

ECE 595: Final project(Spring 2020)
Frank Lin (ID:54)

Exercise 1

(a) We study the minimum-distance attack

(i) Consider general l_2 -distance attack with linear decision boundary as the problem

$$\underset{\mathbf{x}}{\operatorname{argmin}} \|\mathbf{x} - \mathbf{x}_0\|_2^2 \quad \text{subject to} \quad \mathbf{w}^T \mathbf{x} + w_0 = 0$$

by Lagrange multiplier such that the solution is

$$\mathbf{x}^* = \mathbf{x}_0 - \left(\frac{\mathbf{w}^T \mathbf{x}_0 + w_0}{\|\mathbf{w}\|_2^2} \right) \mathbf{w} \quad (1)$$

we obtain

$$\begin{aligned} \mathbf{x}^* &= \mathbf{x}_0 - \left(\frac{\mathbf{w}^T \mathbf{x}_0 + w_0}{\|\mathbf{w}\|_2^2} \right) \mathbf{w} \\ &= \mathbf{x}_0 - \left(\frac{(\boldsymbol{\mu}_j - \boldsymbol{\mu}_t)^T \boldsymbol{\Sigma}^{-1} \mathbf{x}_0 - \frac{1}{2}(\boldsymbol{\mu}_j - \boldsymbol{\mu}_t)^T \boldsymbol{\Sigma}^{-1}(\boldsymbol{\mu}_j + \boldsymbol{\mu}_t) + \log(\pi_j/\pi_t)}{\|\boldsymbol{\Sigma}^{-1}(\boldsymbol{\mu}_j - \boldsymbol{\mu}_t)\|_2^2} \right) \boldsymbol{\Sigma}^{-1}(\boldsymbol{\mu}_j - \boldsymbol{\mu}_t) \\ &= \mathbf{x}_0 - \left(\frac{(\boldsymbol{\mu}_j - \boldsymbol{\mu}_t)^T \boldsymbol{\Sigma}^{-1}[\mathbf{x}_0 - \frac{1}{2}(\boldsymbol{\mu}_j + \boldsymbol{\mu}_t)] + \log(\pi_j/\pi_t)}{\|\boldsymbol{\Sigma}^{-1}(\boldsymbol{\mu}_j - \boldsymbol{\mu}_t)\|_2^2} \right) \boldsymbol{\Sigma}^{-1}(\boldsymbol{\mu}_j - \boldsymbol{\mu}_t). \end{aligned}$$

Consider general l_∞ -distance attack with linear decision boundary as the problem

$$\underset{\mathbf{x}}{\operatorname{argmin}} \|\mathbf{x} - \mathbf{x}_0\|_\infty \quad \text{subject to} \quad \mathbf{w}^T \mathbf{x} + w_0 = 0$$

we let $\mathbf{r} = \mathbf{x} - \mathbf{x}_0$ and $b_0 = -(\mathbf{w}^T \mathbf{x}_0 + w_0)$ to form the equivalent problem: minimize $\|\mathbf{r}\|_\infty$ subject to $\mathbf{w}^T \mathbf{r} = b_0$. We then apply Holder's inequality to have $|b_0| = |\mathbf{w}^T \mathbf{r}| \leq \|\mathbf{w}\|_1 \|\mathbf{r}\|_\infty$ such that the optimal solution is obtained as the equality holds, i.e.,

$$\mathbf{r}^* = \frac{b_0}{\|\mathbf{w}\|_1} \operatorname{sign}(\mathbf{w}) \quad \Rightarrow \quad \mathbf{x}^* = \mathbf{x}_0 - \left(\frac{\mathbf{w}^T \mathbf{x}_0 + w_0}{\|\mathbf{w}\|_1} \right) \operatorname{sign}(\mathbf{w}). \quad (2)$$

We substitute into (2) to obtain

$$\begin{aligned} \mathbf{x}^* &= \mathbf{x}_0 - \left(\frac{\mathbf{w}^T \mathbf{x}_0 + w_0}{\|\mathbf{w}\|_1} \right) \operatorname{sign}(\mathbf{w}) \\ &= \mathbf{x}_0 - \left(\frac{(\boldsymbol{\mu}_j - \boldsymbol{\mu}_t)^T \boldsymbol{\Sigma}^{-1}[\mathbf{x}_0 - \frac{1}{2}(\boldsymbol{\mu}_j + \boldsymbol{\mu}_t)] + \log(\pi_j/\pi_t)}{\|\boldsymbol{\Sigma}^{-1}(\boldsymbol{\mu}_j - \boldsymbol{\mu}_t)\|_1} \right) \operatorname{sign}(\boldsymbol{\Sigma}^{-1}(\boldsymbol{\mu}_j - \boldsymbol{\mu}_t)). \end{aligned}$$

(ii) Given the quadratic decision boundary $g(\mathbf{x})$, we have the gradient

$$\nabla g(\mathbf{x}) = (\mathbf{W}_j - \mathbf{W}_t)\mathbf{x} + (\mathbf{w}_j - \mathbf{w}_t)$$

such that the iterates of solving the DeepFool attack is

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \frac{g(\mathbf{x}^{(k)})}{\|\nabla g(\mathbf{x}^{(k)})\|_2^2} \nabla g(\mathbf{x}^{(k)}). \quad (3)$$

- (iii) Consider a decision boundary $g(x) = 1/x$ such that $dg(x)/dx = -1/x^2$. If we apply the decision boundary to the iterates of DeepFool attack, we have

$$x^{(k+1)} = x^{(k)} - \frac{1/x^{(k)}}{1/(x^{(k)})^4} \left(-\frac{1}{(x^{(k)})^2} \right) = 2x^{(k)},$$

which obviously never converge. The intuitive reason is that the decision boundary has zero $g(x) = 1/x = 0$ when $x \rightarrow \infty$.

- (b) Then we study the maximum-loss attacks

- (i) Consider the maximum loss l_∞ attack of linear case as

$$\operatorname{argmin}_{\mathbf{x}} \mathbf{w}^T \mathbf{x} + w_0 \quad \text{subject to} \quad \|\mathbf{x} - \mathbf{x}_0\|_\infty < \eta. \quad (4)$$

Let $\mathbf{x} = \mathbf{x}_0 + \mathbf{r}$ and $b_0 = \mathbf{w}^T \mathbf{x}_0 + w_0$. Then the original problem (4) becomes

$$\operatorname{argmin}_{\mathbf{r}} \mathbf{w}^T \mathbf{x} + w_0 = \mathbf{w}^T (\mathbf{x}_0 + \mathbf{r}) + w_0 = \mathbf{w}^T \mathbf{r} + b_0 \quad \text{subject to} \quad \|\mathbf{r}\|_\infty < \eta. \quad (5)$$

We apply the Holder's inequality

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

where the equality holds when $\mathbf{r} = -\eta \cdot \operatorname{sign}(\mathbf{w})$ such that the optimal solution of (5) is $\mathbf{r} = -\eta \cdot \operatorname{sign}(\mathbf{w})$ and hence the optimal solution of (4) is $\mathbf{x} = \mathbf{x}_0 - \eta \operatorname{sign}(\mathbf{w})$. Replacing \mathbf{w} , we then have $\mathbf{x} = \mathbf{x}_0 - \eta \operatorname{sign}(\Sigma^{-1}(\boldsymbol{\mu}_j - \boldsymbol{\mu}_t))$.

- (ii) Consider the FGSM attack with Gaussian decision boundary $J(\mathbf{x}) = -g(\mathbf{x})$, we have $\nabla_{\mathbf{x}} J(\mathbf{x}) = -(\mathbf{W}_j - \mathbf{W}_t)\mathbf{x} - (\mathbf{w}_j - \mathbf{w}_t)$ and hence formula becomes

$$\mathbf{x}^* = \mathbf{x}_0 + \eta \operatorname{sign}(\nabla_{\mathbf{x}} J(\mathbf{x}_0)) = \mathbf{x}_0 - \eta \operatorname{sign}((\mathbf{W}_j - \mathbf{W}_t)\mathbf{x}_0 + (\mathbf{w}_j - \mathbf{w}_t)).$$

- (iii) Consider the FGSM attack with Gaussian decision boundary $J(\mathbf{x}) = -g(\mathbf{x})$, we have $\nabla_{\mathbf{x}} J(\mathbf{x}) = -(\mathbf{W}_j - \mathbf{W}_t)\mathbf{x} - (\mathbf{w}_j - \mathbf{w}_t)$ and hence formula becomes

$$\begin{aligned} \mathbf{x}^{(k+1)} &= \mathcal{P}_{[0,1]} \left\{ \mathbf{x}_0 + \eta \operatorname{sign} \left(\nabla_{\mathbf{x}} J(\mathbf{x}^{(k)}) \right) \right\} \\ &= \mathcal{P}_{[0,1]} \left\{ \mathbf{x}_0 + \eta \operatorname{sign} \left((\mathbf{W}_j - \mathbf{W}_t)\mathbf{x}^{(k)} + (\mathbf{w}_j - \mathbf{w}_t) \right) \right\} \end{aligned}$$

- (c) We finally examine the regularization based attacks

- (i) Consider the regularization-based attack in the linear case as

$$\operatorname{argmin}_{\mathbf{x}} \frac{1}{2} \|\mathbf{x} - \mathbf{x}_0\|_2^2 + \lambda(\mathbf{w}^T \mathbf{x} + w_0), \quad (6)$$

we apply the first-order necessary condition to have

$$\frac{\partial}{\partial \mathbf{x}} \frac{1}{2} \|\mathbf{x} - \mathbf{x}_0\|_2^2 + \lambda(\mathbf{w}^T \mathbf{x} + w_0) = \mathbf{x} - \mathbf{x}_0 + \lambda \mathbf{w} = \mathbf{0}$$

such that the optimal solution is

$$\mathbf{x}^* = \mathbf{x}_0 - \lambda \mathbf{w} = \mathbf{x}_0 - \lambda \Sigma^{-1}(\boldsymbol{\mu}_j - \boldsymbol{\mu}_t).$$

(ii) Consider the Carlini-Wagner attack with Gaussian decision boundary as

$$\operatorname{argmin}_{\mathbf{x}} \varphi(\mathbf{x}) \quad \varphi(\mathbf{x}) = \|\mathbf{x} - \mathbf{x}_0\|_2^2 + \lambda \cdot \max(g_j(\mathbf{x}) - g_t(\mathbf{x}), 0) = \|\mathbf{x} - \mathbf{x}_0\|_2^2 + \lambda \cdot \max(-g(\mathbf{x}), 0).$$

where $g(\mathbf{x}) = g_t(\mathbf{x}) - g_j(\mathbf{x})$. Since $\nabla_{\mathbf{x}} g(\mathbf{x}) = (\mathbf{W}_j - \mathbf{W}_t)\mathbf{x} + (\mathbf{w}_j - \mathbf{w}_t)$, we have the gradient

$$\begin{aligned} \nabla_{\mathbf{x}} \varphi(\mathbf{x}) &= 2(\mathbf{x} - \mathbf{x}_0) - \lambda \cdot \mathbb{I}_{g(\mathbf{x}) < 0} \cdot \nabla_{\mathbf{x}} g(\mathbf{x}) \\ &= 2(\mathbf{x} - \mathbf{x}_0) - \lambda \cdot \mathbb{I}_{g(\mathbf{x}) < 0} \cdot [(\mathbf{W}_j - \mathbf{W}_t)\mathbf{x} + (\mathbf{w}_j - \mathbf{w}_t)] \\ &= 2(\mathbf{x} - \mathbf{x}_0) - \lambda \cdot \mathbb{I}_{g(\mathbf{x}) < 0} \cdot \left[\Sigma_j^{-1}(\mathbf{x} - \boldsymbol{\mu}_j) - \Sigma_t^{-1}(\mathbf{x} - \boldsymbol{\mu}_t) \right]. \end{aligned}$$

Therefore the gradient descent iterate is given by

$$\begin{aligned} \mathbf{x}^{(k+1)} &= \mathbf{x}^{(k)} - \alpha_k \nabla_{\mathbf{x}} \varphi(\mathbf{x}^{(k)}) \\ &= \mathbf{x}^{(k)} - 2\alpha_k(\mathbf{x}^{(k)} - \mathbf{x}_0) + \alpha_k \lambda \cdot \mathbb{I}_{g(\mathbf{x}^{(k)}) < 0} \cdot \left[\Sigma_j^{-1}(\mathbf{x}^{(k)} - \boldsymbol{\mu}_j) - \Sigma_t^{-1}(\mathbf{x}^{(k)} - \boldsymbol{\mu}_t) \right] \end{aligned} \tag{7}$$

Exercise 2

- (a) We solve the CW optimization problem by gradient descent obtained in (7) with the regularization coefficient $\lambda = 5$.

```
def gradient_nonoverlapping(pert_vec, img_vec, target_index, lam):
    # Calculate g_j, g_t, determine if patch_vec is already in target class
    grad = 0
    tc = 1-2*target_index
    tg = 2*target_index-1
    g_cat = dsc_fun(pert_vec, mean_cat, inv_cat, res_cat)
    g_grass = dsc_fun(pert_vec, mean_grass, inv_grass, res_grass)
    if (tc*g_cat+tg*g_grass > 0):
        grad = 2*(pert_vec-img_vec)-lam*(tc*inv_cat*(pert_vec-mean_cat)+tg*inv_grass*(pert_vec-mean_grass))
    return grad
```

- (b) We implement the gradient descent algorithm on every non-overlapping patch in the image. The step size is set as $\alpha = 0.0001$ and the stopping criterion is the maximum number of iterations 300 or the norm $\|\mathbf{X}^{(k+1)} - \mathbf{X}^{(k)}\|_2 < 0.001$.

```
def CW_attack_nonoverlapping(img_matrix, target_index, alpha, lam):
    # Start gradient descent
    while itr_num <= 300 and change >= 0.001:
        pert_img_prev = pert_img_curr.copy()
        for i in range(M):
            for j in range(N):
                if i + 8 > M or j + 8 > N or i % 8 != 0 or j % 8 != 0:
                    continue
                # Perturb the patch with top left pixel (i,j) in img_matrix by
                x = img_matrix[i:i+8, j:j+8]
                z = pert_img_prev[i:i+8, j:j+8]
```

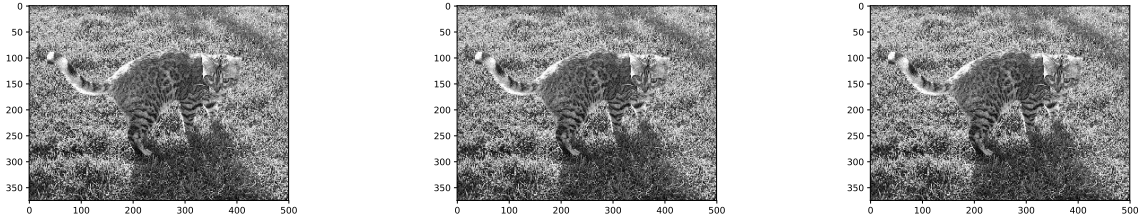
```

img_vec = x.flatten('F').reshape((d,1))
pert_vec = z.flatten('F').reshape((d,1))
grad = gradient_nonoverlapping(pert_vec, img_vec, target_index, lam)
pert_vec = np.clip(pert_vec - alpha*grad, 0.0, 1.0)
pert_img_curr[i:i+8, j:j+8] = np.reshape(pert_vec, (8, 8), order = 'F')
# Process the perturbed image matrix, display perturbed image, etc.
change = np.linalg.norm(pert_img_curr - pert_img_prev)
itr_num += 1
print(itr_num)
return pert_img_curr

```

(c) Consider $\lambda = 1, 5, 10$ we implement the attack as follows.

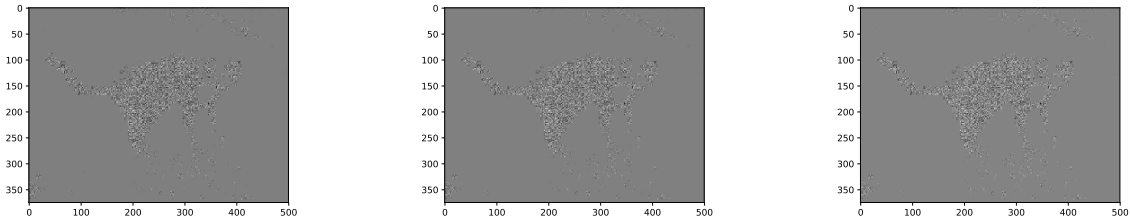
- (i) The final perturbed images for $\lambda = 1, 5, 10$ and target index 0 are given by Fig. 1. The perturbed image looks very similar to the original image. It is hard to point out the differences by eyes.



(a) The perturbed image for $\lambda = 1$. (b) The perturbed image for $\lambda = 5$. (c) The perturbed image for $\lambda = 10$.

Figure 1: The final perturbed images for $\lambda = 1, 5, 10$, change cat to grass.

- (ii) The perturbations added to the original images for $\lambda = 1, 5, 10$, change cat to grass are given in Fig. 2. It can be observed that if we aim to change cat into grass, we add perturbation on the cat part of the classifier's output of the original image.



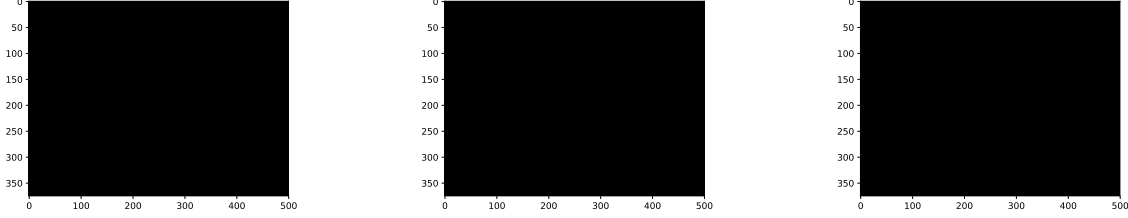
(a) Added perturbation for $\lambda = 1$. (b) Added perturbation for $\lambda = 5$. (c) Added perturbation for $\lambda = 10$.

Figure 2: Added perturbation for $\lambda = 1, 5, 10$ change cat to grass.

- (iii) The Frobenius norm of the perturbations for $\lambda = 1, 5, 10$, change cat to grass are given by Table 1. We can observe that when λ increases, we allow more distance (Frobenius norm) between the perturbed image and original image.
- (iv) The classifier's output on the attack for $\lambda = 1, 5, 10$, change cat to grass are given in Fig. 3. It can be observed that after sufficient large (≤ 300) iterates, all the patches ($46 \times 62 = 2852$) of the perturbed image will be classified as grass.

Table 1: The Frobenius norm of the perturbations for $\lambda = 1, 5, 10$, change cat to grass.

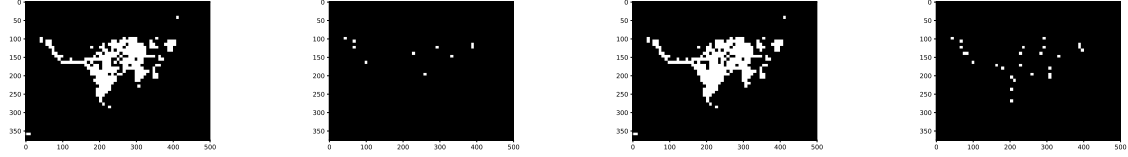
	$\lambda = 1$	$\lambda = 5$	$\lambda = 10$
cat to grass	8.2470	8.9304	9.6988



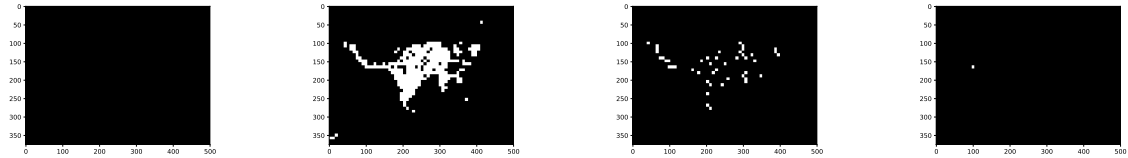
(a) The classifier's output on the at- (b) The classifier's output on the at- (c) The classifier's output on the at-
 tack for $\lambda = 1$. Total $46 \times 62 = 2852$ tack for $\lambda = 5$. Total $46 \times 62 = 2852$ tack for $\lambda = 10$. Total $46 \times 62 = 2852$
 patches classified as grass. patches classified as grass. patches classified as grass.

Figure 3: The classifier's output on the attack for $\lambda = 1, 5, 10$, when change cat to grass.

- (v) The classifier's output on the perturbed images iterates per 50, 10, 5 iterations for $\lambda = 1, 5, 10$ when changing cat to grass are given by Fig. 4. We can observe that during the gradient descend algorithm, as the number of iterates increases, more and more patches of the perturbed image will be classified from cat to grass.



(a) The classifier's output (b) The classifier's output (c) The classifier's output (d) The classifier's output
 on the perturbed images of on the perturbed images of on the perturbed images of on the perturbed images of
 50 iterates for $\lambda = 1$. 100 iterates for $\lambda = 1$. 10 iterates for $\lambda = 5$. 20 iterates for $\lambda = 5$.



(e) The classifier's output (f) The classifier's output (g) The classifier's output (h) The classifier's output
 on the perturbed images of on the perturbed images of on the perturbed images of on the perturbed images of
 30 iterates for $\lambda = 5$. 5 iterates for $\lambda = 10$. 10 iterates for $\lambda = 10$. 15 iterates for $\lambda = 10$ and
 target index 0 (to black).

Figure 4: The classifier's output on the perturbed images iterates per 50, 10, 5 iterations for $\lambda = 1, 5, 10$, change cat to grass.

- (d) When we decrease the step size α , we search from a smaller neighborhood of the original image such that we have smaller initial change but a lower convergence rate. With a smaller step size, the algorithm is more likely to converge.

Exercise 3

- (a) Let $h(\mathbf{P}_i \mathbf{X}) = g_j(\mathbf{P}_i \mathbf{X}) - g_t(\mathbf{P}_i \mathbf{X})$ such that the gradient of $h(\mathbf{P}_i \mathbf{X})$ with respect to \mathbf{X} is

$$\begin{aligned}\nabla_{\mathbf{X}} h(\mathbf{P}_i \mathbf{X}) &= \mathbf{P}_i^T \nabla_{\mathbf{x}} h(\mathbf{x})|_{\mathbf{x}=\mathbf{P}_i \mathbf{X}} = \mathbf{P}_i^T (\nabla_{\mathbf{x}} g_j(\mathbf{x}) - \nabla_{\mathbf{x}} g_t(\mathbf{x}))|_{\mathbf{x}=\mathbf{P}_i \mathbf{X}} \\ &= \mathbf{P}_i^T [(\mathbf{W}_t - \mathbf{W}_j) \mathbf{P}_i \mathbf{X} + (\mathbf{w}_t - \mathbf{w}_j)] = \mathbf{P}_i^T \left[\Sigma_t^{-1} (\mathbf{P}_i \mathbf{X} - \boldsymbol{\mu}_t) - \Sigma_j^{-1} (\mathbf{P}_i \mathbf{X} - \boldsymbol{\mu}_j) \right].\end{aligned}$$

Therefore the gradient with respect to \mathbf{X} is

$$\nabla_{\mathbf{X}} \varphi(\mathbf{X}) = 2(\mathbf{X} - \mathbf{X}_0) + \lambda \sum_{i=1}^L \mathbb{I}_{h(\mathbf{P}_i \mathbf{X}) > 0} \cdot \mathbf{P}_i^T \left[\Sigma_t^{-1} (\mathbf{P}_i \mathbf{X} - \boldsymbol{\mu}_t) - \Sigma_j^{-1} (\mathbf{P}_i \mathbf{X} - \boldsymbol{\mu}_j) \right]$$

- (b) We implement the gradient descent algorithm on the overlapping patch classifier with the iterates

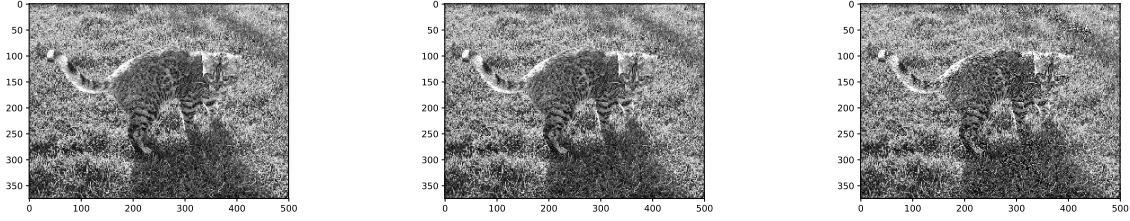
$$\mathbf{X}^{(k+1)} = \mathbf{X}^{(k)} - \alpha \nabla_{\mathbf{X}} \varphi(\mathbf{X}^{(k)}) = \mathbf{X}^{(k)} - 2\alpha(\mathbf{X}^{(k)} - \mathbf{X}_0) - \alpha\lambda \sum_{i=1}^L \mathbb{I}_{h(\mathbf{P}_i \mathbf{X}^{(k)}) > 0} \cdot \nabla_{\mathbf{X}} h(\mathbf{P}_i \mathbf{X}^{(k)})$$

The stopping criterion is the maximum number of iterations 300 or the norm $\|\mathbf{X}^{(k+1)} - \mathbf{X}^{(k)}\|_2 < 0.01$.

```
def CW_attack_overlapping(img_matrix, target_index, alpha, lam):
    # Start gradient descent
    while itr_num <= 300 and change >= 0.01:
        pert_img_prev = pert_img_curr.copy()
        pert_img_curr = pert_img_curr - 2*alpha*(pert_img_prev - img_matrix)
        for i in range(M-8):
            for j in range(N-8):
                # Perturb the patch with top left pixel (i,j) in img_matrix by
                # pert_patch = patch_vec + alpha*grad
                z = pert_img_prev[i:i+8, j:j+8]
                pert_vec = z.flatten('F').reshape((d,1))
                grad = gradient_overlapping(pert_vec, target_index)
                pert_img_curr[i:i+8, j:j+8] += np.reshape(-alpha*lam*grad, (8, 8), order =
                # Process the perturbed image matrix, display perturbed image, etc.
                pert_img_curr = np.clip(pert_img_curr, 0.0, 1.0)
                change = np.linalg.norm(pert_img_curr - pert_img_prev)
                itr_num += 1
    return pert_img_curr
```

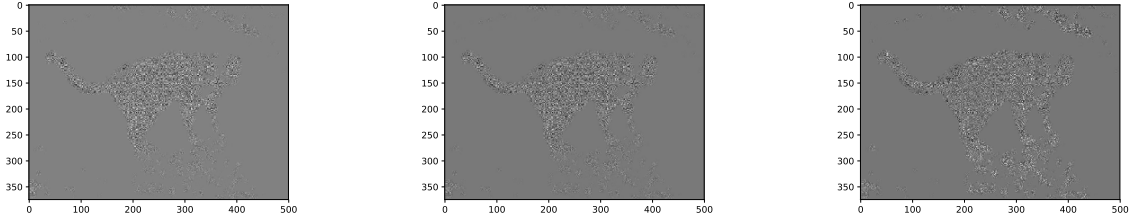
- (c) Consider $\lambda = 0.5, 1, 5$ we implement the attack as follows.

- (i) The final perturbed images for $\lambda = 0.5, 1, 5$ and target index 0 are given by Fig. 5. The perturbed image looks very similar to the original image.
- (ii) The perturbations added to the original images for $\lambda = 0.5, 1, 5$ when changing cat to grass are given in Fig. 6.
- (iii) The Frobenius norm of the perturbations for $\lambda = 0.5, 1, 5$ when changing cat to grass are 14.0593, 16.1586, and 33.6171, respectively. We can observe that when λ increases, we allow more distance (Frobenius norm) between the perturbed image and original image.



(a) The perturbed image for $\lambda = 0.5$. (b) The perturbed image for $\lambda = 1$. (c) The perturbed image for $\lambda = 5$.

Figure 5: The final perturbed images for $\lambda = 0.5, 1, 5$, change cat to grass.



(a) Added perturbation for $\lambda = 0.5$. (b) Added perturbation for $\lambda = 1$. (c) Added perturbation for $\lambda = 5$.

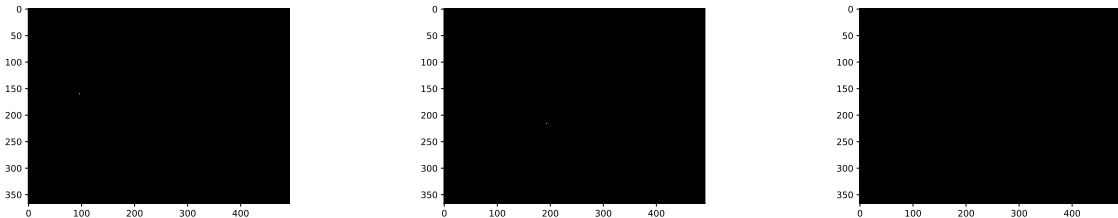
Figure 6: Added perturbation for $\lambda = 0.5, 1, 5$, change cat to grass.

- (iv) The classifier's output on the attack for $\lambda = 0.5, 1, 5$ when changing cat to grass are given in Fig. 7. All the $367 \times 492 = 180564$ patches of the perturbed image will be classified as grass.
- (v) The classifier's output on the perturbed images iterates per 5, 5, 2 iterations for $\lambda = 0.5, 1, 5$ when changing cat to grass are given by Fig. 8, Fig. 9, and Fig. 10, respectively. As iterations continues, more patches of the perturbed image will be classified to grass.

Exercise 4

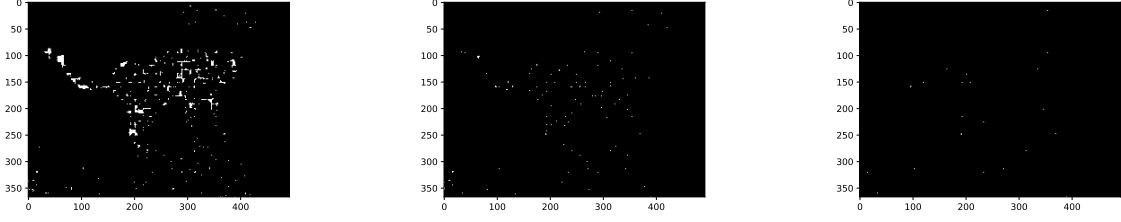
• Method Description

To prevent attackers from launching an adversarial attack to effect the performance of an classifier, we develop a training data preprocessing method to strengthen the data robustness against any malicious data modifications. Different from previous works, instead of trying to

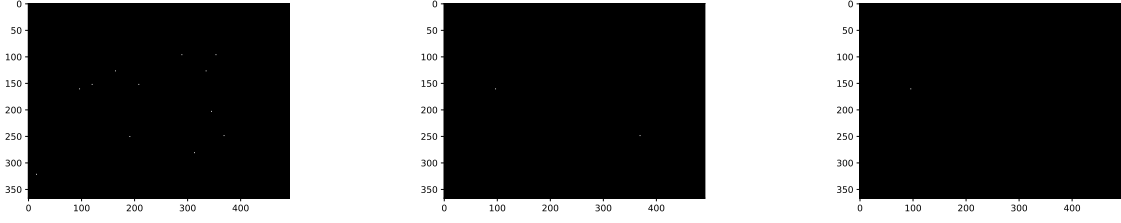


(a) The classifier's output on the at- (b) The classifier's output on the at- (c) The classifier's output on the at-
tack for $\lambda = 0.5$. tack for $\lambda = 1$. tack for $\lambda = 5$.

Figure 7: The classifier's output on the attack for $\lambda = 0.5, 1, 5$.

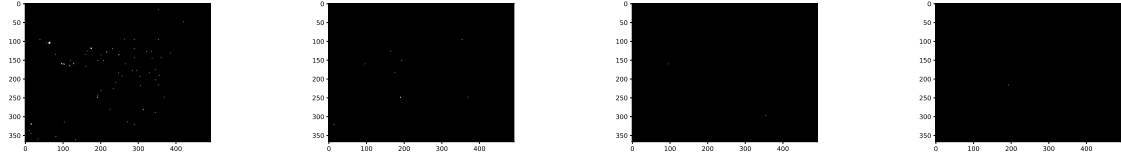


(a) The classifier's output on the per- (b) The classifier's output on the per- (c) The classifier's output on the per-
turbed images of 5 iterates for $\lambda = 0.5$. turbed images of 10 iterates for $\lambda = 0.5$. turbed images of 15 iterates for $\lambda = 0.5$.



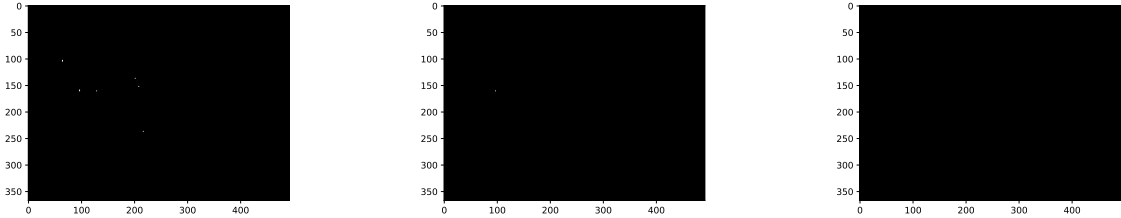
(d) The classifier's output on the per- (e) The classifier's output on the per- (f) The classifier's output on the per-
turbed images of 20 iterates for $\lambda = 0.5$. turbed images of 25 iterates for $\lambda = 0.5$. turbed images of 30 iterates for $\lambda = 0.5$.

Figure 8: The classifier's output on the perturbed images per 5 iterates for $\lambda = 0.5$, change cat to grass.



(a) The classifier's output (b) The classifier's output (c) The classifier's output (d) The classifier's output
on the perturbed images of on the perturbed images of on the perturbed images of on the perturbed images of
5 iterates for $\lambda = 1$. 10 iterates for $\lambda = 1$. 15 iterates for $\lambda = 1$. 20 iterates for $\lambda = 1$.

Figure 9: The classifier's output on the perturbed images per 5 iterates for $\lambda = 1$, change cat to grass.



(a) The classifier's output on the per- (b) The classifier's output on the per- (c) The classifier's output on the per-
turbed images of 2 iterates for $\lambda = 5$. turbed images of 4 iterates for $\lambda = 5$. turbed images of 6 iterates for $\lambda = 5$.

Figure 10: The classifier's output on the perturbed images per 5 iterates for $\lambda = 1$, change cat to grass.

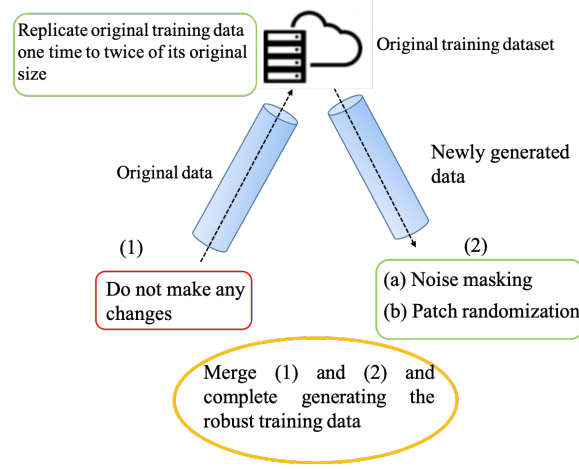


Figure 11: The noise patch-randomization scheme.

recover the damaged data after an attack, we aim to explore a way to effectively intensify the robustness of training data such that the effect of an attack on learning performance is minimized before one is launched.

Consider that the CW attack is a soft-version of the minimum-distance attack, it inherits the idea that, given an input data point x_0 , the attackers want to send x_0 to a targeted class \mathcal{C}_t by heading along a specific direction towards and across the decision boundary. In order to effectively prevent this attack, we want to enlarge the minimum distance between x_0 and any point x in the targeted class \mathcal{C}_t so that the probability for x_0 to be moved to \mathcal{C}_t can be guaranteed to decrease. Therefore, we need to investigate the characteristics of the decision boundary for a Gaussian classifier.

Since we know that the decision boundary of a Gaussian classifier is nonlinear when the covariance matrix of data in different groups is nonequivalent, we can expect that the decision boundary becomes more similar to a linear decision boundary and moves away from any input x_0 as the amount of training data becomes larger. One of the way to do so is by image augmentation. However, it is very important that we preserve the features of data as we increase the training samples. Therefore, we developed a two-step noise patch-randomization scheme to effectively augment the training image as depicted in Fig.11. In the first step, a small randomness of zero mean is added to every pixel of the image. In the second step, a larger but feature-preservative randomness is added to every patch in the image. We first replicate the original data samples (e.g. cat and grass) one time then preform the noise patch-randomization scheme on the newly generated data after replication as follows:

(a) **Noise masking**

A zero-mean Gaussian white noise with a standard deviation of 0.2 is added to the image patches for each epoch, so that each epoch uses a slightly different data set.

(b) **Patch randomization**

Every pixel in the same patch is randomly re-ordered. By doing so, we preserves the feature of the data while provide a large randomness to newly generated training data samples.

• **Comparisons**

To show that our scheme provides a more robust training dataset, we train the classifier

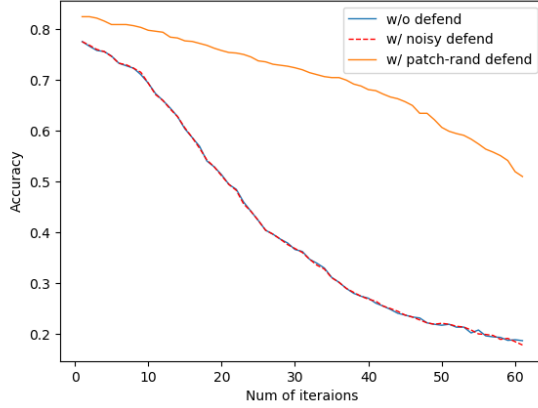


Figure 12: Accuracy w.r.t various defending strategies against CW attack.

(i) with dataset generated by the noise patch-randomization scheme and compare it with the classifier (ii) trained with the original dataset. After showing the robustness of our generated dataset against the original data, we want to further demonstrate the power of patch-randomization. Therefore, we train a classifier (iii) with dataset preprocessed only with noise masking and compare it with classifier (i), which dataset are preprocessed with both noise masking and patch randomization. We compare the performance of the schemes based on accuracy.

The result is demonstrated in Fig.12 with $\alpha = 0.001$ and $\lambda = 0.007$. There are two observations worth noticing. First, it can be observed that the accuracy of our proposed scheme is higher than the other two before the attack is launched. This shows that the patch randomization technique efficiently enlarges the data variability so that the classifier trained with these data has better performance. Second, the decrease rate of accuracy of our proposed scheme is much lower (approximately one-half) than the other two. After 60 iterations of attack, the accuracy obtained by our approach is 50% while the accuracy obtained by the other two remains only 20%. Therefore, we demonstrate the robustness and effectiveness of training data preprocessing of our noise patch-randomization scheme.

• Limitations

In our analysis, the attacker has all the information to our model and data so that they can correctly compute the gradient (White box attack). Even though the classifier generated by our method are robust to the attack compared to the classifier generated from the original model, the attacker can still successfully modify the class of a particular data when enough iterations of attack is performed. For black box attack, when the attacker does not have access to the model parameters, the attacker can still approximate the gradient information by averaging over the data using the its distribution. However, since the patches for every image is highly randomized, we expect that the average approximation of gradient may not be a good estimate such that it takes many iterations for the attacker to finish the update. However, under that condition, the perturbation of the image may be quite obvious. Therefore, this idea of designing a robust training data can be extended to nonlinear neural-network-based classifiers for that the gradient information becomes less useful compared to linear models.