontschieder et al. 2015)	Neural network with soft DTs as output layer	5	10	60,000	NA	0.7	
	TAO (Carreira-Perpinan and Tavallali 2018)	Sparse oblique tree	12	1	10,000	NA	5.69	
	RRF (Deng and Runger 2012)*	RF with regularizer	NA	200	NA	2057.97	4.87	
	DART (Rashmi and Gilad-Bachrach 2015)*	XGBoost with regularizer	6	NA	NA	475.25	2.91	
	kNN	DBT (Zoran et al. 2017)	Boundry tree with integrated DNN	NA	NA	482,630	NA	1.85
NPL (Li et al. 2018)	Neural Prototype Learning	NA	NA	NA	NA	NA	0.47	
CIFAR-10	RF (Breiman 2000)	Ensemble of DT	NA	2000	NA	NA	50.17	
	ANT (Tanno et al. 2019)	Soft DT with integrated neural network	NA	1	1.4M	NA	8.31	
	Tree	ANT (Tanno et al. 2019)	Soft DT with integrated neural network	NA	8	8.7M	NA	7.71
	DF (Zhou and Feng 2017)	Stacked Forest	NA	1000	NA	DF	38.22	
	CondCNN (Ioannou et al. 2016)	Conditional network	NA	NA	NA	NA	15.99	
	RBF SVM (Vinyals et al. 2012)	SVM	NA	NA	NA	NA	21	
	SVM	R2SVM (Vinyals et al. 2012)	Stacked SVM	NA	NA	NA	NA	20.3
DNN+SVM (Tang 2013)	DNN with SVM as last layer	NA	NA	284,106	NA	11.9		
HAR	Tree	RF (Breiman 2000)*	NA	100	NA	13.70	7.2	
	DART (Rashmi and Gilad-Bachrach 2015)*	XGBoost with regularizer	6	NA	NA	32.58	6.37	
	mGBDT (Feng et al. 2018)*	Multi-layer XGBoost	5	NA	NA	399.97	7.57	
	RRF (Deng and Runger 2012)*	RF with regularizer	NA	100	NA	74.50	3.82	

Table 3

The list of abbreviations and their descriptions utilized in this survey.

Abbreviation	Description
AE	Autoencoder
ANT	Adaptive Neural Tree
CNN	Convolutional Neural Network
CondNN	Conditional Neural Network
DART	Dropout Multiple Additive Regression Trees
DBT	Differentiable Boundary Tree
DBN	Deep Belief Network
DCCA	Deep Canonical Correlation Analysis
Deep PCA	Deep principal components analysis
DF	Deep Forest
DGP	Deep Gaussian Processes
DKF	Deep Kalman Filters
DNDT	Deep Network Decision Tree
DNDF	Deep Network Decision Forest
DNN	Deep Neural Network
DSVM	Deep SVM
DT	Decision tree
DTA-LS-SVM	Deep Transfer Additive Kernel Least Square SVM
eForest	Encoder Forest
FC	Fully Connected
FSDT	Frosst Soft Decision Tree
GAF	Generative Adversarial Forest
GAN	Generative Adversarial Network
GRRF	Guided Regularized Random Forest
LMM	Level-wise Mixture Model
mGBDT	Multilayer Gradient Decision Tree
ML-SVM	Multilayer SVM
MLP	Multilayer perceptron
NLP-SVM	Newton Linear Programming SVM
NPL	Neural Prototype Learning
R2-SVM	Random Recursive SVM
RBM	Restricted Boltzmann Machine
RNN	Recurrent Neural Network
RRF	Regularized Random Forest
SCAD-SVM	Smoothly Clipped Absolute Deviation SVM
SDF	Siamese Deep Forest
SNN	Siamese Neural Network
TAO	Tree Alternation Optimization
VAE	Variational Autoencoder