

## mixup: BEYOND EMPIRICAL RISK MINIMIZATION

**Abstract**

- regularizes the neural network to favor simple linear behavior
- reduces the memorization of corrupt labels
- increases the robustness to adversarial examples
- stabilizes the training of generative adversarial networks

**Introduction**

- these neural networks share two commonalities
- 1 minimize their average error over the training data
- 2 size of networks scales linearly with the number of training examples

convergence of ERM is guaranteed as long as the size of the learning machine does not increase with the number of training data  
ERM allows large neural networks to memorize (instead of generalize)  
-> predictions drastically when evaluated on examples just outside the training distribution

### Vicinal Risk Minimization

human knowledge is required to describe a vicinity or neighborhood around each example  
additional virtual examples can be drawn from the vicinity distribution of the training examples  
ex) horizontal reflections, slight rotations, and mild scalings  
-> dataset-dependent, and thus requires the use of expert knowledge  
does not model the vicinity relation across examples of different classes.

### Contribution

$\tilde{x} = \lambda x_i + (1 - \lambda)x_j,$       where  $x_i, x_j$  are raw input vectors  
 $\tilde{y} = \lambda y_i + (1 - \lambda)y_j,$       where  $y_i, y_j$  are one-hot label encodings

### Virtual Training Examples

a new state-of-the-art performance  
increases the robustness of neural networks

### From Empirical Risk Minimization to Mixup

$$\mu(\tilde{x}, \tilde{y} | x_i, y_i) = \frac{1}{n} \sum_j^n \mathbb{E}_{\lambda} [\delta(\tilde{x} = \lambda \cdot x_i + (1 - \lambda) \cdot x_j, \tilde{y} = \lambda \cdot y_i + (1 - \lambda) \cdot y_j)],$$





implementation of mixup training is straightforward, and introduces a minimal computation overhead

$$\tilde{x} = \lambda x_i + (1 - \lambda) x_j,$$

$$\tilde{y} = \lambda y_i + (1 - \lambda) y_j,$$

- 1 combinations of three or more examples with weights sampled from a Dirichlet distribution does not provide further gain,
- 2 single data loader to obtain one minibatch,
- 3 interpolating only between inputs with equal label did not lead to the performance gains

**What is mixup doing?**

Model	Method	Epochs	Top-1 Error	Top-5 Error
ResNet-50	ERM (Goyal et al., 2017)	90	23.5	-
	<i>mixup</i> $\alpha = 0.2$	90	<b>23.3</b>	<b>6.6</b>
ResNet-101	ERM (Goyal et al., 2017)	90	22.1	-
	<i>mixup</i> $\alpha = 0.2$	90	<b>21.5</b>	<b>5.6</b>
ResNeXt-101 32*4d	ERM (Xie et al., 2016)	100	21.2	-
	ERM	90	21.2	5.6
	<i>mixup</i> $\alpha = 0.4$	90	<b>20.7</b>	<b>5.3</b>
ResNeXt-101 64*4d	ERM (Xie et al., 2016)	100	20.4	5.3
	<i>mixup</i> $\alpha = 0.4$	90	<b>19.8</b>	<b>4.9</b>
ResNet-50	ERM	200	23.6	7.0
	<i>mixup</i> $\alpha = 0.2$	200	<b>22.1</b>	<b>6.1</b>
ResNet-101	ERM	200	22.0	6.1
	<i>mixup</i> $\alpha = 0.2$	200	<b>20.8</b>	<b>5.4</b>
ResNeXt-101 32*4d	ERM	200	21.3	5.9
	<i>mixup</i> $\alpha = 0.4$	200	<b>20.1</b>	<b>5.0</b>

Table 1: Validation errors for ERM and *mixup* on the development set of ImageNet-2012.

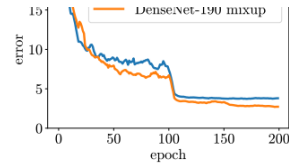
model trained with mixup is more stable in terms of model predictions and gradient norms in-between training samples.

Experiments

Dataset	Model	ERM	<i>mixup</i>
---------	-------	-----	--------------



CIFAR-10	PreAct ResNet-18	5.6	<b>4.2</b>
	WideResNet-28-10	3.8	<b>2.7</b>
	DenseNet-BC-190	3.7	<b>2.7</b>
CIFAR-100	PreAct ResNet-18	25.6	<b>21.1</b>
	WideResNet-28-10	19.4	<b>17.5</b>
	DenseNet-BC-190	19.0	<b>16.8</b>



Model	Method	Validation set	Test set
LeNet	ERM	<b>9.8</b>	<b>10.3</b>
	<i>mixup</i> ( $\alpha = 0.1$ )	10.1	10.8
	<i>mixup</i> ( $\alpha = 0.2$ )	10.2	11.3
VGG-11	ERM	5.0	4.6
	<i>mixup</i> ( $\alpha = 0.1$ )	4.0	3.8
	<i>mixup</i> ( $\alpha = 0.2$ )	<b>3.9</b>	<b>3.4</b>

Label corruption	Method	Test error		Training error	
		Best	Last	Real	Corrupted
20%	ERM	12.7	16.6	0.05	0.28
	ERM + dropout ( $p = 0.7$ )	8.8	10.4	5.26	83.55
	<i>mixup</i> ( $\alpha = 8$ )	<b>5.9</b>	6.4	2.27	86.32
	<i>mixup</i> + dropout ( $\alpha = 4, p = 0.1$ )	6.2	<b>6.2</b>	1.92	85.02
50%	ERM	18.8	44.6	0.26	0.64
	ERM + dropout ( $p = 0.8$ )	14.1	15.5	12.71	86.98
	<i>mixup</i> ( $\alpha = 32$ )	11.3	12.7	5.84	85.71
	<i>mixup</i> + dropout ( $\alpha = 8, p = 0.3$ )	<b>10.9</b>	<b>10.9</b>	7.56	87.90
80%	ERM	36.5	73.9	0.62	0.83
	ERM + dropout ( $p = 0.8$ )	30.9	35.1	29.84	86.37
	<i>mixup</i> ( $\alpha = 32$ )	25.3	30.9	18.92	85.44
	<i>mixup</i> + dropout ( $\alpha = 8, p = 0.3$ )	<b>24.0</b>	<b>24.8</b>	19.70	87.67

Metric	Method	FGSM	I-FGSM
Top-1	ERM	90.7	99.9
	<i>mixup</i>	<b>75.2</b>	99.6
Top-5	ERM	63.1	93.4
	<i>mixup</i>	<b>49.1</b>	95.8

Metric	Method	FGSM	I-FGSM
Top-1	ERM	57.0	57.3
	<i>mixup</i>	<b>46.0</b>	<b>40.9</b>
Top-5	ERM	24.8	18.1
	<i>mixup</i>	<b>17.4</b>	<b>11.8</b>

(a) White box attacks.

(b) Black box attacks.

TABULAR DATA

Dataset	ERM	<i>mixup</i>	Dataset	ERM	<i>mixup</i>
Abalone	74.0	73.6	Htru2	2.0	2.0
Arcene	57.6	<b>48.0</b>	Iris	21.3	<b>17.3</b>
Arrhythmia	56.6	<b>46.3</b>	Phishing	16.3	15.2

Table 4: ERM and *mixup* classification errors on the UCI datasets.

#### STABILIZATION OF GENERATIVE ADVERSARIAL NETWORKS (GANS)

$$\max_g \min_d \mathbb{E}_{x,z,\lambda} \ell(d(\lambda x + (1 - \lambda)g(z)), \lambda).$$

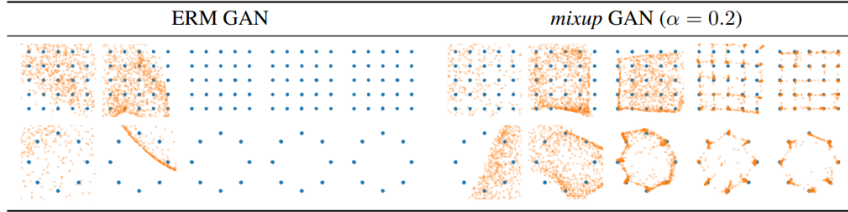


Figure 5: Effect of *mixup* on stabilizing GAN training at iterations 10, 100, 1000, 10000, and 20000.

#### ABLATION STUDIES

Method	Specification	Modified		Weight decay	
		Input	Target	$10^{-4}$	$5 \times 10^{-4}$
ERM		✗	✗	5.53	5.18
<i>mixup</i>	AC + RP	✓	✓	<b>4.24</b>	4.68
	AC + KNN	✓	✓	4.98	5.26
mix labels and latent representations (AC + RP)	Layer 1	✓	✓	4.44	<b>4.51</b>
	Layer 2	✓	✓	4.56	4.61
	Layer 3	✓	✓	5.39	5.55
	Layer 4	✓	✓	5.95	5.43
	Layer 5	✓	✓	5.39	5.15
mix inputs only	SC + KNN (Chawla et al., 2002)	✓	✗	5.45	5.52
	AC + KNN	✓	✗	5.43	5.48
	SC + RP	✓	✗	5.23	5.55
	AC + RP	✓	✗	5.17	5.72
label smoothing (Szegedy et al., 2016)	$\epsilon = 0.05$	✗	✓	5.25	5.02
	$\epsilon = 0.1$	✗	✓	5.33	5.17
	$\epsilon = 0.2$	✗	✓	5.34	5.06
mix inputs +	$\epsilon = 0.05$	✓	✓	5.02	5.40

label smoothing (AC + RP)	$\epsilon = 0.1$	✓	✓	5.08	5.09
	$\epsilon = 0.2$	✓	✓	4.98	5.06
	$\epsilon = 0.4$	✓	✓	5.25	5.39
add Gaussian noise to inputs	$\sigma = 0.05$	✓	✗	5.53	5.04
	$\sigma = 0.1$	✓	✗	6.41	5.86
	$\sigma = 0.2$	✓	✗	7.16	7.24

mixup =

random convex combination of raw inputs

convex combination of one-hot label encodings

ablation study로서 알게 된 점들

First, mixup is the best data augmentation method we test, and is significantly better than the second best method (mix input + label smoothing)

Second, the effect of regularization can be seen by comparing the test error with a small weight decay ( $10^{-4}$ ) with a large one ( $5 \times 10^{-4}$ ).

SMOTE algorithm (Chawla et al., 2002) does not lead to a noticeable gain in performance.

## Related Work

Data augmentation

rotation, translation, cropping, resizing, flipping

random erasing

in speech recognition, noise injection

mixup 단점

only operate among the nearest neighbors within a certain class at the input / feature level

does not account for changes in the corresponding labels

label smoothing

penalizing high-confidence softmax distributions

두 개의 단점

-> independent from the associated feature values

mixup enjoys several desirable aspects of previous data augmentation and regularization schemes

not require significant domain knowledge

## Discussion

with increasingly large  $\alpha$ , the training error on real data increases, while the generalization gap decreases.

do not yet have a good theory for understanding the 'sweet spot' of this bias-variance trade-off

increasing the model capacity would make training error less sensitive to large  $\alpha$  -> giving mixup a more significant advantage.

several possibilities for further exploration

First, is it possible to make similar ideas work on other types of supervised learning problems, such as regression and structured prediction?

Second, can similar methods prove helpful beyond supervised learning?

ex) Can we extend mixup to feature-label extrapolation to guarantee a robust model behavior far away from the training data?