

Input Margins Can Predict Generalization Too

Coenraad Mouton^{1,2,3}, Marthinus Wilhelmus Theunissen^{1,2}, Marelle H Davel^{1,2,4}

¹Faculty of Engineering, North-West University, South Africa

²Centre for Artificial Intelligence Research, South Africa

³South African National Space Agency

⁴National Institute for Theoretical and Computational Sciences, South Africa
{moutoncoenraad, tiantheunissen, marelle.davel}@gmail.com

Abstract

Understanding generalization in deep neural networks is an active area of research. A promising avenue of exploration has been that of margin measurements: the shortest distance to the decision boundary for a given sample or its representation internal to the network. While margins have been shown to be correlated with the generalization ability of a model when measured at its hidden representations (hidden margins), no such link between large margins and generalization has been established for *input margins*. We show that while input margins are not generally predictive of generalization, they can be if the search space is appropriately constrained. We develop such a measure based on input margins, which we refer to as ‘constrained margins’. The predictive power of this new measure is demonstrated on the ‘Predicting Generalization in Deep Learning’ (PGDL) dataset and contrasted with hidden representation margins. We find that constrained margins achieve highly competitive scores and outperform other margin measurements in general. This provides a novel insight on the relationship between generalization and classification margins, and highlights the importance of considering the data manifold for investigations of generalization in DNNs.

1 Introduction

Our understanding of the generalization ability of deep neural networks (DNNs) remains incomplete. Various bounds on the generalization error for classical machine learning models have been proposed based on the complexity of the hypothesis space (Vapnik 1999; Koltchinskii and Panchenko 2002). However, this approach paints an unfinished picture when considering modern DNNs (Zhang et al. 2021). Generalization in DNNs is an active field of study and updated bounds are proposed on an ongoing basis (Arora et al. 2018; Kawaguchi, Kaelbling, and Bengio 2022; Chuang et al. 2021; Lotfi et al. 2022).

A complementary approach to developing theoretical bounds is to develop empirical techniques that are able to predict the generalization ability of certain families of DNN models. The ‘Predicting Generalization in Deep Learning’ (PGDL) challenge, exemplifies such an approach. The challenge was held at NeurIPS 2020 (Jiang et al. 2020) and provides a useful test bed for evaluating *complexity measures*,

where a complexity measure is a scalar-valued function that relates a model’s training data and parameters to its expected performance on unseen data. Such a predictive complexity measure would not only be practically useful but could lead to new insights into how DNNs generalize.

In this work, we focus on classification margins in deep neural classifiers. It is important to note that the term ‘margin’ is, often confusingly, used to refer to 1) output margins (Bartlett, Foster, and Telgarsky 2017), 2) input margins (Sokolić et al. 2017), and 3) hidden margins (Jiang et al. 2018), interchangeably. Here (1) is a measure of the difference in class output values, while (2) or (3) is concerned with measuring the distance from a sample to its nearest decision boundary in either input or hidden representation space, respectively. We limit our focus to input and hidden margins.

While margins measured at the hidden representations of deep neural classifiers have been shown to be predictive of a model’s generalization, this link has not been established for input space margins. We show that, in several circumstances, the classical definition of input margin does *not* predict generalization, but a direction-constrained version of this metric does: a quantity we refer to as *constrained margins*. By measuring margins in directions of ‘high utility’, that is, directions that are expected to be more useful to the classification task, we are able to better capture the generalization ability of a trained DNN.

We make several contributions:

1. Demonstrate the first link between large input margins and generalisation performance, by developing a new input margin-based complexity measure that achieves highly competitive performance on the PGDL benchmark and outperforms several contemporary complexity measures.
2. Show that margins do not necessarily need to be measured at multiple hidden layers to be predictive of generalization, as suggested in (Jiang et al. 2018).
3. Provide a new perspective on margin analysis and how it applies to DNNs, that of finding high utility directions along which to measure the distance to the boundary instead of focusing on finding the shortest distance.

2 Background

This section provides an overview of existing work on 1) measuring classification margins and their relationship to generalization, and 2) the PGDL challenge and related complexity measures.

2.1 Classification Margins and Generalization

Considerable prior work exists on understanding classification margins in machine learning models (Boser, Guyon, and Vapnik 1992; Weinberger and Saul 2009). The relation between margin and generalization is well understood for classifiers such as support vector machines (SVMs) under statistical learning theory (Vapnik 1999). However, the non-linearity and high dimensionality of DNN decision boundaries complicate such analyses, and precisely measuring these margins is considered intractable (Yousefzadeh and O’Leary 2020; Yang et al. 2020).

A popular technique (which we revisit in this work) is to approximate the classification margin using a first-order Taylor approximation. Elsayed et al. (2018) use this method in both the input and hidden space, and then formulate a loss function that maximizes these margins. However, while this results in a measurable increase in margin, it does not result in any significant gains in test accuracy. In a seminal paper, Jiang et al. (2018) utilize the same approximation in order to predict the generalization gap of a set of trained networks by training a linear regression model on a summary of their hidden margin distributions. Natekar and Sharma (2020) demonstrate that this measure can be further improved if margins are measured using the representations of Mixup (Zhang et al. 2018) or augmented training samples. Similarly, Chuang et al. (2021) introduce novel generalization bounds and slightly improve on this metric by proposing an alternative cluster-aware normalization scheme (k -variance (Solomon, Greenewald, and Nagaraja 2022)).

Input margins are generally considered from the point of view of adversarial robustness, and many techniques have been developed to generate adversarial samples on or near the decision boundary. Examples include: the Carlini and Wagner Attack (Carlini and Wagner 2017), Projected Gradient Descent (Madry et al. 2018), and DeepFool (Moosavi-Dezfooli, Fawzi, and Frossard 2016). Some of these studies have investigated the link between adversarial robustness and generalization, often concluding that an inherent trade-off exists (Tsipras et al. 2019; Su et al. 2018; Raghu et al. 2019). However, this conclusion and its intricacies are still being debated (Stutz, Hein, and Schiele 2019).

Yousefzadeh and O’Leary (2020) formulate finding a point on the decision boundary as a constrained minimization problem, which is solved using an off-the-shelf optimization method. While this method is more precise, it comes at a great computational cost. To alleviate this, dimensionality reduction techniques are used in the case of image data to reduce the number of input features.

In this work we propose a modification to the Taylor approximation of the input classification margin (and its iterative alternative DeepFool) in order for it to be more predictive of generalization.

2.2 Predicting Generalization in Deep Learning

The PGDL challenge was a competition hosted at NeurIPS 2020 (Jiang et al. 2020). The objective of this challenge was to design a complexity measure to rank models according to their generalization gap. More precisely, participants only had access to a set of trained models, along with their parameters and training data, and were tasked with ranking the models within each set according to their generalization gap. Each solution was then evaluated on how well its ranking aligns with the true ranking on a held-out set of tasks, which was unknown to the competitors.

In total, there are 550 trained models across 8 different tasks and 6 different image classification datasets, where each task refers to a set of models trained on the same dataset with varying hyperparameters and resulting test accuracy. Tasks 1, 2, 4, and 5 were available for prototyping and tuning complexity measures, while Task 6 to 9 were used as a held-out set. There is no Task 3. The final average score on the test set was the only metric used to rank the competitors. Conditional mutual information (CMI) is used as evaluation metric, which measures the conditional mutual information between the complexity measure and true generalization gap, given that a set of hyperparameter types are observed. This is done in order to prevent spurious correlations resulting from specific hyperparameters, a step towards establishing whether a causal relationship exists.

All models were trained to approximately the same, near zero, training loss. Note that this implies that ranking models according to either their generalization gap or test accuracy is essentially equivalent.

Several interesting solutions were developed during the challenge: In addition to the modification of hidden margins mentioned earlier, the winning team (Natekar and Sharma 2020) developed several prediction methods based on the internal representations of each model. Their best-performing method measures clustering characteristics of hidden layers (using Davies-Bouldin Index (Davies and Bouldin 1979)), and combines this with the model’s accuracy on Mixup-augmented training samples. In a similar fashion, the runners-up based their metrics on measuring the robustness of trained networks to augmentations of their training data (Kashyap, Subramanyam et al. 2021).

After the competition’s completion, the dataset was made publicly available, inspiring further research: Schiff et al. (2021) generated perturbation response curves that ‘capture the accuracy change of a given network as a function of varying levels of training sample perturbation’ and develop statistical measures from these curves. They produced eleven complexity measures with different types of sample Mixup and statistical metrics.

While several of the methods rely on using synthetic samples (e.g. Mixup), Zhang et al. (2022) take this to the extreme and generate an artificial test set using pretrained generative adversarial networks (GANs). They demonstrate that simply measuring the classification accuracy on this synthetic test set is very predictive of a model’s generalization. While practically useful, this method does not make a link between any characteristics of the model and its generalization ability.

3 Theoretical Approach

This section provides a theoretical overview of the proposed complexity measure. We first explain our intuition surrounding classification margins, before mathematically formulating constrained margins.

3.1 Intuition

A correctly classified training sample with a large margin can have more varied feature values, potentially due to noise, and still be correctly classified. However, as we will show, input margins are not generally predictive of generalization. This observation is supported by literature regarding adversarial robustness, where it has been shown that adversarial retraining (which increases input margins) can negatively affect generalization (Tsipras et al. 2019; Raghuathan et al. 2019).

Stutz et al. (2019) provide a plausible reason for this counter-intuitive observation: Through the use of Variational Autoencoder GANs they show that the majority of adversarial samples leave the class-specific data manifold of the samples' class. They offer the intuitive example of black border pixels in the case of MNIST images, which are zero for all training samples. Samples found on the decision boundary which manipulate these border pixels have a zero probability under the data distribution, and they do not lie on the underlying manifold.

We leverage this intuition and argue that any input margin measure that relates to generalization should measure distances along directions that do not rely on spurious features in the input space. The intuition is that, while nearby decision boundaries exist for virtually any given training sample, these nearby decision boundaries are likely in directions which are not inherently useful for test set classification, i.e. they diverge from the underlying data manifold.

More specifically, we argue that margins should be measured in directions of 'high utility', that is, directions that are expected to be useful for characterising a given dataset, while ignoring those of lower utility. In our case, we approximate these directions by defining high utility directions as directions which explain a large amount of variance in the data. We extract these using Principal Component Analysis (PCA). While typically used as a dimensionality reduction technique, PCA can be interpreted as learning a low-dimensional manifold (Hinton, Dayan, and Revow 1997), albeit a locally linear one. In this way, the PCA manifold identifies subspaces that are thought to contain the variables that are truly relevant to the underlying data distribution, which the out-of-sample data is assumed to also be generated from. In the following section, we formalize such a measure.

3.2 Constrained Margins

We first formulate the classical definition of an input margin (Yousefzadeh and O'Leary 2020), before adapting it for our purpose.

Let $f : X \rightarrow \mathbb{R}^{|N|}$ denote a classification model with a set of output classes $N = \{1 \dots n\}$, and $f_k(\mathbf{x})$ the output value of the model for input sample \mathbf{x} and output class k .

For a correctly classified input sample \mathbf{x} , the goal is to find the closest point $\hat{\mathbf{x}}$ on the decision boundary between the true class i (where $i = \arg \max_k (f_k(\mathbf{x}))$) and another class $j \neq i$. Formally, $\hat{\mathbf{x}}$ is found by solving the constrained minimization problem:

$$\arg \min_{\hat{\mathbf{x}} \in [L, U]} \|\mathbf{x} - \hat{\mathbf{x}}\|_2 \quad (1)$$

with L and U the lower and upper bounds of the search space, respectively, such that

$$f_i(\hat{\mathbf{x}}) = f_j(\hat{\mathbf{x}}) \quad (2)$$

for i and j as above.

The margin is then given by the Euclidean distance between the input sample, \mathbf{x} , and its corresponding sample on the decision boundary, $\hat{\mathbf{x}}$. We now adapt this definition in order to define a 'constrained margin'. Let the set $P = \{\mathbf{p}_1, \mathbf{p}_2, \dots, \mathbf{p}_m\}$ denote the first m principal component vectors of the training dataset, that is, the m orthogonal principal components which explain the most variance. Such principal components are straightforward to extract by calculating the eigenvectors of the empirical covariance matrix of the normalized training data, where the data is normalized the same as prior to model training.

We now restrict $\hat{\mathbf{x}}$ to any point consisting of the original sample \mathbf{x} plus a linear combination of these (unit length) principal component vectors, that is, for some coefficient vector $\mathbf{B} = [\beta_1, \beta_2, \dots, \beta_m]$

$$\hat{\mathbf{x}} \triangleq \mathbf{x} + \sum_{i=1}^m \beta_i \mathbf{p}_i \quad (3)$$

Substituting $\hat{\mathbf{x}}$ into the original objective function of Equation (1), the new objective becomes

$$\min_{\mathbf{B}} \left\| \sum_{i=1}^m \beta_i \mathbf{p}_i \right\|_2 \quad (4)$$

such that Equation (2) is approximated within a certain tolerance and $\hat{\mathbf{x}} \in [L, U]$. For this definition of margin, the search space is constrained to a lower-dimensional subspace spanned by the principal components with point \mathbf{x} as origin, and the optimization problem then simplifies to finding a point on the decision boundary within this subspace. By doing so, we ensure that boundary samples that rely on spurious features (that is, in directions of low utility) are not considered viable solutions to Equation (1). Note that this formulation does not take any class labels into account for identifying high utility directions.

While it is possible to solve the constrained minimization problem using a constrained optimizer (Yousefzadeh and O'Leary 2020), we approximate the solution by adapting the previously mentioned first-order Taylor approximation (El-sayed et al. 2018; Huang et al. 2015), which greatly reduces the computational cost. The Taylor approximation of the constrained margin $d(\mathbf{x})$ for a sample \mathbf{x} between classes i and j when using an $L2$ norm is given by

$$d(\mathbf{x}) = \frac{f_i(\mathbf{x}) - f_j(\mathbf{x})}{\| [\nabla_{\mathbf{x}} f_i(\mathbf{x}) - \nabla_{\mathbf{x}} f_j(\mathbf{x})] \mathbf{P}^T \|_2} \quad (5)$$

where \mathbf{P} is the $m \times n$ matrix formed by the top m principal components with n input features.

The derivation of Equation (5) is included in the supplementary material.¹

The value $d(\mathbf{x})$ only approximates the margin and the associated discrepancy in Equation (2) can be large. In order to reduce this to within a reasonable tolerance, we apply Equation (5) in an iterative manner, using a modification of the well-known DeepFool algorithm (Moosavi-Dezfooli, Fawzi, and Frossard 2016). DeepFool was defined in the context of generating adversarial samples with the smallest possible perturbation, which is in effect very similar to finding the nearest point on the decision boundary with the smallest violation of Equation (2).

To extract the DeepFool constrained margin for some sample \mathbf{x} , the Taylor approximation of the margin in the lower-dimensional principal component subspace is calculated between the true class i (assuming that the sample is correctly classified) and all other classes j , individually. The smallest lower-dimensional subspace perturbation is then transformed back to the original feature space. This perturbation is then scaled by a set learning rate and added to the original sample. This process is repeated until the distance changes less than a given tolerance compared to the previous iteration. Note that the dimensionality of the sample $\hat{\mathbf{x}}$ is never reduced – only the search for a perturbation is restricted to the lower-dimensional principal component subspace. The exact process is described in Algorithm 1.

Note that we also clip $\hat{\mathbf{x}}$ according to the minimum and maximum feature values of the dataset after each step (line x in Algorithm 1), which ensures that the point stays within the bound constraints expressed in Equation (1). While this is likely superfluous when generating normal adversarial samples – they are generally very close to the original \mathbf{x} – it is a consideration when the search space is constrained, with clipped margins performing better. (See the supplementary material for an ablation analysis of clipping.)

4 Results

We investigate the extent to which constrained margins are predictive of generalization by comparing the new method with current alternatives. In Section 4.1 we describe our experimental setup. Following this, we do a careful comparison between our metric and existing techniques based on standard input and hidden margins (Section 4.2) and, finally, we compare with other complexity measures (Section 4.3).

4.1 Experimental Setup

For all margin-based measures our indicator of generalization (complexity measure) is the mean margin over 5 000 randomly selected training samples, or alternatively the maximum number available for tasks with less than 5 000 training samples. Only correctly classified samples are considered, and the same training samples are used for all models of the same task. To compare constrained margins to

Algorithm 1: DeepFool constrained margin calculation

Input: Sample \mathbf{x} , classifier f , principal components \mathbf{P}

Parameter: Stopping tolerance δ , Learning rate γ , Maximum iterations max

Output: Distance d_{best} , Equality violation v_{best}

```

1:  $\hat{\mathbf{x}} \leftarrow \mathbf{x}, i \leftarrow \arg \max_k f_k(\mathbf{x}), d \leftarrow 0, v_{best} \leftarrow \infty, c \leftarrow 0$ 
2: while  $c < max$  do
3:   for  $j \neq i$  do
4:      $o_j \leftarrow f_i(\hat{\mathbf{x}}) - f_j(\hat{\mathbf{x}})$ 
5:      $\mathbf{w}_j \leftarrow [\nabla f_i(\hat{\mathbf{x}}) - \nabla f_j(\hat{\mathbf{x}})]\mathbf{P}^T$ 
6:   end for
7:    $l \leftarrow \arg \min_{j \neq i} \frac{|o_j|}{\|\mathbf{w}_j\|_2}$ 
8:    $\mathbf{r} \leftarrow \frac{o_l}{\|\mathbf{w}_l\|_2^2} \mathbf{w}_l \mathbf{P}$ 
9:    $\hat{\mathbf{x}} \leftarrow \hat{\mathbf{x}} + \gamma \mathbf{r}$ 
10:   $\hat{\mathbf{x}} \leftarrow \text{clip}(\hat{\mathbf{x}})$ 
11:   $j \leftarrow \arg \max_{k \neq i} f_k(\hat{\mathbf{x}})$ 
12:   $v \leftarrow |f_i(\hat{\mathbf{x}}) - f_j(\hat{\mathbf{x}})|$ 
13:   $d \leftarrow \|\mathbf{x} - \hat{\mathbf{x}}\|_2$ 
14:  if  $v \geq v_{best}$  or  $|d - d_{best}| < \delta$  then
15:    return  $d_{best}, v_{best}$ 
16:  else
17:     $v_{best} \leftarrow v$ 
18:     $d_{best} \leftarrow d$ 
19:     $c \leftarrow c + 1$ 
20:  end if
21: end while
22: return  $d_{best}, v_{best}$ 

```

input and hidden margins we rank the model test accuracies according to the resulting indicator and calculate the Kendall’s rank correlation (Kendall 1938), as used in (Jiang et al. 2019). This allows for a more interpretable comparison than CMI. (As CMI is used throughout the PGDL challenge, we also include the resulting CMI scores in the supplementary material.) To compare constrained margins to published results of other complexity measures, we measure CMI between the complexity measure and generalization gap and contrast this with the reported scores of other methods.

As a baseline we calculate the **standard input margins** (‘Input’) using the first order Taylor approximation (Equation 5 without the subspace transformation), as we find that it achieves better results than the iterative DeepFool variant and is therefore the stronger baseline; see the supplementary material for a full comparison.

Hidden margins (‘Hidden’) are measured by considering the output (post activation function) of some hidden layer, and then calculating the margin at this representation. This raises the question of which hidden layers to consider for the final complexity measure. Jiang et al. (2018) consider three equally spaced layers, Natekar and Sharma (2020) consider all layers, and Chuang et al. (2021) consider either the first or last layer only. We calculate the mean hidden margin (using the Taylor approximation) for all these variations and find that for the tasks studied here, using the first layer performs best, while the mean over all layers comes in second. We include both results here (a full analysis is included in the

¹The supplementary material can be found at <https://arxiv.org/abs/2308.15466>

| Task | Architecture | Dataset | Constrained | Input | Hidden (1st) | Hidden (all) |
|---------|--------------|------------------------|---------------|---------|---------------|--------------|
| 1 | VGG | CIFAR10 | 0.8049 | 0.0265 | 0.5794 | 0.7825 |
| 2 | NiN | SVHN | 0.8686 | 0.6841 | 0.7037 | 0.8281 |
| 4 | FCN | CINIC10 | 0.6633 | 0.6251 | 0.7958 | 0.2707 |
| 5 | FCN | CINIC10 | 0.2282 | 0.3571 | 0.5427 | 0.1329 |
| 6 | NiN | OxFlowers | 0.8017 | -0.1351 | 0.4427 | 0.2839 |
| 7 | NiN | OxPets | 0.5133 | 0.3215 | 0.3623 | 0.3925 |
| 8 | VGG | FMNIST | 0.6004 | -0.1233 | -0.0656 | 0.1859 |
| 9 | NiN | CIFAR10 (augmented) | 0.8145 | 0.1573 | 0.7097 | 0.4556 |
| Average | | | 0.6617 | 0.2392 | 0.5088 | 0.4165 |

Table 1: Kendall’s rank correlation between mean margin and test accuracy for constrained, standard input, and hidden margins using the first or all layer(s). Models in Task 4 are trained with batch normalization while models in Task 5 are trained without. There is no Task 3.

| Task | Natekar and Sharma | | | Chuang et al. | | Schiff et al. | Ours |
|---------------------|--------------------|--------------|--------------|-----------------------|--------------------------|---------------|------------------------|
| | DBI*LWM | MM† | AM† | kV - Margin 1st† | kV -GN- Margin 1st† | PCA Gi&Mi | Constrained Margin† |
| 1 | 00.00 | 01.11 | 05.73 | 05.34 | 17.95 | 00.04 | 39.49 |
| 2 | 32.05 | 47.33 | 44.60 | 26.78 | 44.57 | 38.08 | 51.98 |
| 4 | 31.79 | 43.22 | 47.22 | 37.00 | 30.61 | 33.76 | 21.44 |
| 5 | 15.92 | 34.57 | 22.82 | 16.93 | 16.02 | 20.33 | 04.93 |
| 6 | 43.99 | 11.46 | 08.67 | 06.26 | 04.48 | 40.06 | 30.83 |
| 7 | 12.59 | 21.98 | 11.97 | 02.07 | 03.92 | 13.19 | 13.26 |
| 8 | 09.24 | 01.48 | 01.28 | 01.82 | 00.61 | 10.30 | 13.48 |
| 9 | 25.86 | 20.78 | 15.25 | 15.75 | 21.20 | 33.16 | 51.46 |
| Test set average | 22.92 | 13.93 | 09.29 | 06.48 | 07.55 | 24.18 | 27.26 |

Table 2: Conditional Mutual Information (CMI) scores for several complexity measures on the PGDL dataset. Acronyms: *DBI*=Davies Bouldin Index, *LWM*=Label-wise Mixup, *MM*=Mixup Margins, *AM*=Augmented Margins, *kV*=*k*-Variance, *GN*=Gradient Normalized, *Gi*=Gini coefficient, *Mi*=Mixup. Test set average is the average over Tasks 6 to 9. There is no Task 3. †Indicates a margin-based measure.

supplementary material, including results using DeepFool). We normalize each layer’s margin distribution by following (Jiang et al. 2018), and divide each margin by the total feature variance at that layer.

Our **constrained margin** complexity measure (‘Constrained’) is obtained using Algorithm 1, although in practice we implement this in a batched manner. Empirically, we find that the technique is not very sensitive with regard to the selection of hyperparameters and a single learning rate ($\gamma = 0.25$) and max iterations ($max = 100$) is used across all experiments. Furthermore, we use the same distance tolerance ($\delta = 0.01$) for all tasks, except for Tasks 4 and 5, which require a smaller tolerance ($\delta = 0.001$). This is because the features for this dataset (CINIC10 for both Tasks 4 and 5) are normalized to be in the range $[0, 1]$, while the features are z-normalized for the datasets of the other tasks. We find that Algorithm 1 generally terminates very quickly after only 2 to 10 steps, depending on the size of the margin. The number of principal components for each dataset is selected by plotting the explained variance (of the train data) per principal component in decreasing order on a logarithmic scale and applying the elbow method using the Kneedle algorithm from Satopaa et al. (2011). This results in a very

low-dimensional search space, ranging from 3 to 8 principal components for the seven unique datasets considered.

In order to prevent biasing our metric to the PGDL test set (Tasks 6 to 9) we did not perform any tuning or development of the complexity measure using these tasks, nor do we tune any hyperparameters per task. The choice of principal component selection algorithm was done after a careful analysis of Tasks 1 to 5 only, see additional details in supplementary material. In terms of computational expense, we find that calculating the entire constrained margin distribution only takes 1 to 2 minutes per model on a single Nvidia A30.

4.2 Margin Complexity Measures

In Table 1 we show the Kendall’s rank correlation obtained when ranking models according to constrained margin, standard input margins, and hidden margins.

It can be observed that standard input margins are not predictive of generalization for most tasks and, in fact, show a negative correlation for some. This unstable behaviour is supported by ongoing work surrounding adversarial robustness and generalization (Tsipras et al. 2019; Su et al. 2018; Ragunathan et al. 2019). Furthermore, we observe a very

large performance gap between constrained and standard input margins, and an increase from 0.24 to 0.66 average rank correlation is observed by constraining the margin search. This strongly supports our initial intuitions.

In the case of hidden margins, performance is more competitive, however, constrained margins still outperform hidden margins on 6 out of 8 tasks. One also observes that the selection of hidden layers can have a very large effect, and the discrepancy between the two hidden-layer selections is significant. Given that our constrained margin measurement is limited to the input space, there are several advantages: 1) no normalization is required, as all models share the same input space, and 2) the method is more robust when comparing models with varying topology, as no specific layers need to be selected.

4.3 Other Complexity Measures

To further assess the predictive power of constrained margins, we compare our method to the reported CMI scores of several other complexity measures. We compare against three solutions from the winning team (Natekar and Sharma 2020), as well as the best solutions from two more recent works (Chuang et al. 2021; Schiff et al. 2021), where that of Schiff et al. (2021) has the highest average test set performance we are aware of. We do not compare against pre-trained GANs (Zhang et al. 2022). The original naming of each method is kept. Of particular relevance are the *MM* and *AM* columns, which are hidden margins applied to Mixup and Augmented samples, as well as *kV*-Margin and *kV*-GN-Margin which are output and hidden margins with *k*-Variance normalization, respectively. The results of this comparison are shown in Table 2.

One observes that constrained margins achieve highly competitive scores, and in fact, outperform all other measures on 4 out of 8 tasks. It is also important to note that the *MM* and *AM* columns show that hidden margins can be improved in some cases if they are measured using the representations of Mixup or augmented training samples. That said, these methods still underperform on average in comparison to constrained input margins, which do not rely on any form of data augmentation.

5 A Closer Look

In this section we do a further analysis of constrained margins. In Section 5.1 we investigate how the performance of constrained margins changes when lower utility subspaces are considered, whereafter we discuss limitations of the method in Section 5.2.

5.1 High to Low Utility

We examine how high utility directions compare to those of lower utility when calculating constrained margins. This allows us to further test our approach, as one would expect that margins measured using the lower-ranked principal components should be less predictive of a model’s performance. We calculate the mean constrained margin using select subsets of 10 contiguous principal components in descending order of explained variance. For example, we calculate the

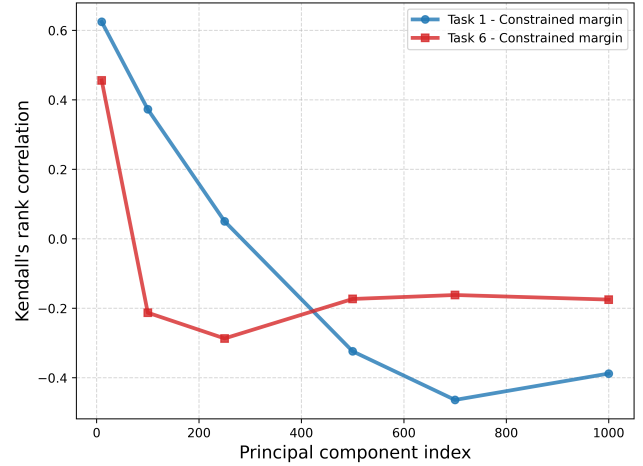


Figure 1: Comparison of predictive performance (Kendall’s rank correlation) of high to low utility directions using subspaces spanned by 10 principal components for Task 1 (blue) and 6 (red). The x-axis indicates the first component in each set of principal components.

constrained margins using components 1 to 10, then 100 to 109, etc. This allows us to calculate the distance to the decision boundary using 10 dimensional subspaces of decreasing utility. We, once again, make use of 5 000 training samples. For this analysis, we select two tasks where there is a large difference between the performance of constrained margins and standard input margins: Tasks 1 and 6. Figure 1 shows the resulting Kendall’s rank correlation for each subset of principal components indexed by the first component in each set (principal component index).

As expected, the first principal components lead to margins that are more predictive of generalization. We see a gradual decrease in predictive power when considering later principal components. For both tasks, we observe that they reach negative correlations when considering the later principal component subspaces. This supports the idea that utilizing the directions of highest utility is a necessary aspect of input margin measurements. After the point shown here (index 1 000), we find that the mean margin increases as DeepFool struggles to find samples on the decision boundary within the bound constraints. Due to this, it is difficult to draw any conclusions from an investigation of the lower-ranked principal components.

5.2 Limitations

It has been demonstrated that our proposed metric performs well and aligns with our initial intuition. However, there are also certain limitations that require explanation. Empirically we observe that, for tasks where constrained margins perform well, they do so across most hyperparameter variations, with the exception of depth. This is illustrated in Figure 2 (left), which shows the mean constrained margin versus test accuracy for Task 1. We observe that sets of networks with two and six convolutional layers, respectively, each exhibit a separate relationship between margin and test accuracy.

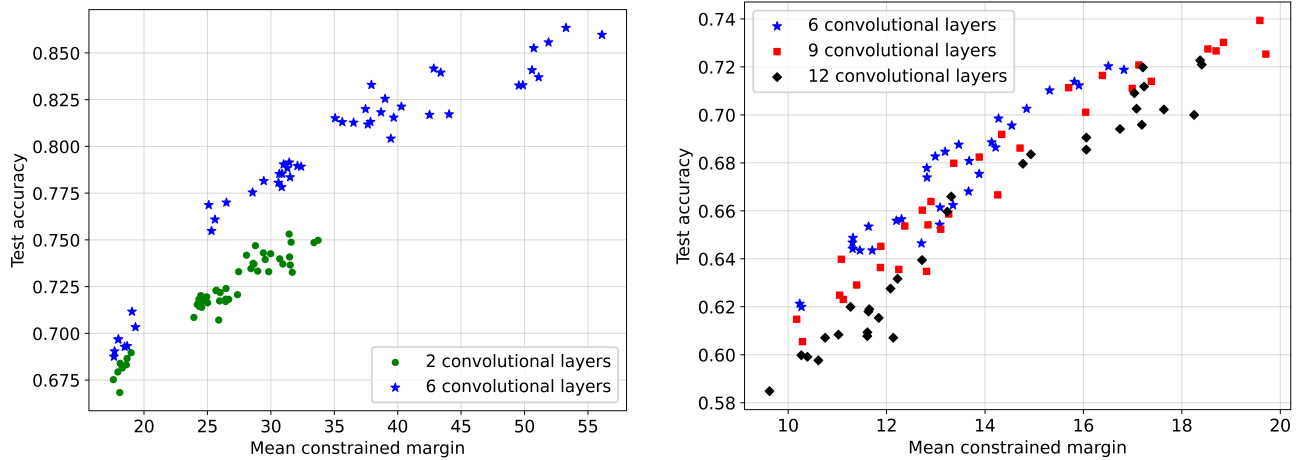


Figure 2: Mean constrained margin versus test accuracy for PGDL Task 1 (left) and 6 (right). Left: Models with 2 (green circle) and 6 (blue star) convolutional layers. Right: Models with 6 (blue star), 9 (red square), and 12 (black diamond) convolutional layers.

This discrepancy is not always as strongly present: for Task 6 all three depth configurations show a more similar relationship, as observed on the right of Figure 2, although the discrepancy is still present. The same trend holds several tasks (Tasks 1, 2, 4, 6, 9). It appears that shallower networks model the input space in a distinctly different fashion than their deeper counterparts.

For tasks such as 5 and 7, where constrained margins perform more poorly, there is no single hyperparameter that appears to be the culprit. We do note that the resulting scatter plots of margin versus test accuracy never show points in the lower right (large margin but low generalization) or upper left (small margin but high generalization) quadrants. It is therefore possible that a larger constrained margin is always beneficial to a model’s generalization, even though it is not always fully descriptive of its performance. Finally, it also possible to construct a hypothetical dataset such that the ideal decision boundary is not in the input space directions of highest variance, i.e. where high variance does not correspond to high utility. However, as evidenced by our results, this scenario does not present itself in natural image datasets. See the supplementary material for a fuller description of such a scenario.

6 Conclusion

We have shown that constraining input margins to high utility subspaces can significantly improve their predictive power i.t.o generalization. Specifically, we have used the principal components of the data as a proxy for identifying these subspaces, which can be considered a rough approximation of the underlying data manifold.

There are several implications to this work. First, we have shown that it is essential that the data manifold be taken into account to relate input margins to generalization. This is an important consideration for probing generalization from the decision boundary perspective. Secondly, several au-

thors have developed techniques to maximize margins during training (Elsayed et al. 2018; Xu et al. 2023); however, these have not resulted in improved generalization. We believe that constrained margin has the potential of being a powerful regularizer, in line with how other complexity measures have been used in the past.

In terms of future work, we know that constraining the search to a warped subspace and using Euclidean distance to measure closeness is equivalent to defining a new distance metric on the original space. We are therefore, in effect, seeking a relevant distance metric to measure the closeness of the decision boundary. Understanding the requirements for such a metric remains an open question. Unfortunately, current approximations and methods for finding points on the decision boundary are largely confined to L_p metrics. The positive results achieved with the current PCA-and-Euclidean-based approach provide strong motivation that this is a useful avenue to pursue.

In conclusion, we propose constraining input margins to make them more predictive of generalization in DNNs. It has been demonstrated that this greatly increases the predictive power of input margins, and also outperforms hidden margins and several other contemporary methods on the PGDL tasks. This method has the benefits of requiring no per-layer normalization, no arbitrary selection of hidden layers, and does not rely on any form of surrogate test set (e.g. data augmentation or synthetic samples).

References

- Arora, S.; Ge, R.; Neyshabur, B.; and Zhang, Y. 2018. Stronger generalization bounds for deep nets via a compression approach. In *International Conference on Machine Learning (ICML)*, 254–263. PMLR.
- Bartlett, P.; Foster, D.; and Telgarsky, M. 2017. Spectrally-normalized margin bounds for neural networks. *Advances in Neural Information Processing Systems*, 30.

- Boser, B. E.; Guyon, I. M.; and Vapnik, V. N. 1992. A training algorithm for optimal margin classifiers. In *Proceedings of the fifth annual workshop on Computational Learning Theory*, 144–152.
- Carlini, N.; and Wagner, D. 2017. Towards Evaluating the Robustness of Neural Networks. In *IEEE Symposium on Security and Privacy*, 39–57. IEEE.
- Chuang, C.-Y.; Mroueh, Y.; Greenewald, K.; Torralba, A.; and Jegelka, S. 2021. Measuring generalization with optimal transport. *Advances in Neural Information Processing Systems*, 34: 8294–8306.
- Davies, D. L.; and Bouldin, D. W. 1979. A cluster separation measure. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, (2): 224–227.
- Elsayed, G.; Krishnan, D.; Mobahi, H.; Regan, K.; and Bengio, S. 2018. Large Margin Deep Networks for Classification. *Advances in Neural Information Processing Systems*, 31.
- Hinton, G. E.; Dayan, P.; and Revow, M. 1997. Modeling the Manifolds of Images of Handwritten Digits. *IEEE Transactions on Neural Networks*, 8(1): 65–74.
- Huang, R.; Xu, B.; Schuurmans, D.; and Szepesvári, C. 2015. Learning with a Strong Adversary. *arXiv preprint arXiv:1511.03034*.
- Jiang, Y.; Foret, P.; Yak, S.; Roy, D. M.; Mobahi, H.; Dziugaite, G. K.; Bengio, S.; Gunasekar, S.; Guyon, I.; and Neyshabur, B. 2020. Neurips 2020 competition: Predicting generalization in deep learning. *arXiv preprint arXiv:2012.07976*.
- Jiang, Y.; Krishnan, D.; Mobahi, H.; and Bengio, S. 2018. Predicting the Generalization Gap in Deep Networks with Margin Distributions. In *International Conference on Learning Representations*.
- Jiang, Y.; Neyshabur, B.; Mobahi, H.; Krishnan, D.; and Bengio, S. 2019. Fantastic Generalization Measures and Where to Find Them. In *International Conference on Learning Representations*.
- Kashyap, D.; Subramanyam, N.; et al. 2021. Robustness to augmentations as a generalization metric. *arXiv preprint arXiv:2101.06459*.
- Kawaguchi, K.; Kaelbling, L. P.; and Bengio, Y. 2022. *Generalization in Deep Learning*. Mathematical Aspects of Deep Learning. Cambridge University Press.
- Kendall, M. G. 1938. A New Measure of Rank Correlation. *Biometrika*, 30(1-2): 81–93.
- Koltchinskii, V.; and Panchenko, D. 2002. Empirical margin distributions and bounding the generalization error of combined classifiers. *The Annals of Statistics*, 30(1): 1–50.
- Lotfi, S.; Finzi, M. A.; Kapoor, S.; Potapczynski, A.; Goldblum, M.; and Wilson, A. G. 2022. PAC-Bayes Compression Bounds So Tight That They Can Explain Generalization. In *Advances in Neural Information Processing Systems*.
- Madry, A.; Makelov, A.; Schmidt, L.; Tsipras, D.; and Vladu, A. 2018. Towards Deep Learning Models Resistant to Adversarial Attacks. In *International Conference on Learning Representations*.
- Moosavi-Dezfooli, S.-M.; Fawzi, A.; and Frossard, P. 2016. DeepFool: a simple and accurate method to fool deep neural networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2574–2582.
- Natekar, P.; and Sharma, M. 2020. Representation based complexity measures for predicting generalization in deep learning. *arXiv preprint arXiv:2012.02775*.
- Raghunathan, A.; Xie, S. M.; Yang, F.; Duchi, J.; and Liang, P. 2019. Adversarial Training Can Hurt Generalization. In *ICML Workshop on Identifying and Understanding Deep Learning Phenomena*.
- Satopaa, V.; Albrecht, J.; Irwin, D.; and Raghavan, B. 2011. Finding a “kneedle” in a haystack: Detecting knee points in system behavior. In *2011 31st international Conference on Distributed Computing systems workshops*, 166–171. IEEE.
- Schiff, Y.; Quanz, B.; Das, P.; and Chen, P.-Y. 2021. Predicting Deep Neural Network Generalization with Perturbation Response Curves. *Advances in Neural Information Processing Systems*, 34: 21176–21188.
- Sokolić, J.; Giryès, R.; Sapiro, G.; and Rodrigues, M. R. 2017. Robust large margin deep neural networks. *IEEE Transactions on Signal Processing*, 65(16): 4265–4280.
- Solomon, J.; Greenewald, K.; and Nagaraja, H. 2022. k-variance: A clustered notion of variance. *SIAM Journal on Mathematics of Data Science*, 4(3): 957–978.
- Stutz, D.; Hein, M.; and Schiele, B. 2019. Disentangling Adversarial Robustness and Generalization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 6976–6987.
- Su, D.; Zhang, H.; Chen, H.; Yi, J.; Chen, P.-Y.; and Gao, Y. 2018. Is Robustness the Cost of Accuracy?—A Comprehensive Study on the Robustness of 18 Deep Image Classification Models. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 631–648.
- Tsipras, D.; Santurkar, S.; Engstrom, L.; Turner, A.; and Madry, A. 2019. Robustness May Be at Odds with Accuracy. In *International Conference on Learning Representations*.
- Vapnik, V. N. 1999. An overview of statistical learning theory. *IEEE Transactions on Neural Networks*, 10(5): 988–999.
- Weinberger, K. Q.; and Saul, L. K. 2009. Distance Metric Learning for Large Margin Nearest Neighbor Classification. *Journal of Machine Learning Research*, 10(9): 207–244.
- Xu, Y.; Sun, Y.; Goldblum, M.; Goldstein, T.; and Huang, F. 2023. Exploring and Exploiting Decision Boundary Dynamics for Adversarial Robustness. In *International Conference on Learning Representations*.
- Yang, Y.; Khanna, R.; Yu, Y.; Gholami, A.; Keutzer, K.; Gonzalez, J. E.; Ramchandran, K.; and Mahoney, M. W. 2020. Boundary thickness and robustness in learning models. *Advances in Neural Information Processing Systems*, 33: 6223–6234.
- Yousefzadeh, R.; and O’Leary, D. P. 2020. Deep learning interpretation: Flip points and homotopy methods. In Lu, J.; and Ward, R., eds., *Proceedings of The First Mathematical and Scientific Machine Learning Conference*, volume 107 of *Proceedings of Machine Learning Research*, 1–26. PMLR.

Zhang, C.; Bengio, S.; Hardt, M.; Recht, B.; and Vinyals, O. 2021. Understanding deep learning (still) requires rethinking generalization. *Communications of the ACM*, 64(3): 107–115.

Zhang, H.; Cisse, M.; Dauphin, Y. N.; and Lopez-Paz, D. 2018. mixup: Beyond Empirical Risk Minimization. In *International Conference on Learning Representations*.

Zhang, Y.; Gupta, A.; Saunshi, N.; and Arora, S. 2022. On Predicting Generalization using GANs. In *International Conference on Learning Representations*.