

ECE 595: Homework 5

Yi Qiao, Class ID 187
(Spring 2019)

Exercise 1: Adversarial Attacks on Gaussian Classifier

(a) minimum-norm attacks

(i) minimum l_2 and l_∞ attack

Since we only have 2 classes, the question becomes

$$\begin{aligned} & \underset{\mathbf{x}}{\text{minimize}} \quad \|\mathbf{x} - \mathbf{x}_0\| \\ & \text{subject to} \quad \mathbf{w}^T \mathbf{x} + w_0 = 0 \end{aligned} \tag{1}$$

using l_2 norm

the problem is the same as

$$\begin{aligned} & \underset{\mathbf{x}}{\text{minimize}} \quad \frac{1}{2} \|\mathbf{x} - \mathbf{x}_0\|_2^2 \\ & \text{subject to} \quad \mathbf{w}^T \mathbf{x} + w_0 = 0 \end{aligned} \tag{2}$$

The lagrangian is

$$\mathcal{L}(\mathbf{x}, \lambda) = \frac{1}{2} \|\mathbf{x} - \mathbf{x}_0\|_2^2 + \lambda(\mathbf{w}^T \mathbf{x} + w_0) \tag{3}$$

Taking the derivative

$$\begin{aligned} \nabla_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \lambda) &= \mathbf{x} - \mathbf{x}_0 + \lambda \mathbf{w} = 0 \\ \frac{\partial}{\partial \lambda} \mathcal{L}(\mathbf{x}, \lambda) &= \mathbf{w}^T \mathbf{x} + w_0 = 0 \end{aligned} \tag{4}$$

$$\begin{aligned} \lambda \mathbf{w} &= \mathbf{x}_0 - \mathbf{x} \\ \lambda \mathbf{w}^T \mathbf{w} &= \mathbf{w}^T \mathbf{x}_0 - \mathbf{w}^T \mathbf{x} \end{aligned} \tag{5}$$

$$\lambda = (\mathbf{w}^T \mathbf{w})^{-1} (\mathbf{w}^T \mathbf{x}_0 + w_0)$$

$$\begin{aligned} \mathbf{x} &= \mathbf{x}_0 - \lambda \mathbf{w} \\ &= \mathbf{x}_0 - \frac{\mathbf{w}(\mathbf{w}^T \mathbf{x}_0 + w_0)}{\|\mathbf{w}\|_2^2} \end{aligned} \tag{6}$$

using l_∞ norm

$$\begin{aligned} & \underset{\mathbf{x}}{\text{minimize}} \quad \|\mathbf{x} - \mathbf{x}_0\|_\infty \\ & \text{subject to} \quad \mathbf{w}^T \mathbf{x} + w_0 = 0 \end{aligned} \tag{7}$$

Let $\mathbf{r} = \mathbf{x} - \mathbf{x}_0$, $b_0 = -(\mathbf{w}^T \mathbf{x}_0 + w_0)$, the problem becomes:

$$\begin{aligned} & \underset{\mathbf{x}}{\text{argmin}} \quad \|\mathbf{x} - \mathbf{x}_0\|_\infty \\ & \text{subject to} \quad \mathbf{w}^T \mathbf{r} = b_0 \end{aligned} \tag{8}$$

The lagrangian is

$$\mathcal{L}(\mathbf{r}, \lambda) = \|\mathbf{r}\|_\infty + \lambda(b_0 - \mathbf{w}^T \mathbf{r}) \quad (9)$$

Taking derivative,

$$\frac{\partial}{\partial \lambda} \mathcal{L}(\mathbf{r}, \lambda) = b_0 - \mathbf{w}^T \mathbf{r} = 0 \quad (10)$$

By Holder's Inequality:

$$\begin{aligned} |b_0| &= |\mathbf{w}^T \mathbf{r}| \leq \|\mathbf{w}\|_1 \|\mathbf{r}\|_\infty \\ \|\mathbf{r}\|_\infty &\geq \frac{|b_0|}{\|\mathbf{w}\|_1} \end{aligned} \quad (11)$$

Consider $\mathbf{r} = \eta \cdot \text{sign}(\mathbf{w})$, for some constant η tbd. We can show that

$$\|\mathbf{r}\|_\infty = \underset{i}{\operatorname{argmax}} |\eta \cdot \text{sign}(w_i)| = |\eta| \quad (12)$$

let $\eta = \frac{b_0}{\|\mathbf{w}\|_1} \cdot \text{sign}(\mathbf{w})$, then we have,

$$\|\mathbf{r}\|_\infty = |\eta| = \frac{b_0}{\|\mathbf{w}\|_1} \quad (13)$$

Lower bound is achieved, thus the solution is,

$$\mathbf{r} = \frac{|b_0|}{\|\mathbf{w}\|_1} \cdot \text{sign}(\mathbf{w}) \quad (14)$$

(ii) DeepFool attack

$$\begin{aligned} &\underset{\mathbf{x}}{\operatorname{argmin}} \|\mathbf{x} - \mathbf{x}_0\|_2^2 \\ &\text{subject to } g(\mathbf{x}) = 0 \end{aligned} \quad (15)$$

First order approximation

$$g(\mathbf{x}) \approx g(\mathbf{x}^{(k)}) + \nabla_{\mathbf{x}} g(\mathbf{x}^{(k)})^T (\mathbf{x} - \mathbf{x}^{(k)}) \quad (16)$$

Then the problem can be approximate by

$$\begin{aligned} &\underset{\mathbf{x}}{\operatorname{argmin}} \|\mathbf{x} - \mathbf{x}_0\|_2^2 \\ &\text{subject to } g(\mathbf{x}^{(k)}) + \nabla_{\mathbf{x}} g(\mathbf{x}^{(k)})^T (\mathbf{x} - \mathbf{x}^{(k)}) = 0 \end{aligned} \quad (17)$$

Let $\mathbf{w}^{(k)} = \nabla_{\mathbf{x}} g(\mathbf{x}^{(k)})$ and $w_0^{(k)} = g(\mathbf{x}^{(k)}) - \nabla_{\mathbf{x}} g(\mathbf{x}^{(k)})^T \mathbf{x}^{(k)}$

Then the problem is equivalent to

$$\begin{aligned} &\underset{\mathbf{x}}{\operatorname{argmin}} \|\mathbf{x} - \mathbf{x}_0\|_2^2 \\ &\text{subject to } (\mathbf{w}^{(k)})^T \mathbf{x} + w_0^{(k)} = 0 \end{aligned} \quad (18)$$

This is the same problem as minimum l_2 norm attack, Thus the solution will be,

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \frac{((\mathbf{w}^{(k)})^T \mathbf{x}^{(k)} + w_0^{(k)}) \mathbf{w}^{(k)}}{\|\mathbf{w}^{(k)}\|_2^2} \quad (19)$$

substitute \mathbf{w} and w_0 back, we get

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \left(\frac{g(\mathbf{x}^{(k)})}{\|\nabla_{\mathbf{x}} g(\mathbf{x}^{(k)})\|^2} \right) \nabla_{\mathbf{x}} g(\mathbf{x}^{(k)}) \quad (20)$$

(iii) An example DeepFool never converge

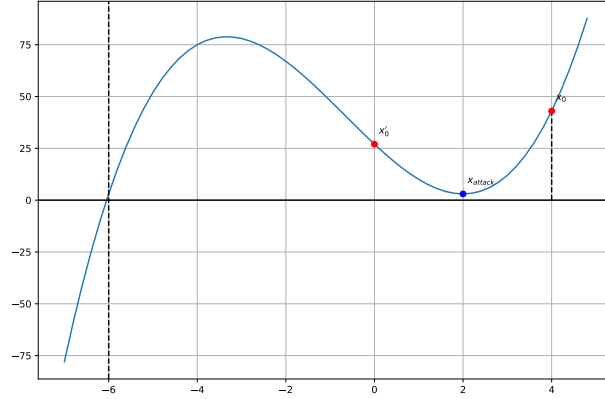


Figure 1: 1D example classifier

Above is a example discriminate function $g(x)$. there are two classes, separated by the vertical line at $g(x) = 0$. If x_0 at the position shown above, the Deep Fool attack will result in x_{attack} in the end, and thus never converge to $g(x) = 0$.

(b) maximum-allowable attack

(i) l_∞ attack in the linear case

The problem is,

$$\begin{aligned} \underset{\mathbf{x}}{\operatorname{argmin}} \quad & \mathbf{w}^T \mathbf{x} + w_0 \\ \text{subject to} \quad & \|\mathbf{x} - \mathbf{x}_0\|_\infty < \eta \end{aligned} \quad (21)$$

let $\mathbf{x} = \mathbf{x}_0 + \mathbf{r}$, $b_0 = (\mathbf{w}^T \mathbf{x}_0 + w_0)$, the problem becomes,

$$\begin{aligned} \underset{\mathbf{r}}{\operatorname{argmin}} \quad & \mathbf{w}^T \mathbf{r} + b_0 \\ \text{subject to} \quad & \|\mathbf{r}\|_\infty < \eta \end{aligned} \quad (22)$$

by Holder's inequality,

$$\mathbf{w}^T \mathbf{r} \geq -\|\mathbf{r}\|_\infty \|\mathbf{w}\|_1 \geq -\eta \|\mathbf{w}\|_1 \quad (23)$$

as shown in the lecture note, the solution

$$\mathbf{r} = -\eta \cdot \operatorname{sign}(\mathbf{w}) \quad (24)$$

(ii) FGSM attack

$$\begin{aligned} \underset{\mathbf{x}}{\operatorname{argmax}} \quad & J(\mathbf{x}, \mathbf{w}) \\ \text{subject to} \quad & \|\mathbf{x} - \mathbf{x}_0\|_\infty \leq \eta \end{aligned} \quad (25)$$

Approximately, $J(\mathbf{x}, \mathbf{w}) = J(\mathbf{x}_0 + \mathbf{r}, \mathbf{w}) \approx J(\mathbf{x}_0, \mathbf{w}) + \nabla_{\mathbf{x}} J(\mathbf{x}_0, \mathbf{w})^T \mathbf{r}$
Then, the problem becomes

$$\begin{aligned} \underset{\mathbf{x}}{\operatorname{argmin}} \quad & J(\mathbf{x}_0, \mathbf{w}) - \nabla_{\mathbf{x}} J(\mathbf{x}_0, \mathbf{w})^T \mathbf{r} \\ \text{subject to } & \|\mathbf{r}\|_{\infty} \leq \eta \end{aligned} \quad \text{in} \quad (26)$$

of which the solution is given by

$$\mathbf{x} = \mathbf{x}_0 + \eta \cdot (\nabla_{\mathbf{x}} J(\mathbf{x}_0, \mathbf{w})) \quad (27)$$

In the problem setup, we get $J(\mathbf{x}) = -g(\mathbf{x})$, thus the solution is

$$\begin{aligned} \mathbf{x} &= \mathbf{x}_0 - \eta \cdot (\nabla_{\mathbf{x}} g(\mathbf{x}_0)) \\ &= \mathbf{x}_0 - \eta \cdot ((\mathbf{W}_j - \mathbf{W}_t)\mathbf{x}_0 + (\mathbf{w}_j - \mathbf{w}_t)) \\ &= (\eta(\mathbf{W}_t - \mathbf{W}_j) + 1)\mathbf{x}_0 + \eta(\mathbf{w}_t - \mathbf{w}_j) \end{aligned} \quad (28)$$

(iii) I-FGSM attack

$$\begin{aligned} J(\mathbf{x}, \mathbf{w}) &= J(\mathbf{x}_0 + \mathbf{r}, \mathbf{w}) \\ &\approx J(\mathbf{x}_0, \mathbf{w}) + \nabla_{\mathbf{x}} J(\mathbf{x}_0, \mathbf{w})^T \mathbf{r} \\ &= J(\mathbf{x}_0, \mathbf{w}) + \nabla_{\mathbf{x}} J(\mathbf{x}_0, \mathbf{w})^T (\mathbf{x} - \mathbf{x}_0) \\ &= J(\mathbf{x}_0, \mathbf{w}) + \nabla_{\mathbf{x}} J(\mathbf{x}_0, \mathbf{w})^T \mathbf{x} - \nabla_{\mathbf{x}} J(\mathbf{x}_0, \mathbf{w})^T \mathbf{x}_0 \end{aligned} \quad (29)$$

$$\begin{aligned} \mathbf{x}^{(k+1)} &= \underset{0 \leq \mathbf{x} \leq 1}{\operatorname{argmax}} J(\mathbf{x}^{(k)}, \mathbf{w}) \text{ subject to } \|\mathbf{x} - \mathbf{x}_0\| \leq \eta \\ &= \underset{0 \leq \mathbf{x} \leq 1}{\operatorname{argmax}} \nabla_{\mathbf{x}} J(\mathbf{x}^{(k)}, \mathbf{w})^T \mathbf{x} \text{ subject to } \|\mathbf{x} - \mathbf{x}_0\| \leq \eta \\ &= \mathcal{P} \left\{ \mathbf{x}^{(k)} + \eta \cdot \operatorname{sign} \left(\nabla_{\mathbf{x}} J(\mathbf{x}^{(k)}, \mathbf{w}) \right) \right\} \end{aligned} \quad (30)$$

(c) Regularization based attack

(i) linear case

$$\underset{\mathbf{x}}{\operatorname{argmin}} \quad \frac{1}{2} \|\mathbf{x} - \mathbf{x}_0\|_2^2 + \lambda(\mathbf{w}^T \mathbf{x} + w_0) \quad (31)$$

Taking derivative

$$\begin{aligned} \nabla_{\mathbf{x}} \quad & \frac{1}{2} \|\mathbf{x} - \mathbf{x}_0\|_2^2 + \lambda(\mathbf{w}^T \mathbf{x} + w_0) \\ &= \mathbf{x} - \mathbf{x}_0 + \lambda \mathbf{w} = 0 \end{aligned} \quad (32)$$

Solve for \mathbf{x} ,

$$\mathbf{x} = \mathbf{x}_0 - \lambda \mathbf{w} \quad (33)$$

(ii)

$$\begin{aligned} & \underset{\mathbf{x}}{\operatorname{argmin}} \varphi(\mathbf{x}), \text{ where} \\ \varphi(\mathbf{x}) &= \|\mathbf{x} - \mathbf{x}_0\|_2^2 + \lambda \zeta(g_j(\mathbf{x}) - g_t(\mathbf{x})), \\ \zeta(y) &= \max(y, 0), \text{ and } j \neq t \end{aligned} \tag{34}$$

Taking the derivative,

$$\nabla \varphi(\mathbf{x}) = 2(\mathbf{x} - \mathbf{x}_0) + \lambda \mathbb{I}\{g_j(\mathbf{x}) - g_t(\mathbf{x}) > 0\} \cdot (\nabla_{\mathbf{x}} g_j(\mathbf{x}) - \nabla_{\mathbf{x}} g_t(\mathbf{x})) \tag{35}$$

Using my favorite gradient descent, we can tell,

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \alpha \nabla_{\mathbf{x}} \varphi(\mathbf{x}^{(k)}) \tag{36}$$

substituting in $g(\mathbf{x}) = \frac{1}{2} \mathbf{x}^T (\mathbf{W}_j - \mathbf{W}_t) \mathbf{x} + (\mathbf{w}_j - \mathbf{w}_t)^T \mathbf{x} + (w_{j,0} - w_{t,0})$, we get

$$\begin{aligned} \nabla_{\mathbf{x}} g(\mathbf{x}) &= (\mathbf{W}_j - \mathbf{W}_t) \mathbf{x} + (\mathbf{w}_j - \mathbf{w}_t) \\ \mathbf{x}^{(k+1)} &= \mathbf{x}^{(k)} - 2\alpha(\mathbf{x}^{(k)} - \mathbf{x}_0) - \alpha \lambda \mathbb{I}\{g(\mathbf{x}) > 0\} \cdot (\nabla_{\mathbf{x}} g(\mathbf{x})) \\ &= \mathbf{x}^{(k)} - 2\alpha(\mathbf{x}^{(k)} - \mathbf{x}_0) - \alpha \lambda \mathbb{I}\{g(\mathbf{x}) > 0\} \cdot ((\mathbf{W}_j - \mathbf{W}_t) \mathbf{x} + (\mathbf{w}_j - \mathbf{w}_t)) \end{aligned} \tag{37}$$

Exercise 2: CW-Attack on Gaussian Classifier - Non-overlapping Patches

(a)&(b) Implementation

refer to code in the back

(c) Results

(i) & (ii) The final perturbed images & their perturbation respectively

(iii) Frobenius norm of the perturbation

λ	1	5	10
Frobenius norm	8.253352243144604	8.928742803780679	9.697938770929634

Table 1: Frobenius norm vs. λ for CW-Attack on Gaussian classifier with non-overlapping patches

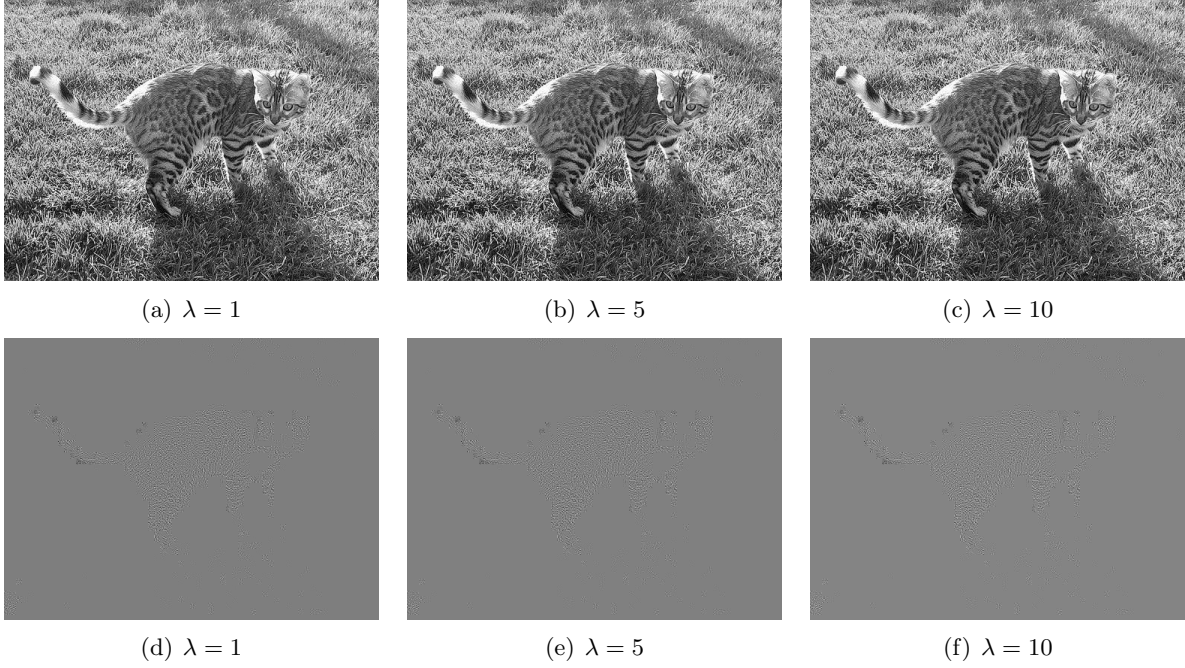


Figure 2: perturbed images and their perturbations with $\lambda = 1, 5, 10$

(iv) The classifiers output

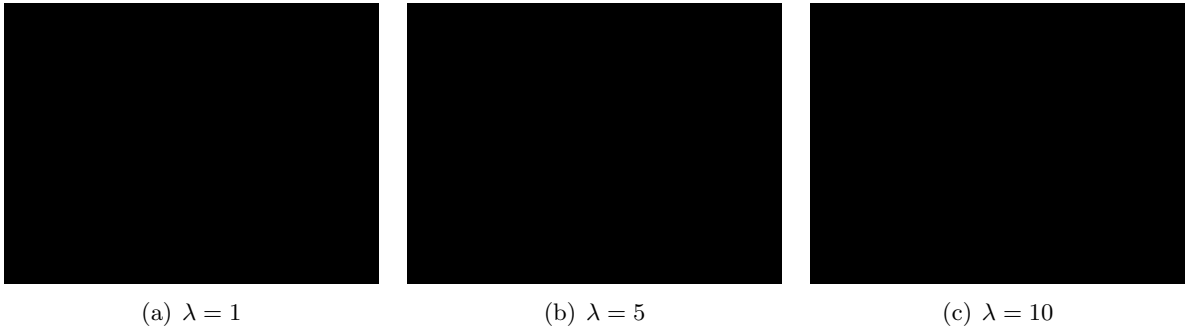


Figure 3: classified perturbed image with $\lambda = 1, 5, 10$

For all three cases, all patches are classified as grass after applying perturbation.

(v) Plots during gradient descent

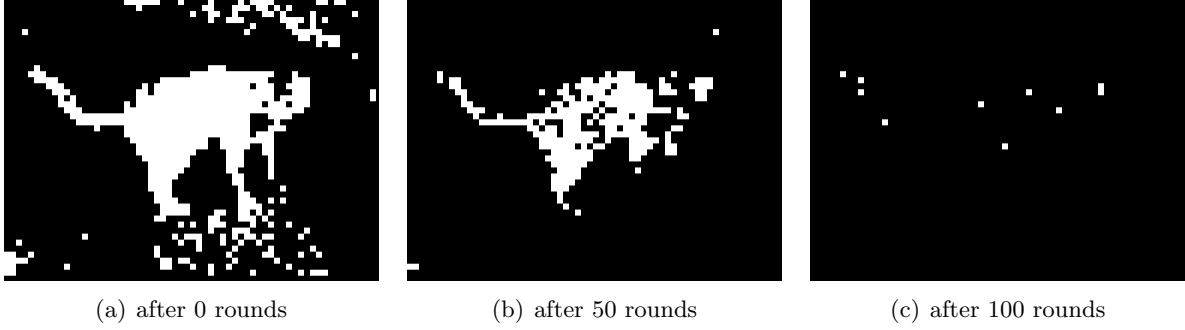


Figure 4: classified perturbed image during gradient descent with $\lambda = 1$

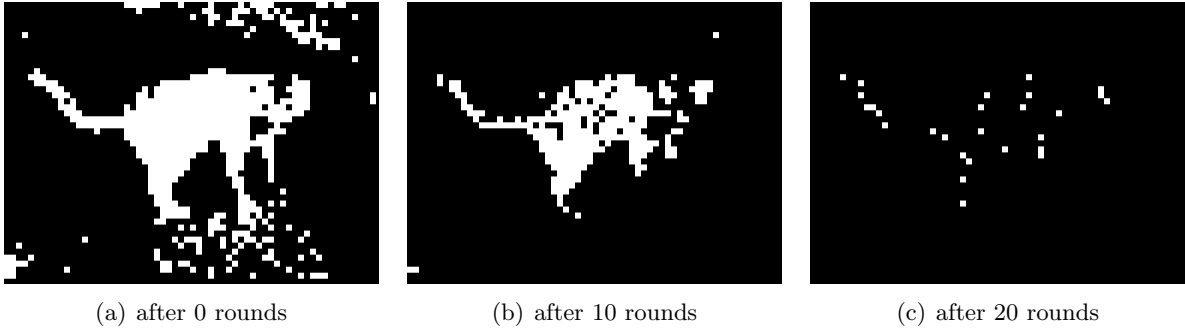


Figure 5: classified perturbed image during gradient descent with $\lambda = 5$

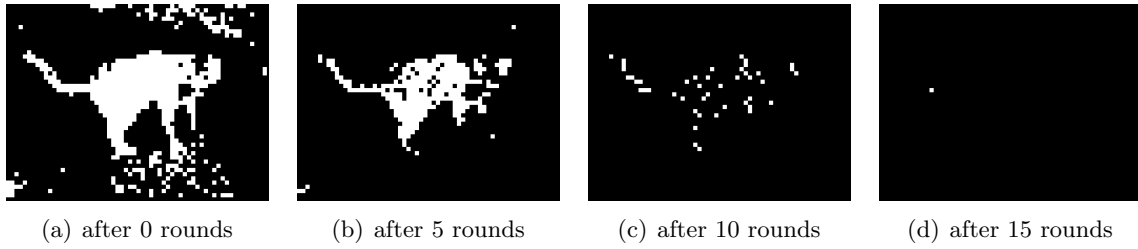


Figure 6: classified perturbed image during gradient descent with $\lambda = 10$

Comments

From the above plots, we can obviously see that increase λ will speed up the attack. However, as λ goes up, the quality of the attack goes down since the Frobenius norm also goes up while achieving the same result.

Exercise 3: CW Attack on Gaussian Classifier - Overlapping Patches

$$\begin{aligned}
 \mathbf{X}^* &= \underset{\mathbf{X}}{\operatorname{argmin}} \sum_{i=1}^L \{ \|\mathbf{P}_i(\mathbf{X} - \mathbf{X}_0)\|_2^2 + \lambda \max(g_j(\mathbf{P}_i\mathbf{X}) - g_t(\mathbf{P}_i\mathbf{X}), 0) \} \\
 &= \underset{\mathbf{X}}{\operatorname{argmin}} \|\mathbf{X} - \mathbf{X}_0\|_2^2 + \lambda \sum_{i=1}^L \max(g_j(\mathbf{P}_i\mathbf{X}) - g_t(\mathbf{P}_i\mathbf{X}), 0)
 \end{aligned} \tag{38}$$

(a) Theory

the gradient of the above function is,

$$\nabla_{\mathbf{X}} = 2(\mathbf{X} - \mathbf{X}_0) + \lambda \mathbb{I}\{g_j(\mathbf{P}_i\mathbf{X}) - g_t(\mathbf{P}_i\mathbf{X}) > 0\} \times (\nabla_{\mathbf{X}}(g_j(\mathbf{P}_i\mathbf{X}) - g_t(\mathbf{P}_i\mathbf{X}))) \tag{39}$$

(b) Implementation

Refer to code in the back.

(c) Results

(i) & (ii) The final perturbed images & their perturbation respectively

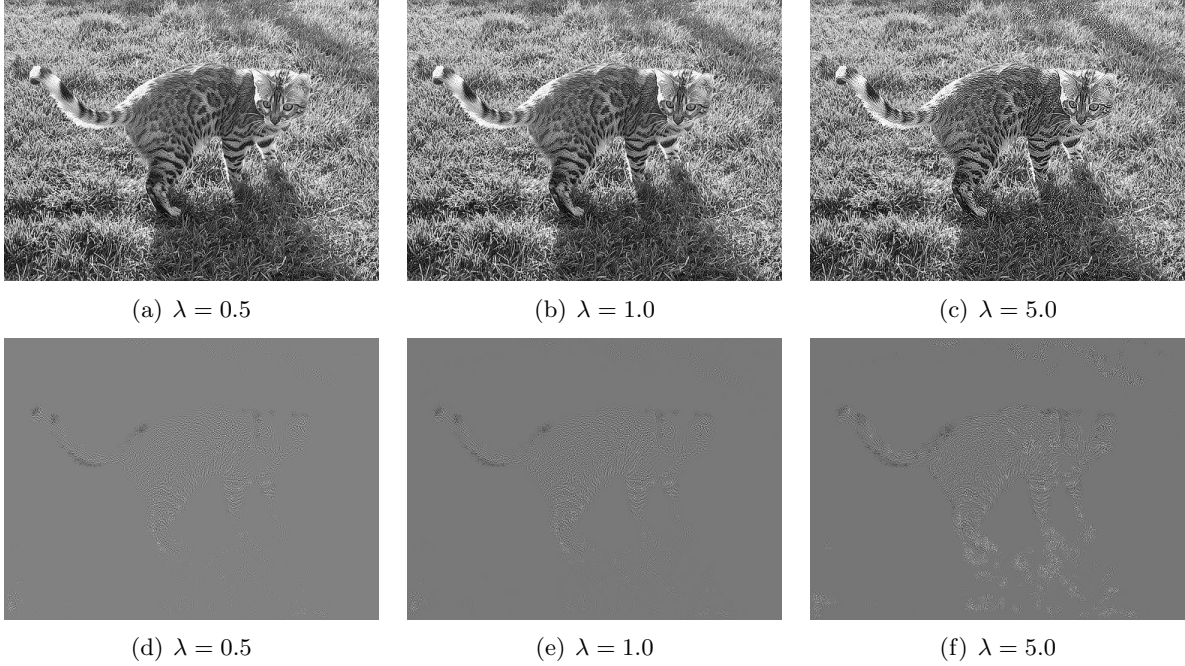


Figure 7: perturbed images and their perturbations with $\lambda = 0.5, 1.0, 5.0$

(iii) Frobenius norm of the perturbation

λ	#grass patches	#cat patches	Frobenius norm
0.5	180561	3	14.059260687697552
1.0	180561	3	16.158563208747122
5.0	180564	0	33.617113348842395

Table 2: testing results for CW-Attack on Gaussian classifier with overlapping patches

(iv) The classifiers output

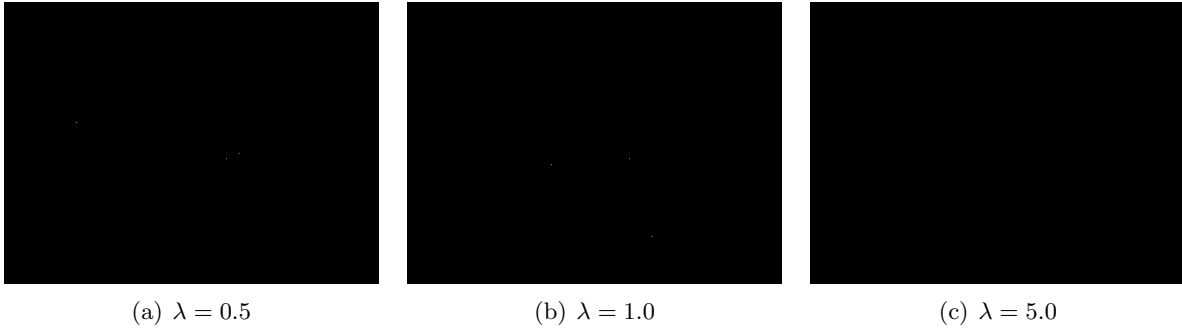


Figure 8: classified perturbed image with $\lambda = 0.5, 1.0, 5.0$

The precise result is in the table above, however, from observation, almost all pixels are classified as grass.

(v) Plots during gradient descent

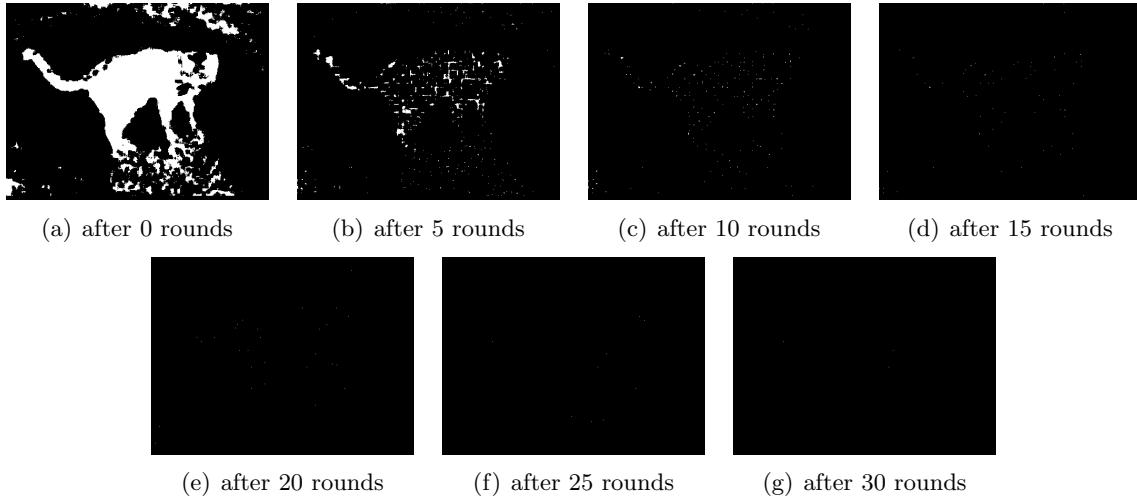


Figure 9: classified perturbed image during gradient descent with $\lambda = 0.5$

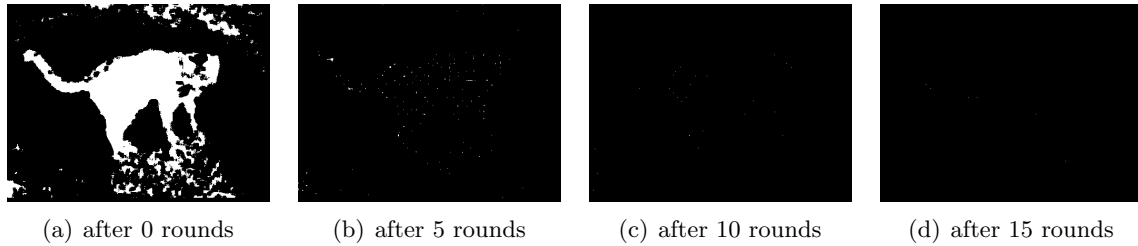


Figure 10: classified perturbed image during gradient descent with $\lambda = 1.0$

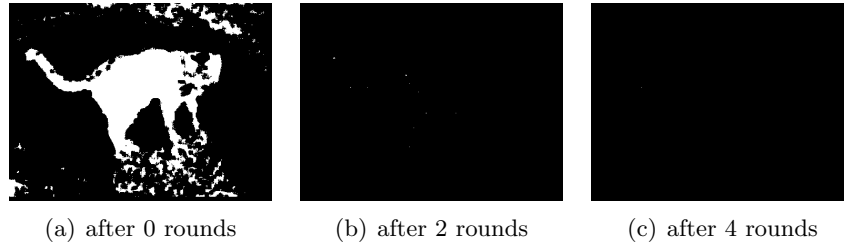


Figure 11: classified perturbed image during gradient descent with $\lambda = 5$

Comments

In general, comparing to the non-overlapping case, the Frobenius norm is significantly larger. However, since the overlapping is a much harder problem, the result is acceptable. As *lambda* goes up, the strength of the attack is approximately the same, less rounds are needed to attack the classifier while sacrificing some quality of the attack.

Code

Exercise 2: CW-Attack on Gaussian Classifier - Non-overlapping Patches

```
#!/usr/bin/env python3
import numpy as np
import scipy.ndimage
from scipy.misc import imread

def readImg(pix):
    return scipy.ndimage.imread(pix, mode='L') / 255

def get_info():
    train_cat = np.matrix(np.loadtxt('../data/train_cat.txt', delimiter=','))
    train_grass = np.matrix(np.loadtxt('../data/train_grass.txt', delimiter=','))

    mu_cat = np.asmatrix(np.mean(train_cat, 1))
    mu_grass = np.asmatrix(np.mean(train_grass, 1))

    cov_cat = np.asmatrix(np.cov(train_cat))
    cov_grass = np.asmatrix(np.cov(train_grass))
```

```

pi_cat      = len(train_cat.T) / (len(train_cat.T) + len(train_grass.T))
pi_grass    = len(train_grass.T) / (len(train_cat.T) + len(train_grass.T))

```

```

return {"mean":mu_cat, "cov":cov_cat, "prior":pi_cat}, {"mean":mu_grass,
↪  "cov":cov_grass, "prior":pi_grass}

```

```

def CWAttack(cat_info, grass_info, img, rounds = 300, alpha = 0.0001, lm = 5):

```

```

    Wcat = np.linalg.inv(cat_info["cov"])
    wcat = -Wcat * cat_info["mean"]
    w0cat = (cat_info["mean"].T * Wcat * cat_info["mean"]) / 2 \
        + np.log(np.linalg.det(cat_info["cov"])) / 2 - np.log(cat_info["prior"])

```

```

    Wgrass = np.linalg.inv(grass_info["cov"])
    wgrass = -Wgrass * grass_info["mean"]
    w0grass = grass_info["mean"].T * Wgrass * grass_info["mean"] / 2 \
        + np.log(np.linalg.det(grass_info["cov"])) / 2 -
    ↪  np.log(grass_info["prior"])

```

```

def __gee(x):

```

```

    # Wt == Wgrass, Wj == Wcat
    ret = x.T * (Wgrass - Wcat) * x / 2 + (wgrass - wcat).T * x + (w0grass -
    ↪  w0cat)
    return ret

```

```

def gradient(z_vector, z_0, lm):

```

```

    # Calculate g_j, g_t, determine if patch_vec is already in target class
    # If patch_vec is in target class, do not add any perturbation (return
    ↪ zero gradient!)
    # Else, calculate the gradient, using results from 1(c)(ii)

```

```

    if __gee(z_vector) > 0:
        return 2 * (z_vector - z_0) + lm * ((Wgrass - Wcat) * z_vector +
    ↪  (wgrass - wcat))

```

```

    else:
        #return 2 * (z_vector - z_0)
        return np.zeros((64,1))

```

```

M,N = img.shape
img_orig = img.copy()
range_i = range(0,M-8,8)

```

```

range_j = range(0,N-8,8)

def classify(img):
    output = np.zeros((M,N))
    for i in range_i:
        for j in range_j:
            z = img[i:i+8, j:j+8]
            z_vector = np.asmatrix(z.flatten('F')).T
            if __gee(z_vector) > 0:
                for ii in range(8):
                    for jj in range(8):
                        output[i+ii][j+jj] = 1
    return output

def cntclass(img):
    num_cat = 0
    num_grass = 0
    for i in range_i:
        for j in range_j:
            z = img[i:i+8, j:j+8]
            z_vector = np.asmatrix(z.flatten('F')).T
            if __gee(z_vector) < 0:
                num_cat+=1
            else:
                num_grass+=1
    return num_grass, num_cat

for r in range(rounds):
    if r % (50 / lm) == 0:
        print(f"round: {r}")
        imsave(f'../pix/CW_Nonoverlap_{lm}_r{r}.png', 255*classify(img))

img_prev = img.copy()
for i in range_i:
    for j in range_j:
        z = img[i:i+8, j:j+8]
        z_vector = np.asmatrix(z.flatten('F')).T
        z_vector_orig = np.asmatrix(img_orig[i:i+8,
            ↪ j:j+8].flatten('F')).T
        z_vector -= alpha * gradient(z_vector, z_vector_orig, lm)
        img[i:i+8, j:j+8] = np.reshape(np.clip(z_vector, 0.0, 1.0),
            ↪ (8,8), order='F')

change = np.linalg.norm(img - img_prev)
if(change < 0.001):

```

```

        print(f"finished, total round: {r}")
        break;

cnts = cntclass(img)
print(f"number of grass patches {cnts[0]}, cat patches {cnts[1]}")
change = np.linalg.norm(img - img_orig)
print(f"Frobenius norm: {change}")

imsave(f'../pix/CW_Nonoverlap_{lm}.png', img)
imsave(f'../pix/CW_Nonoverlap_{lm}_perturbation.png', img-img_orig)
imsave(f'../pix/CW_Nonoverlap_{lm}_classified_pertub.png', classify(img))

if __name__ == "__main__":
    cat_info, grass_info = get_info()
    CWAttack(cat_info, grass_info, readImg('../data/cat_grass.jpg'), lm = 1)
    CWAttack(cat_info, grass_info, readImg('../data/cat_grass.jpg'), lm = 5)
    CWAttack(cat_info, grass_info, readImg('../data/cat_grass.jpg'), lm = 10)

```

Exercise 3: CW-Attack on Gaussian Classifier - Overlapping Patches

```

#!/usr/bin/env python3
import numpy as np
import scipy.ndimage
from scipy.misc import imsave

def readImg(pix):
    return scipy.ndimage.imread(pix, mode='L') / 255

def get_info():
    train_cat = np.matrix(np.loadtxt('../data/train_cat.txt', delimiter=','))
    train_grass = np.matrix(np.loadtxt('../data/train_grass.txt', delimiter=','))

    mu_cat = np.asmatrix(np.mean(train_cat, 1))
    mu_grass = np.asmatrix(np.mean(train_grass, 1))

    cov_cat = np.asmatrix(np.cov(train_cat))
    cov_grass = np.asmatrix(np.cov(train_grass))

    pi_cat = len(train_cat.T) / (len(train_cat.T) + len(train_grass.T))
    pi_grass = len(train_grass.T) / (len(train_cat.T) + len(train_grass.T))

    return {"mean":mu_cat, "cov":cov_cat, "prior":pi_cat}, {"mean":mu_grass,
    ↪ "cov":cov_grass, "prior":pi_grass}

```

```

def CWAttack(cat_info, grass_info, img, rounds = 300, alpha = 0.0001, lm = 0.5):

    Wcat = np.linalg.inv(cat_info["cov"])
    wcat = -Wcat * cat_info["mean"]
    w0cat = (cat_info["mean"].T * Wcat * cat_info["mean"]) / 2 \
        + np.log(np.linalg.det(cat_info["cov"])) / 2 - np.log(cat_info["prior"])

    Wgrass = np.linalg.inv(grass_info["cov"])
    wgrass = -Wgrass * grass_info["mean"]
    w0grass = grass_info["mean"].T * Wgrass * grass_info["mean"] / 2 \
        + np.log(np.linalg.det(grass_info["cov"])) / 2 -
        ↪ np.log(grass_info["prior"])

def __gee(x):

    # Wt == Wgrass, Wj == Wcat
    ret = x.T * (Wgrass - Wcat) * x / 2 + (wgrass - wcat).T * x + (w0grass -
    ↪ w0cat)
    return ret

def gradient(z_vector, lm):
    # Calculate g_j, g_t, determine if patch_vec is already in target class
    # If patch_vec is in target class, do not add any perturbation (return
    ↪ zero gradient!)
    # Else, calculate the gradient, using results from 1(c)(ii)

    if __gee(z_vector) > 0:
        return lm * ((Wgrass - Wcat) * z_vector + (wgrass - wcat))
    else:
        return np.zeros((64,1))

M,N = img.shape
img_orig = img.copy()
range_i = range(0,M-8)
range_j = range(0,N-8)

def classify(img):
    output = np.zeros((M,N))
    for i in range_i:
        for j in range_j:
            z = img[i:i+8, j:j+8]
            z_vector = np.asmatrix(z.flatten('F')).T
            if __gee(z_vector) > 0:
                output[i][j] = 1

```

```

    return output

def cntclass(img):
    num_cat = 0
    num_grass = 0
    for i in range_i:
        for j in range_j:
            z = img[i:i+8, j:j+8]
            z_vector = np.asmatrix(z.flatten('F')).T
            if __gee(z_vector) > 0:
                num_cat+=1
            else:
                num_grass+=1
    return num_grass, num_cat

for r in range(rounds):
    if r % (5 if lm == 0.5 or lm == 1 else 2) == 0:
        print(f"round: {r}")
        imsave(f'../pix/CW_Overlap_{int(lm*10)}_r{r}.png', 255*classify(img))

    img_prev = img.copy()
    grad = np.zeros((M,N))
    for i in range_i:
        for j in range_j:
            z = img[i:i+8, j:j+8]
            z_vector = np.asmatrix(z.flatten('F')).T
            grad[i:i+8, j:j+8] += np.reshape(gradient(z_vector, lm), (8,8),
            ↪ order='F')

    grad += 2 * (img - img_orig)
    img = np.clip(img - alpha * grad, 0.0, 1.0)
    change = np.linalg.norm(img - img_prev)
    if(change < 0.01):
        print(f"finished, total round: {r}")
        break;

cnts = cntclass(img)
print(f"number of grass patches {cnts[0]}, cat patches {cnts[1]}")
change = np.linalg.norm(img - img_orig)
print(f"Frobenius norm: {change}")

imsave(f'../pix/CW_Overlap_{int(lm*10)}.png', img)
imsave(f'../pix/CW_Overlap_{int(lm*10)}_perturbation.png',img-img_orig)
imsave(f'../pix/CW_Overlap_{int(lm*10)}_classified_perturb.png',
↪ classify(img))

```

```
if __name__ == "__main__":  
    cat_info, grass_info = get_info()  
    CWAttack(cat_info, grass_info, readImg('../data/cat_grass.jpg'), lm = 0.5)  
    CWAttack(cat_info, grass_info, readImg('../data/cat_grass.jpg'), lm = 1.0)  
    CWAttack(cat_info, grass_info, readImg('../data/cat_grass.jpg'), lm = 5)
```