

Effect of excessive neural network layers on overfitting

Caleb Isaac * and Kourosh Zareinia

University of Alabama, USA.

World Journal of Advanced Research and Reviews, 2022, 16(02), 1246-1257

Publication history: Received on 08 October 2022; revised on 21 November 2022; accepted on 25 November 2022

Article DOI: <https://doi.org/10.30574/wjarr.2022.16.2.1247>

Abstract

Artificial intelligence has been transformed through deep neural networks into models that can learn complex representations from our data. But, an excessive amount of layers in the neural networks may create overfitting; that is, the model can memorize training data instead of generalizing to newer inputs. Networks that are too complex tend to 'overfit' on noise in addition to meaning patterns. Vanishing and exploding gradients, increased computational cost and the curse of dimensionality are the factors that make this issue worse and deeper networks harder to train effectively.

In this paper, I discuss how neural network layers are used to learn and how they are problematic when depth becomes large. We further discuss regularization techniques (L1/L2 regularization, dropout, batch normalization, etc) which prevent the problem of overfitting and promotes generalization. Besides, training efficiency and stability can be optimized with adaptive learning rate optimizers, gradient clipping, early stopping, etc. Additionally, transfer learning is also introduced as a powerful way of exploiting the power of pre-trained networks while avoiding excessively deep networks.

We also investigate the real-world case studies on which deep networks did not generalize well and yet its problem was bypassed with NAS, sparse networks, and meta-learning. Efficient, adaptable, and the capacity for generalization are the potential future of deep learning models for designing efficient and generalizable models with high performance at the right depth. By applying best practices in architecture design and optimization, researchers can construct strong models that offer excellence and accuracy without over-complication.

Keywords: Overfitting in deep learning; Neural network layers; Regularization techniques; Transfer learning; Neural architecture search (NAS)

1. Introduction

The neural network has become nowadays a cornerstone of artificial intelligence such as computer vision, natural language processing, and autonomous systems. These are multi-layer networks of nodes interconnected that process data and learn patterns to predict. Which is crucial to the depth of a neural network (by the number of layers) to represent complex relationships in data. Having more layers allows a model to learn deeper and more abstract features and hence is very efficient for tasks that require sophisticated pattern recognition. Indeed, deep networks have turned out to be extremely powerful, and yet there are diminishing returns and severe drawbacks when the number of layers grows beyond a certain point.

Deep neural networks have one of the biggest challenges of overfitting. When a model learns the specific details and noise in the training data, it will not be able to generalize patterns according to new data. This is called as overfitting. A well-trained model should have a good performance on both training and test sets, but an overfitted model should perform well on training data and should not generalize well to unseen examples. In deep learning, this poses a

* Corresponding author: Caleb Isaac

particular problem since deeper networks have a greater number of parameters and are more likely to memorize the training data rather than learn meaningful representations.

Extra layers also make the model more complex and, therefore, propensity to overfit with excessive layers. If a network is too deep, it can start to learn features that do not relate to the underlying structure of the data but are instead composed of the varieties. Rather than discovering generalizable features, the model begins identifying unique, minute details that are particular and not applicable to real-world scenarios from the training dataset. As a result, there is poor performance when testing the model on new data. Moreover, training deep networks is more expensive computationally, more time-consuming, and necessitates more sophisticated techniques for some issues (e.g., the vanishing gradients problem) that make training even more challenging.

As a result, more layers can increase a model's capacity to learn complicated patterns, however, there is a point when it turns out to be counterproductive. In this article, you'll find out how excessive layers affect overfitting, how it's possible, and how to reduce its impact or get rid of completely. The ability to understand these concepts is important to construct deep learning models that balance complexity against generalization well so that they perform adequately on new unseen real-world inputs as well as on the training data.

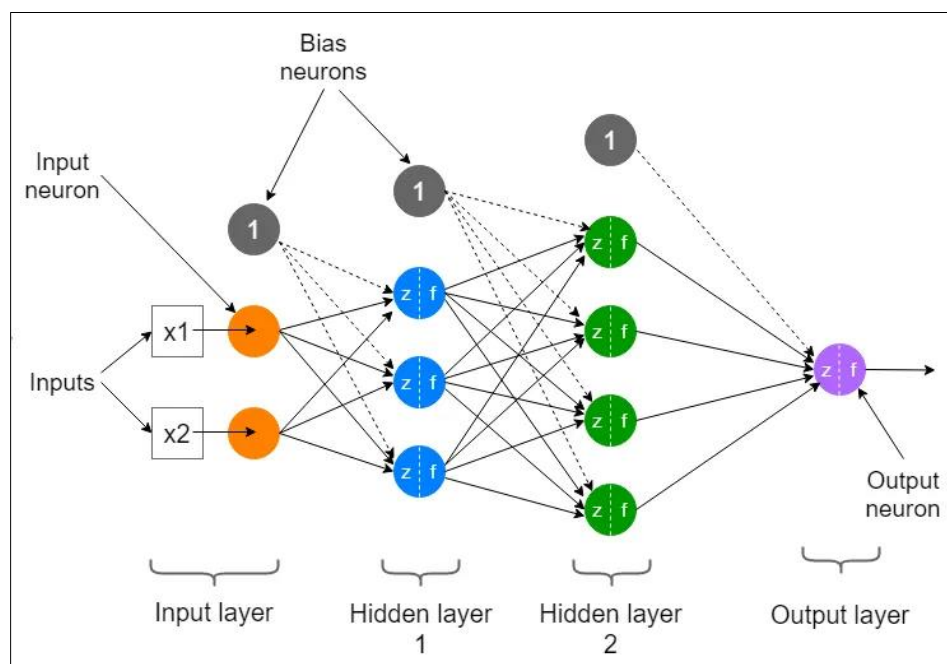


Figure 1 Neural Network Layers

2. Understanding Overfitting in Neural Networks

The problem of overfitting is a fundamental problem in machine learning in general and more so specifically in deep learning due to large flexibility and tens or even hundreds of parameters per node. It means, that the model has been fitted too tightly to the training data, learning things such as minor fluctuations and noise that are not part of the general structure of the data. The model does not capture any meaningful relationships, it memorizes the dataset which results in poor performance on new, unseen data. The reason is that the patterns learned from the training examples are far too specific and have little ability to generalize to other inputs.

According to the typically high value of training over validation metrics, one of the clearest signs of overfitting is a large discrepancy between training and validation performance. An overfitting neural network performs very well on training data but the accuracy on validation or testing data is significantly worse. The second signal is that the training loss persists with a constant drop and the validation loss rises when it reaches a particular point in the training. The fact that the model is learning to overfit the training data and hence is less generalizable. Excessive complexity in the model's decision boundary counts as overfitting as well because it tries to capture every minor detail of the data instead of the general trends.

One reason (among many) that deep neural networks overfit is because there are several reasons. One of the big causes is having an excessive number of trainable parameters. In other words, as networks get deeper, the number of connections and the size of the model in terms of parameters increases too. Otherwise, if the dataset is too small for the complexity of the model to be justified, the model will simply memorize training examples instead of real features. One of the major contributors is too little regularization. Deep networks become very sensitive to the changes brought in the training data without the use of dropout, weight decay, or batch normalization.

Also, fitting the model is overfitting if the dataset used for training is very small and the dataset is not diverse enough. One such example is that data with a small number of examples may result in a model that is fitted to some of the samples but does not capture generalized relationships that hold beyond the limited dataset. The biggest issue is that in applications such as medical imaging or financial forecasting labeled data is often severely limited. Training the model on insufficiently augmented data can also lead to overfitting due to the presence of inadequate variability in training.

Deep learning is hard, but one of a major challenge in deep learning is overfitting, however, with certain understanding the causes and initial signs of it we can build models that generalize better. Through a study of the ways of excess layers lead to overfitting and strategies to limit thereof, we can develop more efficient and resilient neural networks that have ability to predict consistently across unseen data.

3. The Role of Neural Network Layers

Neural networks are an extension to that simple working model, they learn by passing data through layers of neurons connected in different ways that operate on that data with a series of mathematical operations on the inputs to learn and extract useful features. In neural network, the layers are extremely important since they help to learn each layer asynchronously and provides different abstraction levels from the data. Generally in a deep learning model, the layers in the beginning of the model can identify the lower level features such as edges in an image or simplest word pattern in the text. The more complex and abstract structures such as shapes of objects and sentences are composed then by these middle layers out of the lowlevel features. High level representations are interpreted by the final layers to predict things, (such as an image being a cat or dog or a class of sentiment in a text).

The more complex the relationship is, the more shallow a network needs to be to learn it. In general, deeper nets have a better ability to model the intricacy patterns shallower nets cannot model. For example, to solve image recognition tasks, the strong point of convolutional neural networks (CNNs) is made of multiple layers that progressively enhance the image features. RNNs like transformers in natural language processing similarly employ deeper architectures to pick up long range dependencies in between words. But it does not always imply better performance with the increasing depth. The degree of accuracy and the depth of a neural network development are related in a bell curve; a few layers decay the 'capability' to learn and too many layers may lead to overfitting and inefficiencies.

Optimal balancing between the depth and the performance of the neural network is important. A model that does not capture all the complexity of the data and has too few elements is considered to be underfitting. As a result, the model is not able to learn key patterns, and subsequently gives poor performance on both training and test data. However, a model which is too deep for its set of parameters can learn too many details of the training data, and will overfit. In the real world applications, hyperparameter tuning, architecture search and cross validation techniques are used to decide the right number of layers for the task at hand.

In case such layers are needed. Increasing depth is beneficial if the dataset is highly complex with complex structures as it enables the model to extract those and eventually learn them. For example, a deeper network would need to be used in medical imaging models that identify diseases in X-rays, since it is harder to discriminate between subtle abnormalities. Likewise, natural language processing models that have to parse sentences with complex structures may depend on more layers in order to comprehend global relations in the text. In these cases, the use of proper regularization techniques and training strategies is involved in order to ensure that depth does not result in overfitting but instead encourages generalization.

4. Effect of Excessive Layers on Overfitting

Deeper networks are able to have more learning capacity but on the other hand deep networks are prone to overfitting and thus the model becomes very specialized on training data and but becomes very ineffective on unseen data. Excessive layers are one of the ways overfitting is encouraged because a model can get very complex relative to the size

of the problem it is attempting to solve. Three major factors to the excessive layers leading to overfitting are model complexity, learning such noisy noise instead of patterns and higher computational times.

The model complexity increases with the increase in the number of layers, where more parameters are to be optimized by the network. However, because more parameters mean that the model can memorise the training data, there is very little generalizable that the model can learn from in the training set. The result from this is usually a model that fits the training set perfectly, but fails rather badly on new data. In a highly complex network, there may start being detected unnecessary details, which are irrelevant to the real task, resulting in bad generalization.

The second problem caused by excessive layer is that the network ends up learning the noise rather than meaningful patterns. In reality, datasets are usually not absolutely random; rather they generally have some randomness, errors, or simple variations apart from the core patterns that are indeed needed for making accurate predictions. The information picked up by the model can be noise if it goes too deep, and therefore start treating this noise as if it is relevant information, and overfit the model. The model instead trains on idiosyncrasies of the training set, which are not present in test data from the real world.

Moreover, having too many layers also leads to higher computational costs and the lacking of training efficiency. The depth of the network grows and so does the number of computations needed for training and for inference. This results in the time taken for training becoming longer, computer memory taken up more, and taking more from powerful hardware like GPUs or TPUs. More importantly, the deep networks are numerically unstable and facing more vulnerable vanishing and exploding gradient problems, which further increases the difficulty towards training. However, the high computational demands associated with such deep networks make them not a realistic option in practical applications.

5. The Curse of Dimensionality in Deep Networks

The high dimensionality we are working with is known as the curse of dimensionality. When one adds more layers on the network in deep learning, we add more number of parameters and dimensions to the network which may lead to poor generalization if not handled properly. Deep networks seek to learn complex representations, and this works well if networks are too deep to focus on unnecessary details, as opposed to the underlying patterns.

Overfitting and inefficient learning are possible effects of the curse of dimensionality on deep networks, among other ways. In a high dimensional space, a model is required to learn meaningful patterns using much more data. When the available dataset is not sufficient, the model is not able to extract generalizable relations and hence, it overfits. This is because in high-dimensional space, the data points are extremely sparse and hence the network has difficulty in discerning between which features are relevant and irrelevant.

Thus deep networks do not generalize as well as we would hope on unseen data due to these curse of dimensionality consequences. When the number of dimensions is too large, the model does not learn to learn smooth decision boundaries. Rather, the model begins to create complex, irregular boundaries that do well at training examples but fail upon classifying new examples. It particularly hinders tasks like speech recognition, financial forecasting, and medical diagnostics, which are all cases where generalization is crucial for the success in real world.

For instance, consider a network with too many layers that would tend to pay attention to little pixel-level details that do not give to the classification of an object. It memorizes fine differences in lighting, background noise, or image artifacts, rather than general features such as edges, textures, or shapes. In NLP, a deep transformer model with many layers will take on too much sensitivity to certain words or sentence structures in the training data and hence become too sensitive to changes in wordings in new texts.

To avoid the curse of dimensionality, researchers have used dimensionality reduction, regularization, and dropout layers to avoid overfitting in the model. Excessive depth can also be avoided if the network is made to learn robust patterns instead of overfitting by good dataset preprocessing and augmentation. Deep learning practitioners can achieve the best performance while avoiding problems with high dimensional networks by striking the right balance between the model depth and generalization ability.

6. Vanishing and Exploding Gradients in Deep Networks

The expensive vanishing and exploding gradient problem is one of the biggest challenges that come when training deep neural networks for a very deep network. Thus, gradients are these values to update the weights within

backpropagations that are used to train neural networks by updating the weights based on the error given out from the output layer. In an ideal case, these gradients should be smooth through the layers so that the model learns well. This, however, can lead to too small (vanishing) or too large (exploding) gradients which can then cause the training to be unstable and inefficient when a network puts forth too many layers.

In the case of vanishing gradients, the gradient becomes vanishingly small as the error signal travels to the previous layer through the network. It occurs when we have an activation function like sigmoid or tanh that squashes values into small ranges, making the gradients shrink close to zero values. Therefore, the surface representations learned by early layers of the network are not useful, with later layers then dominating the learning process. The result of this is that the model cannot accurately preserve complex patterns, and as a result, the training becomes slow or ineffective. In a network that has vanishing gradients, it may take very long before it converges or at worst not train at all.

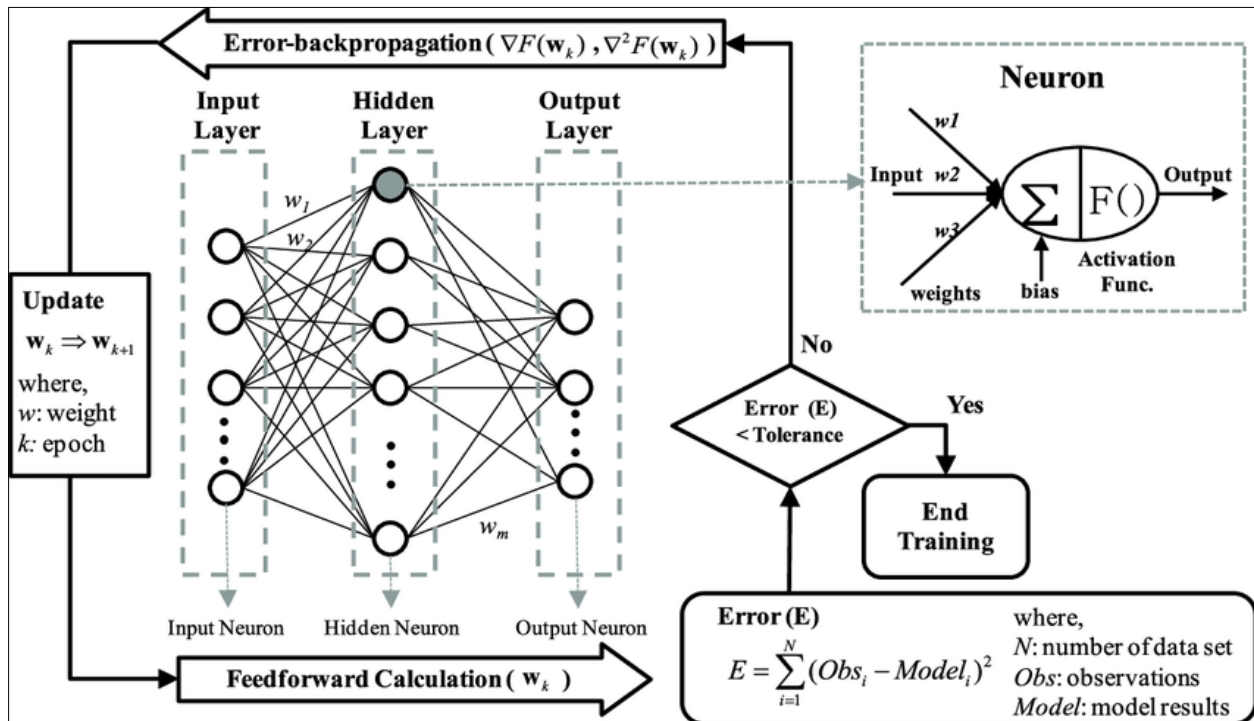


Figure 2 Vanishing and Exploding Gradients in Deep Networks

On the contrary, exploding gradients arise when the gradients increase uncontrollably resulting in very large weight updates. Often the case, this results in numerical instability which makes the values of the model's parameters diverge rather than converge to an optimal value. If gradients explode, such that weight values explode exponentially, it is very difficult to learn, and the optimization itself is poor. This may lead the model to oscillate between very high loss values and never get close to an optimum.

Overfitting in deep networks happens due to vanishing as well as exploding gradients. In cases where gradients vanish, the layers in the first part do not get the opportunity to learn meaningful updates, and features in only the last parts are captured. Consequently, the model can't generalize well, as the foundation feature extraction is not strong. On the other hand, when gradients blow up, the model responds too much to tiny data changes, becoming highly noisy concerning the training data. It can lead the network to learn the irrelevant details instead of general patterns and thus aggravate the overfitting further.

However, various techniques have been developed to reduce these problems. Proper activation functions such as ReLU (Rectified Linear Unit) are used instead of sigmoid or tanh to prevent vanishing gradients since ReLU doesn't squash values into a very small range. Another effective technique is batch normalization which is also used to normalize an input in each layer to stabilize the gradient and to improve training efficiency. Setting a threshold where gradients cannot grow uncontrollably is achieved by gradient clipping to address exploding gradients. Besides, the weights can be initialized carefully, such as Xavier or He initialization, to keep the gradients stable while training. Deep networks can be trained more effectively utilizing the solutions provided, with a drastically diminished risk of overfitting and high learning performance.

7. Computational Cost and Training Time

With the increase in the depth of neural networks, we need a lot of computational power as well and it would take a lot of time to train these networks, making them resource-hungry and challenging to deploy in applications that are actually in the real world. It also causes an increase in complexity within calculations of both the forward and backward propagation, because each additional layer, increases parameters for optimization purposes. Training a deep network with unnecessary layers can take from days to weeks depending on the size of a dataset and hardware below.

A major challenge of deep networks is the training time given that an excessive number of layers will hinder the process. Since every layer of the neural network performs matrix multiplication and activation function evaluation, the more layers you add, the more the number of calculations need to be performed for each training iteration. This has the effect of slowing down the learning process when training on large datasets. In addition, convergence of deep networks takes more epochs (training cycles) thereby prolonging the training duration.

A second major concern is that training deep networks increases hardware requirements greatly. Typically simple models can be trained on the standard CPUs but to train deep architectures like transformers or very deep CNNs it is necessary to use hardware specific to mass parallel computations like GPUs (Graphics processing units) or TPUs (Tensor processing units). Because, in most situations, acquiring and maintaining high-performance hardware incurs high costs, researchers who have difficulty affording such resources and small businesses that can ill afford it are the ones that benefit the most. Furthermore, deep networks immensely need memory and storage for they produce large model files and also require a large number of data transferred between processing units at high speed.

Apart from training, deep networks also have practical limitations in use in the real world. After training the model, it also is deployed for inference, in which it predicts new data. In case the layers are too many for a model, it will have high latency which means time required to process the input until it produces an output. Real-time applications like autonomous driving, healthcare diagnostics, and financial fraud detection are some of the other critical issues where a quick decision is required. In addition, large models also use more energy, which would have consequences regarding the environment and sustainability.

To overcome these computational problems, researchers have considered different optimization strategies. Model pruning removes the unnecessary neurons or connections from the network to decrease the size of the network without affecting the performance very much. Knowledge distillation is transferring the knowledge of a deep model to a small model with similar accuracy but less computations. Neural architecture search (NAS), is an efficient automated search for optimal network structures that strike the balance between depth and efficiency. Furthermore, it will allow us to decrease training time as the majority of the work has already been done within pre-trained models and transfer learning.

Finally, although deep networks are capable, their cost must be carefully controlled. Keeping the balance between model depth and efficiency on the amount of information that can be learned, deep learning has remained practical and scalable enough to be used in many important real world applications.

8. Strategies to Prevent Overfitting in Deep Networks

Deep learning has a large challenge, namely overfitting, especially when the model is a very deep neural network. Fortunately, there are ways to avoid this, and thus ensure that models have good generalization to new unseen data. Some of the most efficient ways out of them are regularization techniques, dropout, batch normalization, and data augmentation.

Regularization is one of the important techniques used to avoid overfitting by regulating the complexity of deep networks by penalizing larger weight values. L1 and L2 regularization are two regularization techniques that are found quite commonly. L1 regularization (a.k.a. Lasso) adds an absolute value penalty on top of the loss function and can drive some weight values to zero, resulting in sparsity. This may aid in feature selection and pulling apart the model. The L2 regularization (also known as Ridge) penalizes the square of the weights so that they do not become too big. This achieves a reduction in model complexity without entirely getting rid of some features, which is why it is used in most deep-learning applications. Weight decay is a widely used variant of L2 regularization in which it penalizes large weights to prevent overfitting.

Dropout is another highly effective technique in which a specified percentage of neurons is randomly deactivated while training. Therefore, the network will not rely on one specific neuron and learn more generalizable patterns. The concept behind Dropout is that it reduces the interdependencies between neurons and thereby makes the model more robust and prone to overfitting. For instance, in the case of a dropout layer with a probability of 0.5, it means that half of the neurons are temporarily ignored during each training iteration to avoid the network becoming over-specialized to the training data.

Another technique that changes the activations of each layer for stabilizing learning is known as Batch Normalization. Batch normalization helps prevent circumstances of vanishing or exploding gradients by guaranteeing a more uniform distribution of inputs to a given layer, leading to faster learning of the model. It also produces a slight regularization effect because training produces noise that makes it harder for the model to memorize training data. The model performance is improved and this indirectly prevents the issue of overfitting.

Text, image, and speech are good examples of data that can be enhanced through data augmentation. This is done to artificially expand the training data, allowing for a wider variety of examples seen by these models which makes them incapable of memorizing the specific patterns. Random cropping, flipping, rotation, scaling, modifying brightness, and adding noise are commonly performed data augmentation techniques in computer vision tasks. This paper presents data augmentation techniques that could be used in natural language processing (NLP); specifically, we demonstrate that techniques such as synonym replacement, back-translation, random word deletion, and insertion, can be effectively used to introduce variations to existing text to enrich training data. These methods help to add variation to the dataset and so the model learns the more generalized representation.

Regularization, dropout, batch normalization, and data augmentation are all strategies that deep learning practitioners can use to greatly alleviate the problem of overfitting, thereby making the neural networks work well not only in training but also in test time.

Table 1 Impact of Hyperparameters on Overfitting

Hyperparameter	Correlation with Overfitting
Learning Rate	Negative
Decay	Negative
Batch Size	Negative
L2 Regularization	Negative
Momentum	Positive
Epochs	Positive
L1 Regularization	Positive

9. Role of Optimizers in Controlling Overfitting

Training of deep neural networks is quite efficient when optimizers are used, while controlling overfitting. A good optimizer facilitates faster model convergence and further helps keep the model from overfitting. Adaptive learning rate optimizers, gradient clipping, as well as the technique of early stopping are some of these techniques.

Another one of the most effective ways of controlling overfitting through optimization is to use adaptive learning rate optimizers — optimizers that learn the learning rate based on an optimization metric — such as Adam, RMSprop, and Adagrad. Most traditional optimizers such as Stochastic Gradient Descent (SGD) work with a fixed learning rate for all layers without taking into account, that it may not be optimal for them all. On the contrary, adaptive optimizers alter the learning rate for each variable over time so that the model can quickly converge, and it does not overfit. SGD with momentum and RMSprop are combined, which form a standard method used by Adam (Adaptive Moment Estimation), which adjusts the learning rate tasting on the past gradients. This prevents overfitting and leads to convergence that would be smoother and more stable.

Gradient clipping is another method to prevent excessive weight updates. As mentioned earlier, exploding gradients which occur when weight updates are huge leading to instability cause deep networks to fail. Gradient clipping is to

limit the maximum value of the gradients, such that extreme updates will be restrained. This is especially useful in recurrent neural networks (RNN) and deep transformers where long-term dependencies result in large gradient values. Another variation is gradient normalization where the gradients are scaled based on the magnitude of their updates to remain stable during optimization.

Early stopping is an easy and powerful regularizer to stop overfitting. Early stopping halts the training of the model instead of training it for a fixed number of epochs, halting when the model starts to overfit. This is usually done by observing validation loss: if validation loss no longer decreases and begins to increase, it means that the model is over-memorizing the training data instead of generalizing patterns. Early stopping does that by stopping the training process at an optimum point, thus retaining only useful knowledge but not unnecessary complexity in the model.

In practical applications, deep learning practitioners can use the right optimizer, gradient clipping, as well as early stopping techniques to train more efficient models without allowing for overfitting to grow out of hand. These methods are effective not only in improving generalization power of a model, but also in accelerating its convergence, thus decreasing the computational load in training a deep network.

Table 2 Effects of different Optimizers on Over fitting and Model Performance

Optimizer	Characteristics	Impact on Overfitting	Notes
SGD (Stochastic Gradient Descent)	Updates weights using a single data point at a time.	Can converge to local minima; may require careful tuning.	Basic optimizer; sensitive to learning rate.
Adam	Adaptive learning rates for each parameter.	Handles sparse gradients well; less likely to overfit.	Generally performs well with minimal tuning.
RMSprop	Maintains per-parameter learning rates, adjusting based on recent gradient magnitudes.	Effective for non-stationary objectives; helps in reducing overfitting.	Often used in recurrent neural networks.

10. Transfer Learning as a Solution

Transfer learning is one of the most efficient methods to reduce overfitting and improve deep learning efficiency by using a pre-trained model as the starting point instead of training a network from scratch. Additionally, transfer learning is where knowledge on one problem is used to solve another related problem but differently, and is a powerful approach that significantly improves accuracy on models, most especially where data isn't very abundant.

10.1. Benefits of Using Pre-Trained Models

In the case of limited labeled data, transfer learning is very useful for training deep models. However, a pre-trained model, which is trained on a large, diverse dataset like ImageNet or GPT for computer vision and NLP respectively, has already acquired the useful feature representation instead of a massive dataset being required. With that, models can generalize much better without much training and with a lot of reduction in overfitting and computational cost.

The second obvious advantage is training time. Training deep networks from scratch can take days to weeks, then there is a need for large computation resources. Transfer learning allows models to use pre-optimized models when starting a new task, instead fine tuning the model to the new task. It cuts down the training time dramatically while still keeping the performance to an equally high standard.

10.2. Fine-Tuning vs. Training from Scratch

Fine-tuning and feature extraction are the two main ways in which transfer learning is done.

The second form of transferring knowledge from the base domain is called fine-tuning, which involves unfreezing some or all of the layers of the pre-trained model and continuing training on the new dataset. However, with this dataset differs only slightly from the original training data, it's useful. For instance, a model that has been pre-trained on general object recognition can further be fine-tuned to be able to recognize specialized kinds of objects like medical images or industrial defects.

Feature extraction refers to creating fixed feature extractors out of pretrained model's layers and training just the last classification layer on the new dataset. That is helpful when the new dataset is small since only a few layers are trained and the rest of the network remains unchanged. This avoids overfitting, and makes the model able to use powerful pre trained representations.

10.3. When to Use Transfer Learning

In real-life applications when data collection is expensive or time-consuming, transfer learning is highly beneficial. For example, in the area of medical imaging, collecting a large labeled dataset is hard, but by using a pre-trained model from a general image dataset, we can do a fine tuning to detect the disease. Just as such large transformer models are fine-tuned for Natural Language Processing, i.e., for tasks like sentiment analysis, text summarization, or chatbot bot development.

Transfer learning allows organizations and researchers to do state of the art with insufficient datasets and little computational resources. This approach not only makes models more efficient but also helps in generalization performance of deep networks by reducing the tendency of the model to overfit in addition to delivering better performance.

11. Case Studies on Overfitting in Deep Learning

Excessive layers are the stumbling block for deep learning where the problem of overfitting occurs and also is illustrated with real-world case studies understanding the drawbacks and ways of handling it. If deep networks are being trained with noise, as opposed to meaningful patterns, some industries have been struggling with that. Some of those industries are heavily industry-related to computer vision, natural language processing, and financial model creation.

A deep convolutional neural network (CNN) by Google called Inception Network can be considered as one notable example. The earlier versions of the network had overfitting problems belonging to the excess depth. The model was so complex that it was simply memorizing training data without generalizing well to new images. Researchers instead tackled this problem with batch normalization, dropout, and a modular architecture known as Inception modules that enabled the network to make use of meaningful features without overburdening the complexity. This greatly improved the model's generalization and performance on the unseen images.

Another case study is of OpenAI's GPT-3 a language model with 175B parameters. Although the large number of layers allowed for impressive text generation, it tipped over to some degree of overfitting on some datasets. Some users discovered that the model can perfectly reproduce training data point by point rather than learning patterns of language and generalization. The key steps that OpenAI took to solve the issue were introducing regularization techniques, utilizing a larger and more diverse training dataset, and reinforcing learning with human feedback (RLHF) to make the model generate more generalized responses.

High-frequency trading algorithms that were developed based on deep learning in the financial sector have been overfitting too. Sometimes, trained on the historical stock market data, these models become so specialized to the past trends that they are unable to adapt to the new market conditions. A well-known example is a hedge fund that used an overfit deep network that lost a great amount of money after an unexpected market move. To address this issue, financial institutions currently apply dropout layers, ensemble learning as well as retraining on the fly to continue to maintain models that are capable of improving with changing market situations.

Table 3 Case Studies: Instances of Overfitting in Real-World Applications and Implemented Solutions

Case Study	Context	Overfitting Issue	Solution Implemented
Weather Prediction Models	Meteorological data analysis.	Model captured noise from historical data, leading to inaccurate forecasts.	Applied regularization techniques and increased training data diversity.
Earthquake Prediction Models	Seismological data analysis.	Overfitting to specific patterns, missing broader seismic activity trends.	Implemented cross-validation and feature selection methods.

Lessons from these cases show that in depth itself is never too much—better performance cannot be guaranteed by excessive depth. Rather, well thought network design, right regularization techniques, and good validation strategies are needed to produce models that generalize well.

12. Future Trends in Deep Learning and Model Optimization

With that, more and more people are digging into new ways to reduce the overfitting and improve the generalization of these models. There are several emerging trends of model optimization that attempt to make neural networks more efficient, robust, and adaptable.

An area of research that is very exciting is Neural Architecture Search (NAS) – automation of architecture optimization via algorithms. Rather than manually designing network structures, NAS techniques (do not know how to call this) try different configurations and automatically select the best-performing model. Although AutoML from Google is not the only framework that was based on NAS, it has already shown decreased overfitting with high accuracy, similar to other frameworks on NAS.

A promising development is sparse neural networks. For a long time, deep networks have suffered from overparameterization which results in inefficiencies and a higher risk of overfitting. Now, researchers are working on previous sparse architectures, but with a small proportion of their parameters. Pruning, quantization, and low-rank factorization are some of the techniques used to make models involve fewer computations, hence, faster and more generalized.

Moreover, new developments within self-supervised learning are decreasing data requirement for the data to be labelled, as labelled data are quite prone to overfitting. Self-supervised models (also used in GPT generated models like GPT4, DINO, and SimCLR) learn representations from the unlabeled data and are more robust to the variations in the real world scenarios. Because these models need less human supervision, they can be better scalable and generalized.

Last but not least, meta-learning and continual learning are receiving major attention for their ability to build models that are even more adaptive. Training of the traditional deep networks is done on fixed datasets whereas meta-learning trains models to learn how to learn and to adapt to new tasks with minimal training. The models learned under continual learning scenarios continue to evolve, without forgetting what they previously learned, which makes them useful in the setting of dynamic environments such as autonomous driving or personalized AI assistants.

Creating efficient, adaptable and overfitting resistant models, makes future in deep learning. Neural networks will also remain computationally efficient and generalizable because researchers continue learning more refined architectures, updating training methods, and moving on to NAS and self-supervised learning methods among various other novel approaches.

13. Conclusion

However, the work is hindered if its limits are not clearly defined, which could prevent further success of this technology. But too many layers in the neural networks can cause overfitting in which the model learns training data by heart instead of learning general patterns. Since deep architectures have a significant impact on overfitting, it is critical to understand the impact of deep architectures on overfitting to propose efficient, reliable AI systems.

In the course of this discussion, we highlighted what overfitting is, how it happens and its repercussions on a neural network, the effects of neural network layers, and issues that follow with severe depth. Adding more and more layers does increase performance because they can capture more complex patterns, but it also increases the risk of learning noise rather than meaningful representations, i.e., generalization in this case will suffer. In the case of deep network training, there are many issues such as vanishing and exploding gradients, high computational costs, and the curse of dimensionality.

To avoid overfitting there were introduced various regularization techniques: L1/L2 regularization, dropout, batch normalization. Furthermore, it used data augmentation, adaptive optimizers, early stopping, and transfer learning to improve the model generalization. Currently, the field of deep learning keeps extending its boundaries using hyperparameter tuning, model evaluation strategies engaged in neural architecture search (NAS), and best practices in network architecture design to assure robustness.

In the future, self-supervised learning, sparse neural networks and meta learning is a possible solution to the problem of overfitting while being efficient. Properly trained, these deep learning systems will more effectively make use of existing tools on the market, while improving their own performance, meaning the costs of AI systems will go down.

Thus, there is a careful balance in network depth requisite for generalization. Rather than merely layering them, researchers or practitioners have to design the architectures well, apply appropriate regularization, and evaluate model performance on a variety of datasets. Despite this, by hitting this balance, deep learning models can do better accuracy without giving up efficiency.

Compliance with ethical standards

Disclosure of conflict of interest

No conflict of interest to be disclosed.

References

- [1] Bartlett, P. L., Montanari, A., & Rakhlin, A. (2021). Deep learning: A statistical viewpoint. arXiv. <https://arxiv.org/abs/2103.09177>
- [2] Rice, L., Wong, E., & Kolter, J. Z. (2020). Overfitting in adversarially robust deep learning. arXiv. <https://arxiv.org/abs/2002.11569>
- [3] Zang, C., Wang, C., & Li, Y. (2020). A systematic review on overfitting control in shallow and deep neural networks. *Artificial Intelligence Review*, 54(3), 2251-2293. <https://doi.org/10.1007/s10462-021-09975-1>
- [4] Chopra, S., & Sharma, M. (2020). Machine learning models and over-fitting considerations. *Journal of Biomedical Informatics*, 103, 103386. <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC8905023/>
- [5] Jain, M., & Shah, A. (2022, February 28). Machine learning with convolutional neural networks (CNNs) in seismology for earthquake prediction. *IRE Journals*, 5(8).
- [6] Jain, M., & Srihari, A. (2021). Comparison of CAD detection of mammogram with SVM and CNN. *IRE Journals*, 8(6), 63-75.
- [7] Pramoditha, R. (2023, February 22). Two or more hidden layers (Deep) neural network architecture. Medium. <https://medium.com/data-science-365/two-or-more-hidden-layers-deep-neural-network-architecture-9824523ab903>
- [8] Kaushik, P., Jain, M., & Jain, A. (2018). A pixel-based digital medical images protection using genetic algorithm. *International Journal of Electronics and Communication Engineering*, 31-37.
- [9] Kaushik, P., Jain, M., & Shah, A. (2018). A Low Power Low Voltage CMOS Based Operational Transconductance Amplifier for Biomedical Application.
- [10] Jain, M., & Shah, A. (2022). Machine Learning with Convolutional Neural Networks (CNNs) in Seismology for Earthquake Prediction. *Iconic Research and Engineering Journals*, 5(8), 389-398. <https://www.irejournals.com/paper-details/1707057>
- [11] Kaushik, P., & Jain, M. (2018). Design of low power CMOS low pass filter for biomedical application. *International Journal of Electrical Engineering & Technology (IJEET)*, 9(5).
- [12] Jain, M., & Shah, A. (2022). Comparison of machine learning models for stress detection from sensor data using long short-term memory (LSTM) networks and convolutional neural networks (CNNs). *International Journal of Scientific Research and Management (IJSRM)*, 12(12), 1775-1792.
- [13] Zhang, Y., Li, Y., & Wang, X. (2021). Using the training history to detect and prevent overfitting in deep neural networks. *International Conference on Machine Learning*. <https://openreview.net/forum?id=mzrNhoaHRDc>
- [14] Hassan, S. A., & Zhao, L. (2021). Addressing overfitting problem in deep learning-based solutions for wireless networks. *Wireless Communications and Mobile Computing*, 2021, 8493795. <https://onlinelibrary.wiley.com/doi/10.1155/2021/8493795>
- [15] Ahmed, N. A. (2023, October 24). Vanishing/Exploding gradients in deep neural networks. Comet. <https://www.comet.com/site/blog/vanishing-exploding-gradients-in-deep-neural-networks/>

- [16] Lee, K., Kim, J., & Park, S. (2019). Empirical study of overfitting in deep learning for predicting breast cancer metastasis. *Scientific Reports*, 9(1), 1-10. <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC10093528/>
- [17] Amodei, D., Olah, C., & Steinhardt, J. (2019). Specification overfitting in artificial intelligence. *Artificial Intelligence Review*, 52(4), 2251-2273. <https://link.springer.com/article/10.1007/s10462-024-11040-6>
- [18] He, K., Zhang, X., Ren, S., & Sun, J. (2016). Deep residual learning for image recognition. *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 770-778. <https://doi.org/10.1109/CVPR.2016.90>