Differential Operators for Structured Adversarial Examples

Hossein Mobahi

Joint work with Calvin Luo Samy Bengio

Google Research

April 22, 2020

Hossein Mobahi Differential Operators for Structured Adversarial Examples

# Data Augmentation



- Data augmentation is an important ingredient for achieving SOTA on image classification tasks.
- Currently incorporated into training by sampling; may not scale to larger space of transformations.

### Adversarial Examples



Original Image Adversarial Perturbation Very rich set of perturbations, that often fall outside of natural transformations.

## Goal

# Model space of transformations, then use adversarial perturbation



**Original Image** 



#### Raw Adversarial Perturbation



Recolorized Image

Restricted Perturbation

Hossein Mobahi

Differential Operators for Structured Adversarial Examples

# Notation

- Given a loss function L(w; z, l) that takes as its input the model parameters w, input image z, and label l.
- Suppose components of *z* (i.e. pixels) are sampled from a continuous image flow *Z*(*x*, *y*, *t*) at a specific *t*,

$$\boldsymbol{z} \triangleq \left[ Z(x_1, y_1, t), \dots, Z(x_n, y_m, t) \right].$$
(1)

Thus loss can be written as,

$$L(Z(x_1, y_1, t), \dots, Z(x_n, y_m, t)),$$
 (2)

(w and l dropped for brevity).

# Generating Adversarial Examples

 Seek a perturbation in the image that maximally increases the loss.

$$\frac{\partial}{\partial t}L = \langle \nabla_{\boldsymbol{z}} L(\boldsymbol{z}), [\dot{Z}(x_1, y_1, t), \dots, \dot{Z}(x_n, y_m, t)] \rangle.$$
 (3)

We will demonstrate a few examples to construct Z.

# In a geometric transform, points may move, but cannot change color/brightness.



Image from Michael Black's PhD thesis

#### Consider brightness constancy assumption,

$$\frac{d}{dt}Z(x(t), y(t), t) = 0, \qquad (4)$$

which is equivalent to,

$$Z_x(x(t), y(t), t)\dot{x}(t)$$
(5)

+ 
$$Z_y(x(t), y(t), t)\dot{y}(t)$$
 (6)

$$+ \dot{Z}(x(t), y(t), t) \tag{7}$$

Applying this yields,

$$egin{aligned} rac{\partial}{\partial t}L &= -\langle 
abla_{oldsymbol{z}}L(oldsymbol{z})\,,\,[ & & Z_x(x_1,y_1,t)\dot{x}(x_1,y_1,t)\ & & +Z_y(x_1,y_1,t)\dot{y}(x_1,y_1,t) \end{aligned}$$

$$Z_{x}(x_{n}, y_{m}, t)\dot{x}(x_{n}, y_{m}, t) + Z_{y}(x_{n}, y_{m}, t)\dot{y}(x_{n}, y_{m}, t)]\rangle.$$
(9)

We are interested in this dynamics at the reference point t = 0. For brevity, denote  $Z_x(x_i, y_j, t = 0)$  by  $Z_{x,i,j}$  and  $\dot{x}((x_i, y_j), t = 0)$  by  $\dot{x}_{i,j}$ . Also denote each component of  $\nabla_{\boldsymbol{x}} L(\boldsymbol{x})$  by  $\ell_{i,j}$ . Thus in matrix form,

$$rac{\partial}{\partial t}L = - \begin{pmatrix} oldsymbol{d}_x \end{pmatrix}^T \begin{pmatrix} \dot{oldsymbol{x}} \\ \dot{oldsymbol{y}} \end{pmatrix}^T \begin{pmatrix} \dot{oldsymbol{x}} \\ \dot{oldsymbol{y}} \end{pmatrix}$$
 (10)

where,

$$N \triangleq nm$$
 (11)

$$\boldsymbol{d}_{\boldsymbol{x} \boldsymbol{N} \times \boldsymbol{J}} \triangleq \begin{pmatrix} \ell_{1,1} Z_{\boldsymbol{x},1,1} & \dots & \ell_{m,n} Z_{\boldsymbol{x},m,n} \end{pmatrix}$$
(12)

$$d_{y_{N\times 1}} \triangleq (\ell_{1,1}Z_{y,1,1} \ldots \ell_{m,n}Z_{x,m,n})$$
(13)

$$\dot{\boldsymbol{x}}_{N \times 1} \triangleq (\dot{x}_{1,1} \ldots \dot{x}_{m,n})$$
 (14)

$$\dot{\boldsymbol{y}}_{N\times 1} \triangleq (\dot{y}_{1,1} \ldots \dot{y}_{m,n})$$
 (15)

Hossein Mobahi Differential Operators for Structured Adversarial Examples

We seek motion field ( $\dot{x}$  and  $\dot{y}$ ) that maximizes  $\frac{\partial}{\partial t}L$ . To keep velocity bounded, restrict  $||\dot{x}||$  and  $||\dot{y}||$ . Also to have coherency in the field impose smoothness.

$$\begin{aligned} (\dot{\boldsymbol{x}}^*, \dot{\boldsymbol{y}}^*) &= \arg\min_{(\dot{\boldsymbol{x}}, \dot{\boldsymbol{y}})} & \begin{pmatrix} \boldsymbol{d}_x \\ \boldsymbol{d}_y \end{pmatrix}^T \begin{pmatrix} \dot{\boldsymbol{x}} \\ \dot{\boldsymbol{y}} \end{pmatrix} + \frac{1}{2} \alpha (\|\dot{\boldsymbol{x}}\|^2 + \|\dot{\boldsymbol{y}}\|^2) \\ &+ \frac{1}{2} \beta (\|\boldsymbol{D}_x \dot{\boldsymbol{x}}\|^2 + \|\boldsymbol{D}_y \dot{\boldsymbol{x}}\|^2) \\ &+ \frac{1}{2} \beta (\|\boldsymbol{D}_x \dot{\boldsymbol{y}}\|^2 + \|\boldsymbol{D}_y \dot{\boldsymbol{y}}\|^2) \,, \end{aligned}$$

where  $D_x$  and  $D_y$  are differential operators along x and y direction, represented in matrix form. For brevity, let  $P \triangleq D_x^T D_x + D_y^T D_y$ . We obtain,

$$(\dot{\boldsymbol{x}}^*, \dot{\boldsymbol{y}}^*) = \arg\min_{(\dot{\boldsymbol{x}}, \dot{\boldsymbol{y}})} \begin{pmatrix} \boldsymbol{d}_x \\ \boldsymbol{d}_y \end{pmatrix}^T \begin{pmatrix} \dot{\boldsymbol{x}} \\ \dot{\boldsymbol{y}} \end{pmatrix}$$
 (16)

$$- \frac{1}{2}\alpha(\dot{\boldsymbol{x}}^T\boldsymbol{I}\dot{\boldsymbol{x}} + \dot{\boldsymbol{y}}^T\boldsymbol{I}\dot{\boldsymbol{y}})$$
(17)

$$- \frac{1}{2}\beta(\dot{\boldsymbol{x}}^{T}\boldsymbol{P}\dot{\boldsymbol{x}}+\dot{\boldsymbol{y}}^{T}\boldsymbol{P}\dot{\boldsymbol{y}}). \tag{18}$$

#### Letting $\gamma \triangleq \frac{\beta}{\alpha}$ , we can express

 $(\dot{x})$ 

 $(\dot{x}^*,$ 

$$(\hat{x}, \dot{y}^*) = \arg\min_{(\dot{x}, \dot{y})} \qquad \begin{pmatrix} d_x \\ d_y \end{pmatrix}^T \begin{pmatrix} \dot{x} \\ \dot{y} \end{pmatrix}$$
 (19)

$$\frac{1}{2}\alpha(\dot{\boldsymbol{x}}^T\boldsymbol{I}\dot{\boldsymbol{x}}+\dot{\boldsymbol{y}}^T\boldsymbol{I}\dot{\boldsymbol{y}})$$
(20)

+ 
$$\frac{1}{2}\beta(\dot{\boldsymbol{x}}^T\boldsymbol{P}\dot{\boldsymbol{x}}+\dot{\boldsymbol{y}}^T\boldsymbol{P}\dot{\boldsymbol{y}}).$$
 (21)

as,

$$(\mathbf{y}^*) = \arg\min_{(\dot{x}, \dot{y})} \qquad \begin{pmatrix} \mathbf{d}_x \\ \mathbf{d}_y \end{pmatrix}^T \begin{pmatrix} \dot{x} \\ \dot{y} \end{pmatrix}$$
(22)

+ 
$$\frac{1}{2}\alpha((\dot{\boldsymbol{x}}^T\boldsymbol{I}\dot{\boldsymbol{x}}+\dot{\boldsymbol{y}}^T\boldsymbol{I}\dot{\boldsymbol{y}})$$
 (23)

+ 
$$\gamma(\dot{\boldsymbol{x}}^T \boldsymbol{P} \dot{\boldsymbol{x}} + \dot{\boldsymbol{y}}^T \boldsymbol{P} \dot{\boldsymbol{y}})).$$
 (24)

This is a convex objective function in  $(\dot{x}, \dot{y})$  and its minimizer can be obtained by zero crossing its gradient,

$$d_x + \alpha (I + \gamma P) \dot{x} = 0 \tag{25}$$

$$d_y + \alpha (I + \gamma P) \dot{y} = 0.$$
(26)

Therefore,

$$\dot{\boldsymbol{x}} = -\frac{1}{\alpha} (\boldsymbol{I} + \gamma \boldsymbol{P})^{-1} \boldsymbol{d}_{\boldsymbol{x}}$$
(27)

$$\dot{\boldsymbol{y}} = -\frac{1}{\alpha} (\boldsymbol{I} + \gamma \boldsymbol{P})^{-1} \boldsymbol{d}_y.$$
 (28)

Note that the matrix  $(I + \gamma P)^{-1}$  solely depends on regularizer, and *does not depend on data* (data being  $d_x$  and  $d_y$ ). Thus, one can compute  $(I + \gamma P)^{-1}$  offline given image dimensions and trade off parameter  $\gamma$ .

Once the motion field  $(\dot{x}, \dot{y})$  is computed, it can be used to warp *I* to a new form via the motion field.

# Example

#### Examples of ResNet trained on ImageNet:



Garbage truck. Confidence 60.95%





Minibus. Confidence 33.90%



Redbone. Confidence of 44.95 Boxer. Confidence 19.34% The flow warped the back of the truck to be raised - a trait common in minibuses. Similarly, for the dog image, it warps the white cloth onto the dog's leg (redbones are normally entirely red, boxers are commonly red dogs that feature patches of white fur).

### Photometric

Goal is to recolorize the image; preserve contours, object boundaries, edges, etc., but okay to change colors in the image. Consider some notion of edge strength, such as,

$$E(x, y, t) \triangleq Z_{Ax}(x, y, t)^2 + Z_{Ay}(x, y, t)^2$$
 (29)

$$+Z_{Bx}(x,y,t)^2 + Z_{By}(x,y,t)^2$$
(30)

$$+Z_{Cx}(x,y,t)^2 + Z_{Cy}(x,y,t)^2.$$
 (31)

We penalize changes in edge strength anywhere that has already weak edges.

$$P(t) \triangleq \int_{\mathcal{X}} \int_{\mathcal{Y}} \left( \frac{1}{E(x,y,0)} \frac{\partial}{\partial t} E(x,y,t) \right)^2 dx \, dy \, (32)$$

$$\int_{\mathcal{Y}} \left( \frac{1}{E(x,y,0)} \left( Z_{Ax}(x,y,t) \dot{Z}_{Ax} + Z_{Ay}(x,y,t) \dot{Z}_{Ay} \right) + Z_{Bx}(x,y,t) \dot{Z}_{Bx} + Z_{By}(x,y,t) \dot{Z}_{By} + Z_{Cx}(x,y,t) \dot{Z}_{Cx} + Z_{Cy}(x,y,t) \dot{Z}_{Cy} \right) + Z_{Cx}(x,y,t) \dot{Z}_{Cx} + Z_{Cy}(x,y,t) \dot{Z}_{Cy} \quad (35)$$

$$(36)$$

Differential Operators for Structured Adversarial Examples

Hossein Mobahi

### Photometric

# We are interested in P(t) at t = 0 and thus drop dependency of the on t from the notation,

$$P = \int_{\mathcal{X}} \int_{\mathcal{Y}} \left( \frac{1}{E(x,y)} \left( \qquad Z_{A_x}(x,y) \dot{Z}_{A_x} + Z_{A_y}(x,y) \dot{Z}_{A_y} \right) \right)$$
(37)

+ 
$$Z_{B_x}(x,y)\dot{Z_B_x} + Z_{B_y}(x,y)\dot{Z_B_y}$$
 (38)

$$+ Z_{C_x}(x,y)\dot{Z_C_x} + Z_{C_y}(x,y)\dot{Z_C_y}$$
(39)

$$\left( \int_{-\infty}^{\infty} dx \, dy \, . \right)$$

This integral can be approximated numerically, e.g. sampling over *evenly spaced* grid of steps  $\delta_x$  and  $\delta_y$ ,

$$P pprox \hat{P} \triangleq \delta_x \delta_y \| \boldsymbol{M}_A \dot{\boldsymbol{z}}_A + \boldsymbol{M}_B \dot{\boldsymbol{z}}_B + \boldsymbol{M}_C \dot{\boldsymbol{z}}_C \|^2 \,,$$
 (41)

where,

$$\begin{aligned} \boldsymbol{M}_{\alpha N \times N} & \triangleq & \operatorname{diag}(\boldsymbol{z}_{\alpha x} \div \boldsymbol{e}) \boldsymbol{D}_{x} + \operatorname{diag}(\boldsymbol{z}_{\alpha y} \div \boldsymbol{e}) \boldsymbol{D}_{y} \\ \boldsymbol{e}_{N \times 1} & \triangleq & \left\| \left[ \boldsymbol{z}_{Ax} \, | \, \boldsymbol{z}_{Ay} \, | \, \boldsymbol{z}_{Bx} \, | \, \boldsymbol{z}_{By} \, | \, \boldsymbol{z}_{Cx} \, | \, \boldsymbol{z}_{Cy} \right]_{N \times 6} \right\|, \end{aligned}$$

and ÷ denotes element-wise division.

### Photometric

# Let $z \triangleq [z_A^T | z_B^T | z_C^T]^T$ . The goal is to find $\dot{z}$ that maximally increases the loss under this penalty.

$$\dot{\boldsymbol{z}}_{3N\times 1}^* \triangleq \arg\min_{\dot{\boldsymbol{z}}} - \langle \nabla_{\boldsymbol{z}} L(\boldsymbol{z}), \, \dot{\boldsymbol{z}} \rangle \, + \, \frac{1}{2} \lambda \hat{P} \,. \tag{42}$$

This is a convex objective function in  $\dot{z}$  and its solution can be obtain by zero crossing the gradient,

$$-\nabla_{\boldsymbol{z}_A} L(\boldsymbol{z}) + \lambda \boldsymbol{M}_A^T (\boldsymbol{M}_A \dot{\boldsymbol{z}}_A + \boldsymbol{M}_B \dot{\boldsymbol{z}}_B + \boldsymbol{M}_C \dot{\boldsymbol{z}}_C) = \boldsymbol{0}$$
(43)

$$-\nabla_{\boldsymbol{z}_B} L(\boldsymbol{z}) + \lambda \boldsymbol{M}_B^T (\boldsymbol{M}_A \dot{\boldsymbol{z}}_A + \boldsymbol{M}_B \dot{\boldsymbol{z}}_B + \boldsymbol{M}_C \dot{\boldsymbol{z}}_C) = \boldsymbol{0}$$
(44)

$$-\nabla_{\boldsymbol{z}_{C}}L(\boldsymbol{z}) + \lambda \boldsymbol{M}_{C}^{T}(\boldsymbol{M}_{A}\dot{\boldsymbol{z}}_{A} + \boldsymbol{M}_{B}\dot{\boldsymbol{z}}_{B} + \boldsymbol{M}_{C}\dot{\boldsymbol{z}}_{C}) = \boldsymbol{0}.$$
 (45)

This can be written more compactly as,

$$\boldsymbol{M}_{3N\times 3N} \dot{\boldsymbol{z}}_{3N\times 1} = \frac{1}{\lambda} \nabla_{\boldsymbol{z}} L(\boldsymbol{z}) \,, \tag{46}$$

where,

$$\boldsymbol{M}_{AMAAM} \triangleq \begin{pmatrix} \boldsymbol{M}_{A}^{T}\boldsymbol{M}_{A} & \boldsymbol{M}_{A}^{T}\boldsymbol{M}_{B} & \boldsymbol{M}_{A}^{T}\boldsymbol{M}_{C} \\ \boldsymbol{M}_{B}^{T}\boldsymbol{M}_{A} & \boldsymbol{M}_{B}^{T}\boldsymbol{M}_{B} & \boldsymbol{M}_{B}^{T}\boldsymbol{M}_{C} \\ \boldsymbol{M}_{C}^{T}\boldsymbol{M}_{A} & \boldsymbol{M}_{C}^{T}\boldsymbol{M}_{B} & \boldsymbol{M}_{C}^{T}\boldsymbol{M}_{C} \end{pmatrix} .$$
(47)

Hossein Mobahi

Differential Operators for Structured Adversarial Examples

#### Thus,

$$\dot{\boldsymbol{z}}_{3N\times 1}^{*} = \frac{1}{\lambda} \boldsymbol{M}_{3N\times 3N}^{-1} \big( \nabla_{\boldsymbol{z}} \boldsymbol{L}(\boldsymbol{z}) \big)_{3N\times 1},$$
(48)

While the matrix M depends on the data (images), it has nothing to do with the loss function L (thus independent of learning model). Thus for each image z, the associated matrix  $M^{-1}$  can be computed once. Then the computed  $M^{-1}$  can the be used for *any learning model* in the future. Clearly  $M^{-1}$  captures some intrinsic properties of the data that

are useful later for photometric transforms.

- There are still challenges using this in practice: although *M*<sup>-1</sup> can be computed once, applying such giant matrix to the gradient vector can be expensive.
- We use SVD to extract leading eigenvectors of M<sup>-1</sup> and only use those to achieve an efficient approximation.

# **Photometric Examples**

















#### Original Recolorized Big Eigenvec.

Hossein Mobahi Differential Operators for Structured Adversarial Examples

## Experiments

	CIFAR-10	CIFAR-100	STL-10
BasicAug	95.56%	77.41%	84.60%
AutoAug	95.70%	77.41%	84.46%
RandAug	94.95%	74.56%	88.75%
Flow + BasicAug	95.79%	76.82%	85.41%
Flow + AutoAug	95.77%	77.26%	85.39%
Flow + RandAug	93.22%	70.14%	87.96%
Recolor + BasicAug	95.80%	77.45%	86.11%
Recolor + AutoAug	95.86%	76.97%	85.21%
Recolor + RandAug	94.78%	74.18%	88.84%

Table: Reported accuracy of a Wide-ResNet-32-10 model trained from scratch using the listed augmentations

	CIFAR-10	CIFAR-100	STL-10
No Augmentation	97.52%	85.68%	97.55%
AutoAugment	97.44%	84.09%	97.32%
RandAugment	96.39%	80.38%	97.14%
Flow	97.60%	85.87%	97.61%
Recolorization	97.34%	85.25%	97.61%

Table: Reported accuracy after finetuning on a EfficientNet-B0 checkpoint for 350 epochs.

# Directions for Future Work

Most important directions for future Work

- Additional transformations that can be modeled in similar way.
- Learning the linear operators from another training set.

#### **Thank You!**

Hossein Mobahi Differential Operators for Structured Adversarial Examples