

Improving the Robustness of Deep Neural Networks via Stability Training

Stephan Zheng
Google, Caltech

stzheng@caltech.edu

Yang Song
Google

yangsong@google.com

Thomas Leung
Google

leungt@google.com

Ian Goodfellow
Google

goodfellow@google.com

Abstract

In this paper we address the issue of output instability of deep neural networks: small perturbations in the visual input can significantly distort the feature embeddings and output of a neural network. Such instability affects many deep architectures with state-of-the-art performance on a wide range of computer vision tasks. We present a general stability training method to stabilize deep networks against small input distortions that result from various types of common image processing, such as compression, rescaling, and cropping. We validate our method by stabilizing the state-of-the-art Inception architecture [11] against these types of distortions. In addition, we demonstrate that our stabilized model gives robust state-of-the-art performance on large-scale near-duplicate detection, similar-image ranking, and classification on noisy datasets.

1. Introduction

Deep neural networks learn feature embeddings of the input data that enable state-of-the-art performance in a wide range of computer vision tasks, such as visual recognition [3, 11] and similar-image ranking [13]. Due to this success, neural networks are now routinely applied to vision tasks on large-scale *un-curated* visual datasets that, for instance, can be obtained from the Internet. Such un-curated visual datasets often contain small distortions that are undetectable to the human eye, due to the large diversity in formats, compression, and manual post-processing that are commonly applied to visual data in the wild. **These lossy image processes do not change the correct ground truth labels and semantic content of the visual data, but can significantly confuse feature extractors, including deep neural networks.** Namely, when presented with a pair of indistinguishable images, state-of-the-art feature extractors can produce two significantly different outputs.

In fact, current feature embeddings and class labels are not robust to a large class of small perturbations. Recently, it has become known that intentionally engineered imperceptible perturbations of the input can change the class label

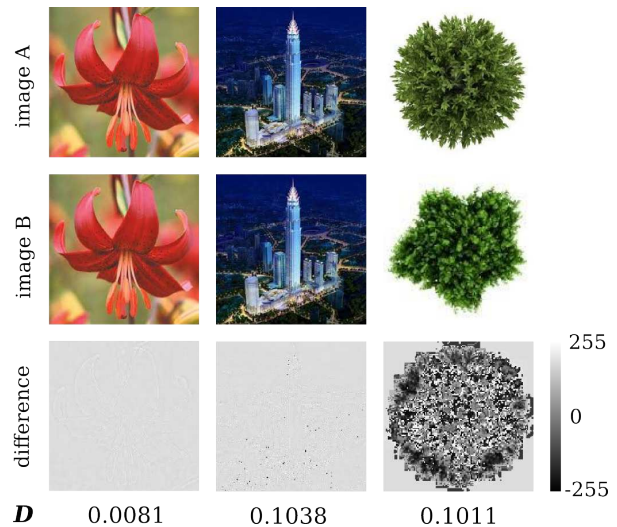


Figure 1: Near-duplicate images can confuse state-of-the-art neural networks due to feature embedding instability. Left and middle columns: near-duplicates with small (left) and large (middle) feature distance. Image A is the original, image B is a JPEG version at quality factor 50. Right column: a pair of dissimilar images. In each column we display the pixel-wise difference of image A and image B, and the feature distance D [13]. Because the feature distances of the middle near-duplicate pair and the dissimilar image pair are comparable, near-duplicate detection using a threshold on the feature distance will confuse the two pairs.

output by the model [1, 12]. A scientific contribution of this paper is the demonstration that these imperceptible perturbations can also occur without being contrived and widely occur due to compression, resizing, and cropping corruptions in visual input.

As such, output instability poses a significant challenge for the large-scale application of neural networks because high performance at large scale requires robust performance on noisy visual inputs. Feature instability complicates tasks such as near-duplicate detection, which is essential for large-scale image retrieval and other applications. In near-duplicate detection, the goal is to detect whether two given images are visually similar or not. When neural networks



Figure 2: Visually similar video frames can confuse state-of-the-art classifiers: two neighboring frames are visually indistinguishable, but can lead to very different class predictions. The class score for ‘fox’ is significantly different for the left frame (27%) and right frame (63%), which causes only the fox in the right image to be correctly recognized, using any reasonable confidence threshold (e.g. > 50%).

are applied to this task, there are many failure cases due to output instability. For instance, Figure 1 shows a case where a state-of-the-art deep network cannot distinguish a pair of near-duplicates [13] and a pair of dissimilar images.

Analogously, class label instability introduces many failure cases in large-scale classification and annotation. For example, unstable classifiers can classify neighboring video-frames inconsistently, as shown in Figure 2. In this setting, output instability can cause large changes in label scores of a state-of-the-art convolutional neural network on consecutive video-frames that are indistinguishable.

The goal of this paper is to propose a general approach to stabilize machine learning models, in particular deep neural networks, and make them more robust to visual perturbations. To this end, we introduce a fast and effective *stability training* technique that makes the output of neural networks significantly more robust, while maintaining or improving state-of-the-art performance on the original task. To do so, our method operates through two mechanisms: 1) introducing an additional stability training objective and 2) training on a large class of distorted copies of the input. The goal of this approach is to force the prediction function of the model to be more constant around the input data, while preventing underfitting on the original learning objective. In summary, our contributions are as follows:

- We propose stability training as a general technique that improves model output stability while maintaining or improving the original performance. Our method is fast in practice and can be used at a minimal additional computational cost.
- We validate our method by stabilizing state-of-the-art classification and ranking networks based on the Inception architecture [11, 13]. We evaluate on three tasks: near-duplicate image detection, similar-image ranking, and image classification.
- We show the impact of stability training by visualizing

what perturbations the model has become robust to.

- Finally, we show that stabilized networks offer robust performance and significantly outperform unstabilized models on noisy and corrupted data.

2. Related work

Adversarial examples. Recently, several machine learning algorithms were found to have extreme instability against *contrived* input perturbations [12] called adversarial examples. An open question remained as to whether such small perturbations that change the class label could occur without intentional human intervention. In this work, we document that they do in fact occur. Previous work has shown that training a classifier to resist adversarial perturbation can improve its performance on both the original data and on perturbed data [1, 6]. We extend this approach by training our feature embeddings to resist the naturally occurring perturbations that are far more common in practice.

Furthermore, our work differs drastically from [7], which is about how a model responds to intentionally contrived inputs that don’t resemble the original data at all. In contrast, in this paper we consider the stability to practically widely occurring perturbations.

Data augmentation. A natural strategy to improve label stability is to augment the training data with *hard positives*, which are examples that the prediction model does not classify correctly with high confidence, but that are visually similar to easy positives. Finding such hard positives in video data for data augmentation has been used in [5, 4, 8] and has been found to improve predictive performance and consistency. As such, data augmentation with hard positives can confer output stability on the classes of perturbations that the hard positives represent. However, our work differs from data augmentation in two ways. Firstly, we take a general approach by proposing a method that intends to make model performance more robust to various types of natural perturbations. Secondly, our proposed method does not use the extra generated samples as training examples for the original prediction task, but only for the stability objective.

3. Stability training

We now present our stability training approach, and how it can be applied to learn robust feature embeddings and class label predictions.

3.1. Stability objective

Our goal is to stabilize the output $f(x) \in \mathbb{R}^m$ of a neural network \mathcal{N} against small natural perturbations to a natural image $x \in [0, 1]^{w \times h}$ of size $w \times h$, where we normalize all pixel values. Intuitively, this means that we want to formulate a training objective that flattens f in a small neighborhood of any natural image x : if a perturbed copy x' is close

to x , we want $f(x)$ to be close to $f(x')$, that is

$$\forall x' : d(x, x') \text{ small} \Leftrightarrow D(f(x), f(x')) \text{ small}. \quad (1)$$

Here d is the distance on $[0, 1]^{w \times h}$ and D is an appropriate distance measure in feature space.

Given a training objective L_0 for the original task (e.g. classification, ranking), a reference input x and a perturbed copy x' , we can implement the stability objective (1) as:

$$L(x, x'; \theta) = L_0(x; \theta) + \alpha L_{\text{stability}}(x, x'; \theta), \quad (2)$$

$$L_{\text{stability}}(x, x'; \theta) = D(f(x), f(x')), \quad (3)$$

where α controls the strength of the stability term and θ denotes the weights of the model \mathcal{N} . The stability objective $L_{\text{stability}}$ forces the output $f(x)$ of the model to be similar between the original x and the distorted copy x' . Note that *our approach differs from data augmentation*: we do not evaluate the original loss L on the distorted inputs x' . This is required to achieve both output stability and performance on the original task, as we explain in 3.2.

Given a training dataset \mathcal{D} , stability training now proceeds by finding the optimal weights θ^* for the training objective (2), that is, we solve

$$\theta^* = \underset{\theta}{\operatorname{argmin}} \sum_{x_i \in \mathcal{D}, d(x_i, x'_i) < \epsilon} L(x_i, x'_i; \theta). \quad (4)$$

To fully specify the optimization problem, we firstly need a mechanism to generate, for each training step, for each training sample x_i , a random perturbed copy x'_i . Secondly, we need to define the distance D , which is task-specific.

3.2. Sampling perturbed images x'

Sampling using Gaussian noise. During training, at every training step we need to generate perturbed versions x' of a clean image x to evaluate the stability objective (3).

A natural approach would be to augment the training data with examples with explicitly chosen classes of perturbation that the model should be robust against. However, it is hard to obtain general robustness in this way, as there are many classes of perturbations that cause output instability, and model robustness to one class of perturbations does not confer robustness to other classes of perturbations.

Therefore, we take a general approach and use a sampling mechanism that adds pixel-wise uncorrelated Gaussian noise ϵ to the visual input x . If k indexes the raw pixels, a new sample is given by:

$$x'_k = x_k + \epsilon_k, \quad \epsilon_k \sim \mathcal{N}(0, \sigma_k^2), \quad \sigma_k > 0, \quad (5)$$

where σ_k^2 is the variance of the Gaussian noise at pixel k . In this work, we use uniform sampling $\sigma_k = \sigma$ to produce unbiased samples of the neighborhood of x , using the variance σ^2 as a hyper-parameter to be optimized.



Figure 3: Examples of reference and distorted training images used for stability training. Left: an original image x . Right: a copy x' perturbed with pixel-wise uncorrelated Gaussian noise with $\sigma = 0.06$, in normalized pixel values. During stability training, we use dynamically sampled copies x' together with the stability loss (3) to flatten the prediction function f around the original image x .

Gaussian noise strength σ	0.0	0.1	0.2
Triplet ranking score @ top-30	7,312	6,300	5,065

Table 1: Underfitting by data augmentation with Gaussian noise on an image ranking task (higher score is better), see section 5.2 for details. The entry with $\sigma = 0.0$ is the model without data augmentation.

Preventing underfitting. Augmenting the training data by adding uncorrelated Gaussian noise can potentially simulate many types of perturbations. Training on these extra samples could in principle lead to output robustness to many classes of perturbations. However, we found that training on a dataset augmented by Gaussian perturbation leads to underfitting, as shown in Table 1. To prevent such underfitting, we do *not* evaluate the original loss L_0 on the perturbed images x' in the full training objective (2), but only evaluate the stability loss (3) on both x and x' . This approach differs from data augmentation, where one would evaluate L_0 on the extra training samples as well. It enables achieving both output stability and maintaining high performance on the original task, as we validate empirically.

3.3. Stability for feature embeddings

We now show how stability training can be used to obtain stable feature embeddings. In this work, we aim to learn feature embeddings for robust similar-image detection. To this end, we apply stability training in a ranking setting. The objective for similar-image ranking is to learn a feature representation $f(x)$ that detects visual image similarity [13]. This learning problem is modeled by considering a *ranking triplet* of images (q, p, n) : a *query* image q , a *positive* image p that is visually similar to q , and a *negative* image n that is less similar to q than p is.

The objective is to learn a feature representation f that

respects the triplet ranking relationship in feature space, that is,

$$D(f(q), f(p)) + g < D(f(q), f(n)), \quad g > 0, \quad (6)$$

where g is a margin and D is the distance. We can learn a model for this objective by using a hinge loss:

$$L_0(q, p, n) = \max(0, g + D(f(q), f(p)) - D(f(q), f(n))). \quad (7)$$

In this setting, a natural choice for the similarity metric D is the L_2 -distance. The stability loss is,

$$L_{\text{stability}}(x, x') = \|f(x) - f(x')\|_2. \quad (8)$$

To make the feature representation f stable using our approach, we sample triplet images (q', p', n') close to the reference (q, p, n) , by applying (5) to each image in the triplet.

3.4. Stability for classification

We also apply stability training in the classification setting to learn stable prediction labels for visual recognition. For this task, we model the likelihood $P(\mathbf{y}|x; \theta)$ for a labeled dataset $\{(x_i, \hat{\mathbf{y}}_i)\}_{i \in \mathcal{I}}$, where $\hat{\mathbf{y}}$ represents a vector of ground truth binary class labels and i indexes the dataset. The training objective is then to minimize the standard cross-entropy loss

$$L_0(x; \theta) = - \sum_j \hat{y}_j \log P(y_j | x; \theta), \quad (9)$$

where the index j runs over classes. To apply stability training, we use the KL-divergence as the distance function D :

$$L_{\text{stability}}(x, x'; \theta) = - \sum_j P(y_j | x; \theta) \log P(y_j | x'; \theta), \quad (10)$$

which measures the correspondence between the likelihood on the natural and perturbed inputs.

4. Implementation

4.1. Network

Base network. In our experiments, we use a state-of-the-art convolutional neural network architecture, the Inception network [11] as our base architecture. Inception is formed by a deep stack of composite layers, where each composite layer output is a concatenation of outputs of convolutional and pooling layers. This network is used for the classification task and as a main component in the triplet ranking network.

Triplet ranking network. Triplet ranking loss (7) is used to train feature embeddings for image similarity and for

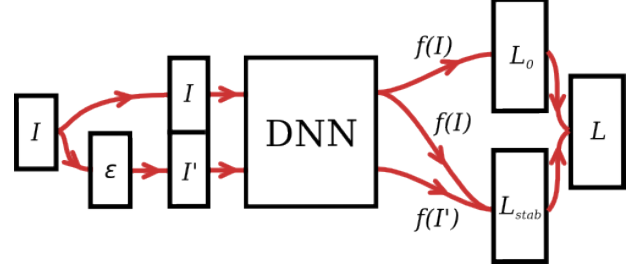


Figure 4: The architecture used to apply stability training to any given deep neural network. The arrows display the flow of information during the forward pass. For each input image I , a copy I' is perturbed with pixel-wise independent Gaussian noise ϵ . Both the original and perturbed version are then processed by the neural network. The task objective L_0 is only evaluated on the output $f(I)$ of the original image, while the stability loss $L_{\text{stability}}$ uses the outputs of both versions. The gradients from both L_0 and $L_{\text{stability}}$ are then combined into the final loss L and propagated back through the network. For triplet ranking training, three images are processed to compute the triplet ranking objective.

near duplicate image detection, similar to [13]. This network architecture uses an Inception module (while in [13], a network like [3] is used) to process every input image x at full resolution and uses 2 additional low-resolution towers. The outputs of these towers map into a 64-dimensional L_2 -normalized embedding feature $f(x)$. These features are used for the ranking task: for each triplet of images (q, p, n) , we use the features $(f(q), f(p), f(n))$ to compute the ranking loss and train the entire architecture.

Stability training. It is straightforward to implement stability training for any given neural network by adding a Gaussian perturbation sampler to generate perturbed copies of the input image x and an additional stability objective layer. This setup is depicted in Figure 4.

4.2. Distortion types

To demonstrate the robustness of our models after stability training is deployed, we evaluate the ranking, near-duplicate detection and classification performance of our stabilized models on both the original and transformed copies of the evaluation datasets. To generate the transformed copies, we apply visual perturbations that widely occur in real-world visual data and that are a result of lossy image processes.

JPEG compression. JPEG compression is a commonly used lossy compression method that introduces small artifacts in the image. The extent and intensity of these artifacts can be controlled by specifying a quality level q . In this work, we refer to this as JPEG- q .

Thumbnail resizing. Thumbnails are smaller versions of a reference image and obtained by downscaling the original image. Because convolutional neural networks use a

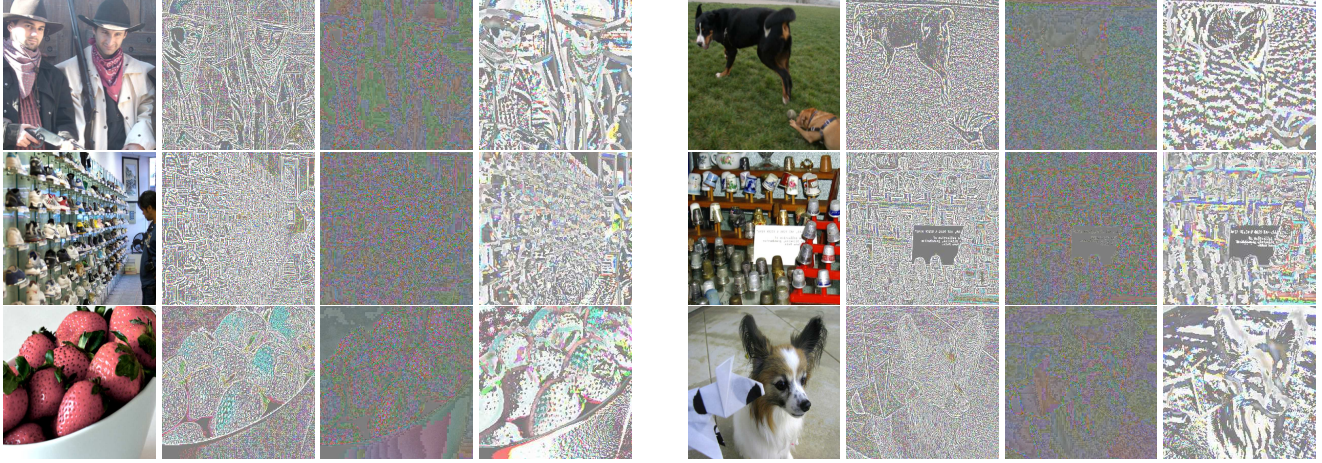


Figure 5: Examples of natural distortions that are introduced by common types of image processing. From left to right: original image (column 1 and 5), pixel-wise differences from the original after different forms of transformation: thumbnail downscaling to 225×225 (column 2 and 6), JPEG compression at quality level 50% (column 3 and 7) and random cropping with offset 10 (column 4 and 8). For clarity, the JPEG distortions have been up-scaled by $5\times$. Random cropping and thumbnail resizing introduce distortions that are structured and resemble the edge structure of the original image. In contrast, JPEG compression introduces more unstructured noise.

fixed input size, both the original image and its thumbnail have to be rescaled to fit the input window. Downscaling and rescaling introduces small differences between the original and thumbnail versions of the network input. In this work we refer to this process as THUMB- A , where we down-scale to a thumbnail with A pixels, preserving the aspect ratio.

Random cropping. We also evaluated the performance on perturbations coming from random crops of the original image. This means that we take large crops with window size $w' \times h'$ of the original image of size $w \times h$, using an offset $o > 0$ to define $w' = w - o$, $h' = h - o$. The crops are centered at random positions, with the constraint that the cropping window does not exceed the image boundaries. Due to the fixed network input size, resizing the cropped image and the original image to the input window introduces small perturbations in the visual input, analogous to thumbnail noise. We refer to this process as CROP- o , for crops with a window defined by offset o .

4.3. Optimization

To perform stability training, we solved the optimization problem (2) by training the network using mini-batch stochastic gradient descent with momentum, dropout [10], RMSprop and batch normalization [2]. To tune the hyper-parameters, we used a grid search, where the search ranges are displayed in Table 2.

As stability training requires a distorted version of the original training example, it effectively doubles the training batch-size during the forward-pass, which introduces a significant extra computational cost. To avoid this over-

Hyper-parameter	Start range	End range
Noise standard deviation σ	0.01	0.4
Regularization coefficient α	0.001	1.0
Learning rate λ	0.001	0.1

Table 2: Hyper-parameter search range for the stability training experiments.

head, in our experiments we first trained the network on the original objective $L_0(x; \theta)$ only and started stability training with $L(x, x'; \theta)$ only in the fine-tuning phase. Additionally, when applying stability training, we only fine-tuned the final fully-connected layers of the network. Experiments indicate that this approach leads to the same model performance as applying stability training right from the beginning and training the whole network during stability training.

5. Experiments

Here we present experimental results to validate our stability training method and characterize stabilized models.

- Firstly, we evaluate stabilized features on near-duplicate detection and similar-image ranking tasks.
- Secondly, we validate our approach of stabilizing classifiers on the ImageNet classification task.

We use training data as in [13] to train the feature embeddings for near-duplicate detection and similar-image ranking. For the classification task, training data from ImageNet are used.

5.1. Near-duplicate detection

Detection criterion. We used our stabilized ranking feature to perform near-duplicate detection. To do so, we define the detection criterion as follows: given an image pair (a, b) , we say that

$$a, b \text{ are near-duplicates} \iff \|f(a) - f(b)\|_2 < T, \quad (11)$$

where T is the near-duplicate detection threshold.

Near-duplicate evaluation dataset. For our experiments, we generated an image-pair dataset with two parts: one set of pairs of near-duplicate images (true positives) and a set of dissimilar images (true negatives).

We constructed the near-duplicate dataset by collecting 650,000 images from randomly chosen queries on Google Image Search. In this way, we obtained a representative sample of un-curated images. We then combined every image with a copy perturbed with the distortion(s) from section 4.2 to construct near-duplicate pairs. For the set of dissimilar images, we collected 900,000 random image pairs from the top 30 Image Search results for 900,000 random search queries, where the images in each pair come from the *same* search query.

5.1.1 Experimental results

Precision-recall performance. To analyze the detection performance of the stabilized features, we report the near-duplicate precision-recall values by varying the detection threshold in (11). Our results are summarized in Figure 6. The stabilized deep ranking features outperform the baseline features for all three types of distortions, for all levels of fixed recall or fixed precision. Although the baseline features already offer very high performance in both precision and recall on the near-duplicate detection task, the stabilized features significantly improve precision across the board. For instance, recall increases by 1.0% at 99.5% precision for thumbnail near-duplicates, and increases by 3.0% at 98% precision for JPEG near-duplicates. This improved performance is due to the improved robustness of the stabilized features, which enables them to correctly detect near-duplicate pairs that were confused with dissimilar image pairs by the baseline features, as illustrated in Figure 1.

Feature distance distribution. To analyze the robustness of the stabilized features, we show the distribution of the feature distance $D(f(x), f(x'))$ for the near-duplicate evaluation dataset in Figure 7, for both the baseline and stabilized deep ranking feature. Stability training significantly increases the feature robustness, as the distribution of feature distances becomes more concentrated towards 0. For instance, for the original feature 76% of near-duplicate image pairs has feature distance smaller than 0.1, whereas this is 86% for the stabilized feature, i.e. the stabilized feature is significantly more similar for near-duplicate images.

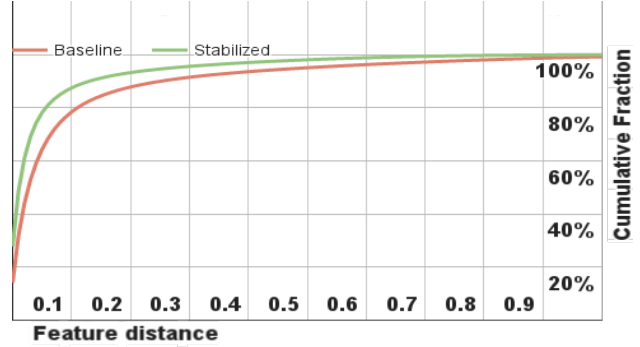


Figure 7: Cumulative distribution of the deep ranking feature distance $D(f(x_i), f(x'_i)) = \|f(x_i) - f(x'_i)\|_2$ for near-duplicate pairs (x_i, x'_i) . Red: baseline features, 76% of distribution < 0.1 . Green: stabilized features using stability training with $\alpha = 0.1, \sigma = 0.2$, 86% of distribution < 0.1 . The feature distances are computed over a dataset of 650,000 near-duplicate image pairs (reference image and a JPEG-50 version). Applying stability training makes the distribution of $D(f(x), f(x'))$ more concentrated towards 0 and hence makes the feature f significantly more stable.

Stabilized feature distance. We also present our qualitative results to visualize the improvements of the stabilized features over the original features. In Figure 8 we show pairs of images and their JPEG versions that were confusing for the un-stabilized features, i.e. that lay far apart in feature space, but whose stabilized features are significantly more close. This means that they are correctly detected as near-duplicates for much more aggressive, that is, lower detection thresholds by the stabilized feature, whereas the original feature easily confuses these as dissimilar images. Consistent with the intuition that Gaussian noise applies a wide range of types of perturbations, we see improved performance for a wide range of perturbation types. Importantly, this includes even localized, structured perturbations that do not resemble a typical Gaussian noise sample.

5.2. Similar image ranking

The stabilized deep ranking features (see section 3.3) are evaluated on the similar image ranking task. Hand-labeled triplets from [13]¹ are used as evaluation data. There are 14,000 such triplets. The ranking score-at-top- K ($K = 30$) is used as evaluation metric. The ranking score-at-top- K is defined as

$$\begin{aligned} \text{ranking score @top-}K &= \\ &\# \text{ correctly ranked triplets} - \# \text{ incorrectly ranked triplets,} \end{aligned} \quad (12)$$

where only triplets whose positive or negative image occurs among the closest K results from the query image are con-

¹<https://sites.google.com/site/imagesimilaritydata/>.

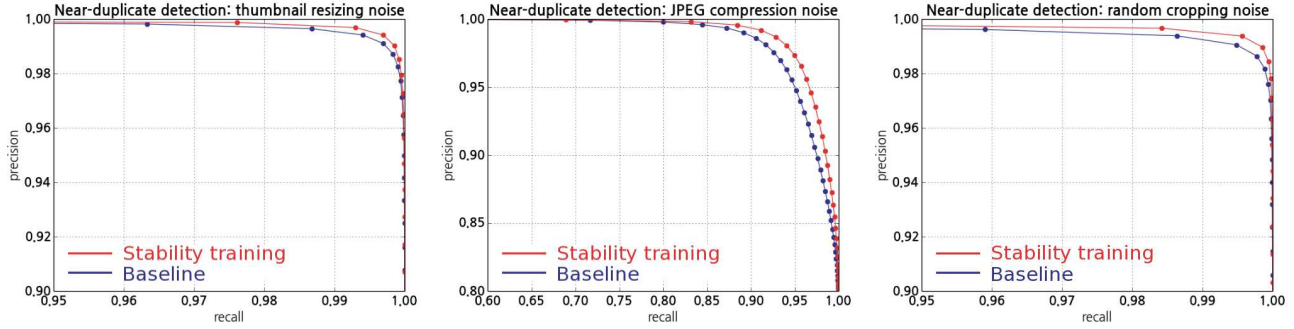


Figure 6: Precision-recall performance for near-duplicate detection using feature distance thresholding on deep ranking features. We compare Inception-based deep ranking features (blue), and the same features with stability training applied (red). Every graph shows the performance using near-duplicates generated through different distortions. Left: THUMB-50k. Middle: JPEG-50. Right: CROP-10. Across the three near-duplicate tasks, the stabilized model significantly improves the near-duplicate detection precision over the baseline model.

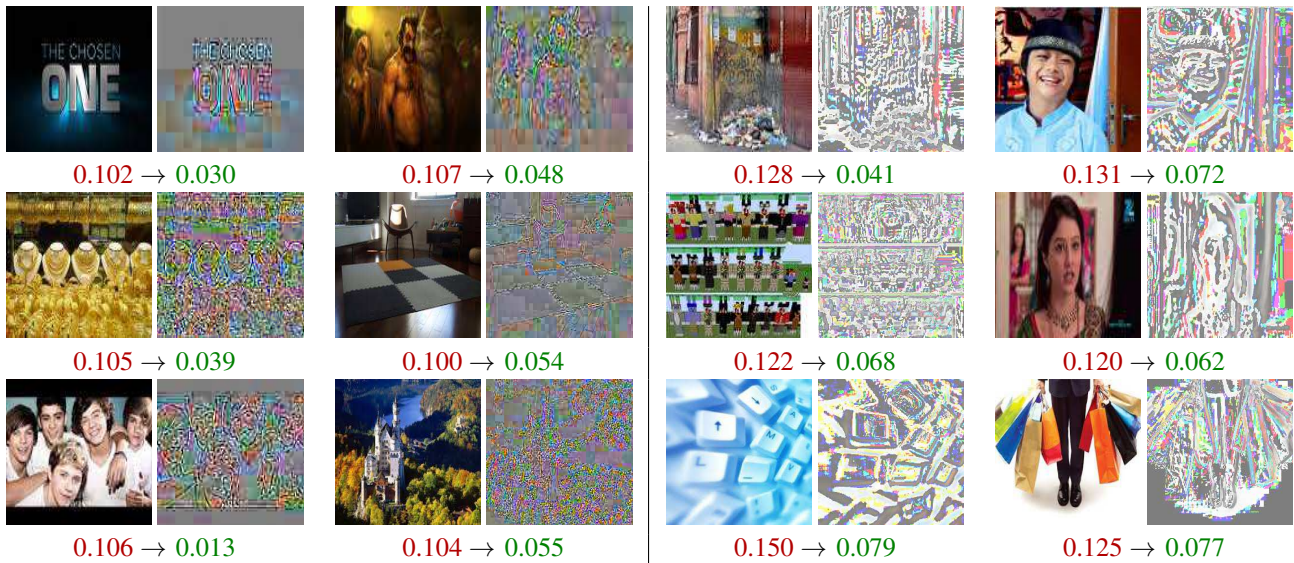


Figure 8: Examples of near-duplicate image pairs that are robustly recognized as near-duplicates by stabilized features (small feature distance), but easily confuse un-stabilized features (large feature distance). Left group: using JPEG-50 compression corruptions. Right group: random cropping CROP-10 corruptions. For each image pair, we display the reference image x , the difference with its corrupted copy $x - x'$, and the distance in feature space $D(f(x), f(x'))$ for the un-stabilized (red) and stabilized features (green).

sidered. This metric measures the ranking performance on the K most relevant results of the query image. We use this evaluation metric because it reflects better the performance of similarity models in practical image retrieval systems as users pay most of their attentions to the results on the first few pages.

5.2.1 Experimental results.

Our results for triplet ranking are displayed in Table 3. The results show that applying stability training improves the ranking score on both the original and transformed versions of the evaluation dataset. The ranking performance of the baseline model degrades on all distorted versions of

the original dataset, showing that it is not robust to the input distortions. In contrast, the stabilized network achieves ranking scores that are higher than the ranking score of the baseline model on the *original* dataset.

5.3. Image classification

In the classification setting, we validated stability training on the ImageNet classification task [9], using the Inception network [11]. We used the full classification dataset, which covers 1,000 classes and contains 1.2 million images, where 50,000 are used for validation. We evaluated the classification precision on both the original and a JPEG-50 version of the validation set. Our benchmark results are in Table 4.

Distortion	Deep ranking	Deep ranking + ST
Original	7,312	7,368
JPEG-50	7,286	7,360
THUMB-30k	7,160	7,172
CROP-10	7,298	7,322

Table 3: Ranking score @top-30 for the deep ranking network with and without stability training (higher is better) on distorted image data. Stability training increases ranking performance over the baseline on all versions of the evaluation dataset. We do not report precision scores, as in [13], as the ranking score @top-30 agrees more with human perception of practical similar image retrieval.

Precision @top-5	Original	JPEG-50	JPEG-10
Szegedy et al [11]	93.3%		
Inception	93.9%	92.4%	83.0%
Stability training	93.6%	92.7%	88.3%
Precision @top-1			
Inception	77.8%	75.1%	61.1%
Stability training	77.9%	75.7%	67.9%

Table 4: Classification evaluation performance of Inception with stability training, evaluated on the original and JPEG versions of ImageNet. Both networks give similar state-of-the-art performance on the original evaluation dataset (note that the performance difference on the original dataset is within the statistical error of 0.3% [9]). However, the stabilized network is significantly more robust and outperforms the baseline on the distorted data.

Applying stability training to the Inception network makes the class predictions of the network more robust to input distortions. On the original dataset, both the baseline and stabilized network achieve state-of-the-art performance. However, the stabilized model achieves higher precision on the distorted evaluation datasets, as the performance degrades more significantly for the baseline model than for the stabilized model. For high distortion levels, this gap grows to 5% to 6% in top-1 and top-5 precision.

Robust classification on noisy data. We also evaluated the effectiveness of stability training on the classification performance of Inception on the ImageNet evaluation dataset with increasing JPEG corruption. In this experiment, we collected the precision @top-1 scores at convergence for a range of the training hyper-parameters: the regularization coefficient α and noise standard deviation σ . A summary of these results is displayed in Figure 9.

At the highest JPEG quality level, the performance of the baseline and stabilized models are comparable, as the visual distortions are small. However, as the JPEG distortions become stronger, the stabilized model starts to significantly

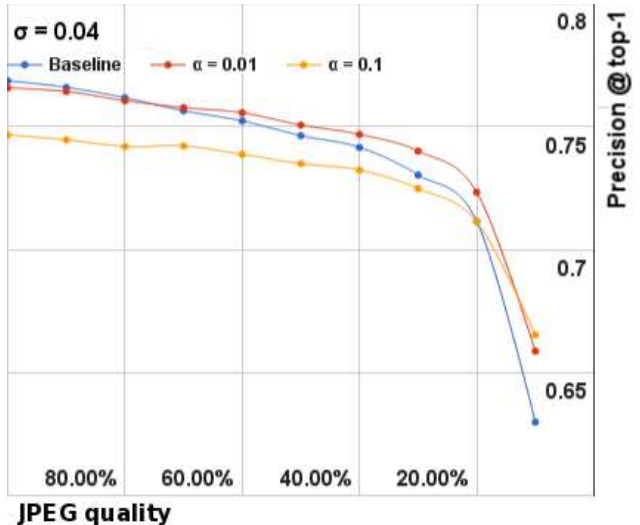


Figure 9: A comparison of the precision @ top-1 performance on the ImageNet classification task for different stability training hyper-parameters α , using JPEG compressed versions of the evaluation dataset at decreasing quality levels, using a fixed $\sigma = 0.04$. At the highest JPEG quality level, the baseline and stabilized models perform comparably. However, as the quality level decreases, the stabilized model starts to significantly outperform the baseline model.

outperform the baseline model. This qualitative behavior is visible for a wide range of hyper-parameters, for instance, using $\alpha = 0.01$ and $\sigma = 0.04$ results in better performance already below the 80% quality level.

6. Conclusion

In this paper we proposed stability training as a lightweight and effective method to stabilize deep neural networks against natural distortions in the visual input. Stability training makes the output of a neural network more robust by training a model to be constant on images that are copies of the input image with small perturbations. As such, our method can enable higher performance on noisy visual data than a network without stability training. We demonstrated this by showing that our method makes neural networks more robust against common types of distortions coming from random cropping, JPEG compression and thumbnail resizing. Additionally, we showed that using our method, the performance of stabilized models is significantly more robust for near-duplicate detection, similar-image ranking and classification on noisy datasets.

References

- [1] I. J. Goodfellow, J. Shlens, and C. Szegedy. Explaining and Harnessing Adversarial Examples. *arXiv:1412.6572 [cs, stat]*, 2014. 1, 2

- [2] S. Ioffe and C. Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *Proceedings of the 32nd International Conference on Machine Learning, ICML 2015, Lille, France, 6-11 July 2015*, pages 448–456, 2015. 5
- [3] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. In F. Pereira, C. J. C. Burges, L. Bottou, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems 25*, pages 1097–1105. Curran Associates, Inc., 2012. 1, 4
- [4] A. Kuznetsova, S. J. Hwang, B. Rosenhahn, and L. Sigal. Expanding object detector’s horizon: Incremental learning framework for object detection in videos. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2015. 2
- [5] I. Misra, A. Shrivastava, and M. Hebert. Watch and learn: Semi-supervised learning for object detectors from video. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2015. 2
- [6] T. Miyato, S.-i. Maeda, M. Koyama, K. Nakae, and S. Ishii. Distributional Smoothing with Virtual Adversarial Training. *arXiv:1507.00677 [cs, stat]*, July 2015. 2
- [7] A. Nguyen, J. Yosinski, and J. Clune. Deep neural networks are easily fooled: High confidence predictions for unrecognizable images. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2015. 2
- [8] A. Prest, C. Leistner, J. Civera, C. Schmid, and V. Ferrari. Learning object class detectors from weakly annotated video. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2012. 2
- [9] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, A. C. Berg, and L. Fei-Fei. ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision (IJCV)*, pages 1–42, April 2015. 7, 8
- [10] N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov. Dropout: A simple way to prevent neural networks from overfitting. *Journal of Machine Learning Research*, 15(1):1929–1958, Jan. 2014. 5
- [11] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich. Going deeper with convolutions. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2015. 1, 2, 4, 7, 8
- [12] C. Szegedy, W. Zaremba, I. Sutskever, J. Bruna, D. Erhan, I. Goodfellow, and R. Fergus. Intriguing properties of neural networks. *arXiv:1312.6199 [cs]*, Dec. 2013. 1, 2
- [13] J. Wang, Y. Song, T. Leung, C. Rosenberg, J. Wang, J. Philbin, B. Chen, and Y. Wu. Learning fine-grained image similarity with deep ranking. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2014. 1, 2, 3, 4, 5, 6, 8