



程序代写代做 CS编程辅导



Wee Adversarial Machine Learning – Vulnerabilities (Part I)

WeChat: cstutorcs

Assignment Project Exam Help

COMP90073
Email: tutorcs@163.com
Security Analytics

QQ: 749389476
Yi Han, CIS

<https://tutorcs.com>
Semester 2, 2021

程序代写代做 CS 编程辅导

- Week 9: Adversarial Machine Learning – Vulnerabilities
 - Definition + example
 - Classification
 - Evasion attacks
 - Gradient-descent based approaches
 - Automatic differentiation
 - Real-world example
 - Poisoning attacks
 - Transferability



WeChat: estutors

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

程序代写代做 CS 编程辅导

- Week 9: Adversarial Machine Learning – Vulnerabilities
 - Definition + example
 - Classification
 - Evasion attacks
 - Gradient-descent based approaches
 - Automatic differentiation
 - Real-world example
 - Poisoning attacks
 - Transferability



WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

程序代写代做 CS编程辅导

What is Adversarial Machine Learning (AML)?



“Adversarial machine learning is a technique employed in the field of machine learning which attempts to fool models through **malicious input**.” – Wikipedia

WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

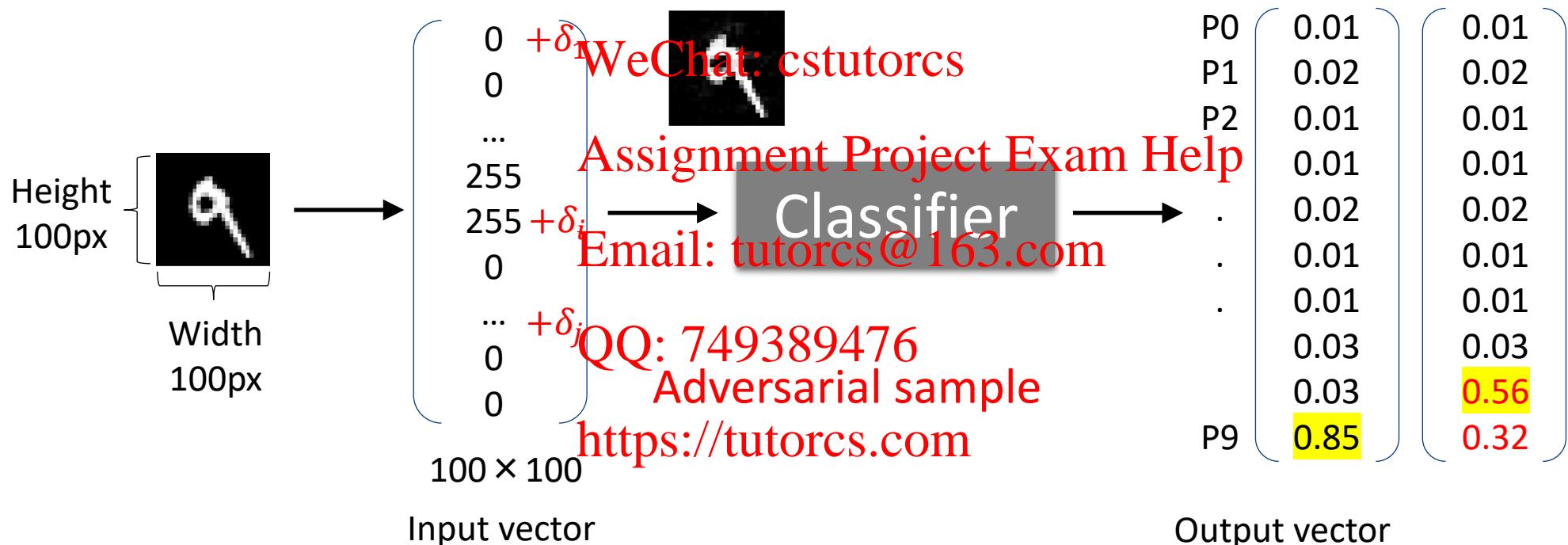
<https://tutorcs.com>

Examples

程序代写代做 CS编程辅导

- Test-time attack

- Image classifier C: inputs $X \rightarrow \{0, 1, 2, \dots, 9\}$,



Examples

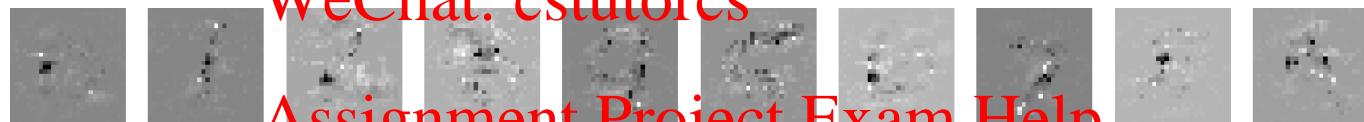
- Test-time attack

程序代写代做 CS编程辅导

Original images



Perturbation



WeChat: cstutorcs
Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

Adversarial samples



Classifier output

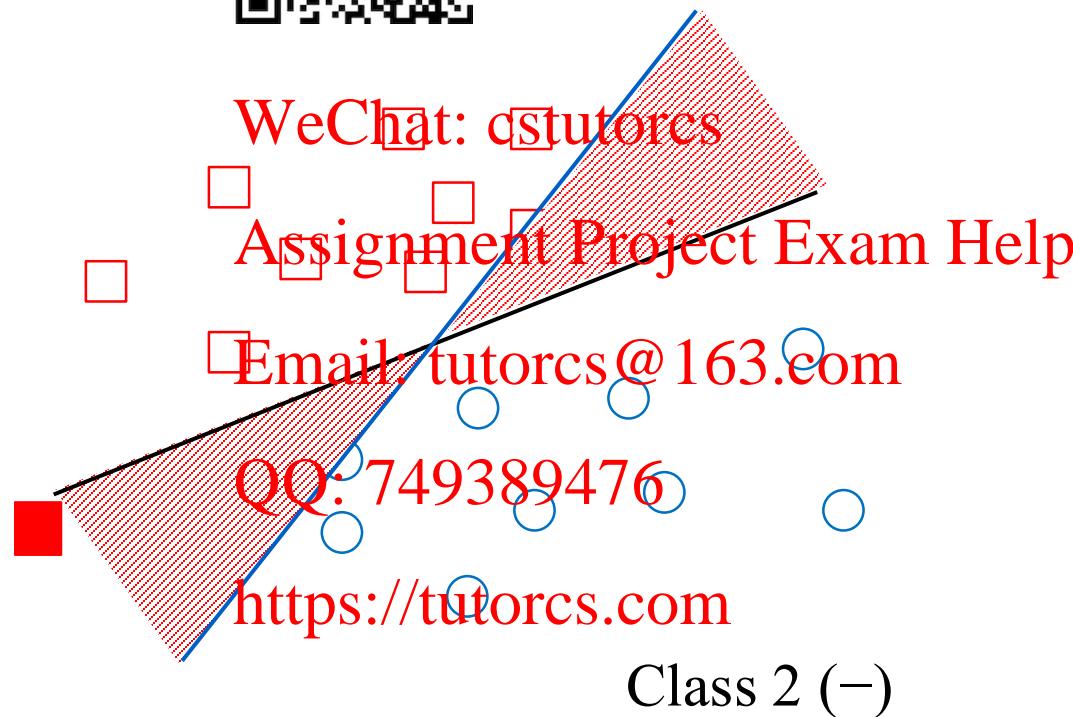
6 2 4 8 9 6 0 2 5 4

Examples

程序代写代做 CS编程辅导

- Training-time attack
 - Insert extra training samples to maximise the loss

Class 1



Examples

- Huge amount of attention
 - Mission-critical tasks



程序代写代做 CS编程辅导
WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>



程序代写代做 CS 编程辅导

- Week 9: Adversarial Machine Learning – Vulnerabilities

- Definition + example

- Classification

- Evasion attacks

- Gradient-descent based approaches

- Automatic differentiation

- Real-world example

- Poisoning attacks

- Transferability



WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

程序代写代做 CS编程辅导

- Classification [1]
 - Exploratory vs. Causative
 - Exploratory/evasive – time
 - Causative/poisonous – long time
- Integrity vs. Availability – security violation
 - Integrity: harmful instances to pass filters
 - Availability: denial of service, benign instances to be filtered
- Targeted vs. Indiscriminate/Untargeted – specificity
 - Targeted: misclassified as a specific class
 - Indiscriminate/untargeted: misclassified as any other class
- White-box vs. Black-box – attacker information
 - White-box: full knowledge of the victim model
 - Black-box: no/minimum knowledge of the model

WeChat: cstutors

Assignment Project Exam Help

Email: tutors@163.com

QQ: 749389476

<https://tutors.com>

程序代写代做 CS 编程辅导

- Week 9: Adversarial Machine Learning – Vulnerabilities

- Definition + examples
- Classification
- Evasion attacks
 - Gradient-descent based approaches
 - Automatic differentiation
 - Real-world example
- Poisoning attacks
- Transferability



WeChat: estutors

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

Evasion attacks (definition)

程序代写代做 CS编程辅导

- Evasion attack

- Aim: minimum perturbation δ to the input x , in order to cause model C to misclassify



$$x \rightarrow x + \delta \quad (x, \delta \in [0,1]^d)$$

WeChat: cstutorcs
such that (s.t.)

Assignment Project Exam Help
 $C(x + \delta) \neq C(x)$ Indiscriminate

OR
Email: tutorcs@163.com
 $C(x + \delta) = l_{target}$ Targeted
QQ: 749389476

<https://tutorcs.com>

程序代写代做 CS编程辅导

- Evasion attack

- Formulated as an optimization problem

$$\arg \min_{\delta} \|\delta\| \quad (1)$$

s. t. $C(x + \delta) \neq C(x)$
WeChat: cstutorcs] Highly non-linear
 OR $C(x + \delta) = l_{target}$
Assignment Project Exam Help

- p -norm: $\|\delta\|_p = \left(\sum_{t=1}^d |\delta_t|^p \right)^{1/p}$
 - $\|\delta\|_1 = \sum_{i=1}^d |\delta_i|$, $\|\delta\|_2 = \sqrt{\sum_{i=1}^d |\delta_i|^2}$, $\|\delta\|_\infty = \max_i |\delta_i|$
 - E.g., $\delta = \langle 1, 2, 3, -4 \rangle$, $\|\delta\|_1 = 10$, $\|\delta\|_2 = \sqrt{30}$, $\|\delta\|_\infty = 4$

Evasion attacks (definition)

程序代写代做 CS编程辅导

Transform (1) to the following problem [2]:

$$\arg \min_{\delta \in [0,1]^d} \|\delta\|$$



$$f(x + \delta)$$

Indiscriminate

$$\arg \min_{\delta \in [0,1]^d} \|\delta\|$$



$$f_{target}(x + \delta)$$

Targeted

WeChat: cstutorcs

Assignment Project Exam Help

Objective function f : how close the

Email: tutorcs@163.com

the prediction and the target are, e.g.,
the cross entropy loss function

QQ: 749389476

<https://tutorcs.com>

Evasion attacks (definition)

- Indiscriminate attack: $\arg \min_{\delta \in \mathcal{B}(x, d)} \|\delta\| - c \cdot f_{true}(x + \delta)$

四



WeChat: cstutorcs⁰

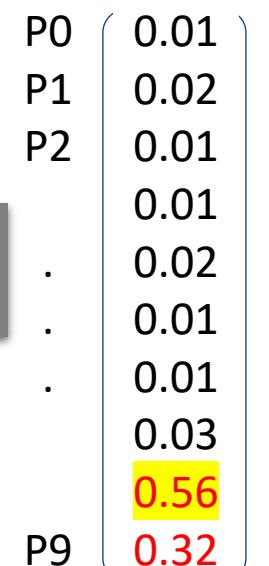
Assignment Project

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

Ground truth
(One-hot vector)



Evasion attacks (definition)

- Targeted attack: $\arg \min_{\delta \in \mathcal{S}^{target}} \|\delta\| + c \cdot f_{target}(x + \delta)$

A QR code located in the top right corner of the page, which links to additional resources for the Tutor CS program.

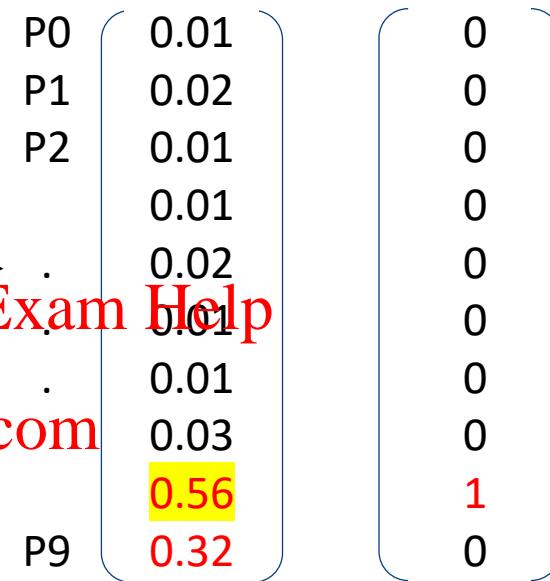
WeChat: cstutorcs

Classifier → Assignment Project Exam Help

Misclassified as “8” (target class)

QQ: 749389476

<https://tutorcs.com>



Prediction

Target

(One-hot vector)

Evasion attacks (definition)

程序代写代做 CS编程辅导

Transform (1) to the following problem [2]:

$$\arg \min_{\delta \in [0,1]^d} \|\delta\| \text{ s.t. } f(x + \delta)$$



Indiscriminate

$$\arg \min_{\delta \in [0,1]^d} \|\delta\| \text{ s.t. } f_{target}(x + \delta)$$

Targeted

WeChat: cstutorcs

How to find the minimum perturbation δ ? Assignment Project Exam Help

Email: tutorcs@163.com

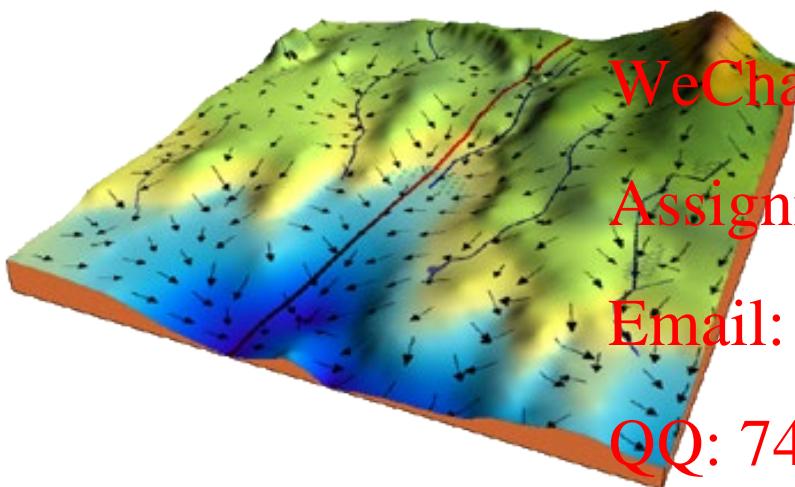
QQ: 749389476

<https://tutorcs.com>

程序代写代做 CS编程辅导

Gradient descent

- Gradient: a vector that points in the direction of greatest increase of a function

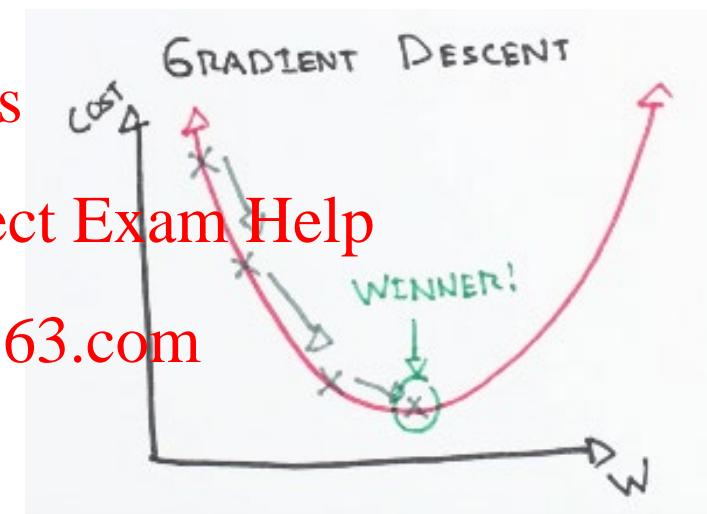


WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476



<https://tutorcs.com>

https://ml-cheatsheet.readthedocs.io/en/latest/gradient_descent.html

Evasion attacks (gradient descent)

程序代写代做 CS 编程辅导

$$\arg \min_{\delta \in [-\epsilon, \epsilon]^d} \|\delta\| - c \cdot f_{true}(x + \delta)$$

$$\arg \max_{\delta \in [-\epsilon, \epsilon]^d} + c \cdot f_{target}(x + \delta)$$

Indiscriminate

Targeted

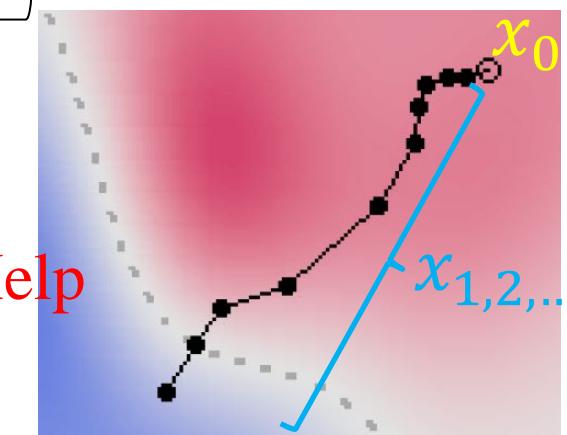
- Start with the initial input x_0
- Repeat $x_i \leftarrow x_{i-1} - \alpha \frac{\partial J}{\partial x_{i-1}}, i > 0$
- Until (1) $C(x_i) \neq C(x_0)$ (or $C(x_i) = l_{target}$), or → success
- (2) $\|\delta\| = \|x_i - x_0\| > \epsilon$, or
- (3) $i \geq i_{max}$, or <https://tutorcs.com>
- (4) $|J(x_i) - J(x_{i-1})| \leq \Delta$

WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476



→ failure

Evasion attacks (gradient descent)

程序代写代做 CS 编程辅导

$$\arg \min_{\delta \in [-\epsilon, \epsilon]^d} \|\delta\| - c \cdot f_{true}(x + \delta)$$



$$\arg \min_{\delta \in [-\epsilon, \epsilon]^d} + c \cdot f_{target}(x + \delta)$$

Indiscriminate

Targeted

J

WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

How to design the objective function?

<https://tutorcs.com>

程序代写代做 CS编程辅导

- Fast gradient sign method (FGSM) [3]:

$$\arg \min_{\delta \in [0,1]^d} \text{[X]} - c \cdot f_{true}(x + \delta)$$



$$\arg \min_{\delta \in [0,1]^d} \text{[X]} + c \cdot f_{target}(x + \delta)$$

WeChat: cstutorcs

$$\arg \min_{\delta \in [0,1]^d} -loss_{true}(x + \delta)$$

$f =$ cross
entropy loss

$$\arg \min_{\delta \in [0,1]^d} loss_{target}(x + \delta)$$

- Single step ϵ : fast rather than optimal

Assignment Project Exam Help

$$x' \leftarrow x + \epsilon \cdot sign\left(\frac{\partial loss_{true}}{\partial x}\right)$$

Email: tutorcs@163.com

QQ: 749389476

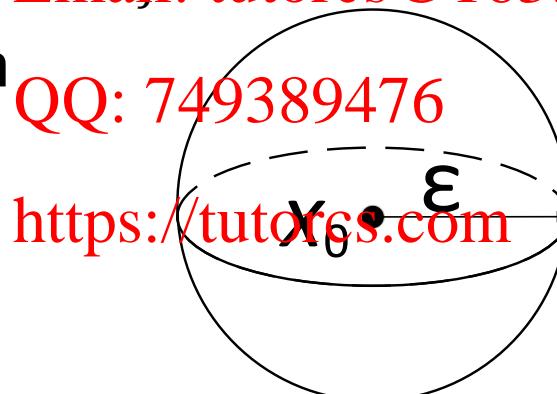
OR $x' \leftarrow x - \epsilon \cdot sign\left(\frac{\partial loss_{target}}{\partial x}\right)$

<https://tutorcs.com>

- Not meant to produce the minimal adversarial perturbations

程序代写代做 CS编程辅导

- Iterative gradient sign [4]
 - Single step $\epsilon \rightarrow$ many smaller steps α
 - $x_i \leftarrow \text{clip}_\epsilon \left(x_{i-1} + \alpha \cdot \frac{\partial f_{true}}{\partial x_{i-1}} \right)$ OR
 - $x_i \leftarrow \text{clip}_\epsilon \left(x_{i-1} - \alpha \cdot \text{sign} \left(\frac{\partial f_{target}}{\partial x_{i-1}} \right) \right)$
 - clip_ϵ : make sure that $x_{i,j}$ is within the range of $[x_{0j} - \epsilon, x_{0j} + \epsilon]$



QQ: 749389476

<https://tutorcs.com>

WeChat: cstutorcs
Assignment Project Exam Help

Email: tutorcs@163.com

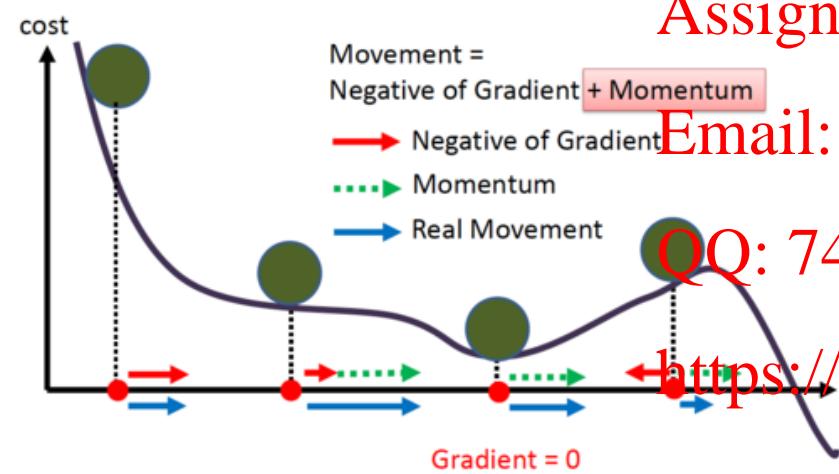
程序代写代做 CS编程辅导

- Momentum iterative fast gradient sign method

$$g_i = \mu \cdot g_{i-1} + \frac{\nabla_x}{\|\nabla_x\|}$$

- Momentum overcomes problems of vanilla gradient descent

- Get stuck in local minima
- Oscillation

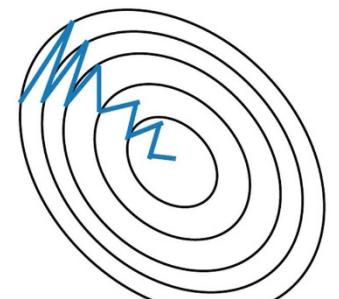


Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>



<https://medium.com/analytics-vidhya/momentum-rmsprop-and-adam-optimizer-5769721b4b19>

<https://eloquentarduino.github.io/2020/04/stochastic-gradient-descent-on-your-microcontroller/>

程序代写代做 CS编程辅导

C & W attack [2]



$$\arg \min_{\delta \in [0,1]^d} f(x + \delta)$$

$C(x + \delta) = l_{target}$ if and only if $f(x + \delta) \leq 0$

WeChat: cstutorcs

$C(x + \delta) \neq l_{target} \Leftrightarrow f(x + \delta) > 0$

Assignment Project Exam Help
Consistent with the definition of function f : how close
the prediction and the target are
Email: tutorcs@163.com
QQ: 749389476

<https://tutorcs.com>

程序代写代做 CS编程辅导
 $C(x + \delta) = l_{target}$ if and only if $f(x + \delta) = f(x') \leq 0$

- Option 1: $f(x') = \max\left(\max_{i \neq t} (F(x')_i - F(x')_t), 0\right)$
- Option 2: $f(x') = \log(1 + \exp(-\max_i (F(x')_i - F(x')_t))) - \log(2)$
- Option 3: $f(x') = \max(0.5 - F(x')_t, 0)$

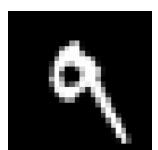


$F(x)$: output vector for x , i.e., probabilities of the input x belonging to each class. For example:

WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

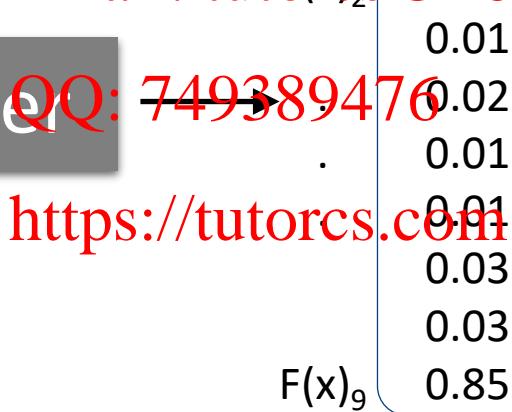


Classifier

QQ: 749389476

<https://tutorcs.com>

$F(x)_9 = 0.85$



程序代写代做 CS编程辅导

- CleverHans

- Do not use the latest version
- Download from: <https://github.com/tensorflow/cleverhans/releases/tag/v3.0.1>
- Prerequisite:

- Python3 (<https://www.python.org/downloads/>)
- Tensorflow (<https://www.tensorflow.org/install/>)
- Python 3.5/3.6/3.7 and TensorFlow {1.8, 1.10, 1.14}



- Installation:

- cd cleverhans
- pip install -e .

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

WeChat: cstutorcs

Assignment Project Exam Help

程序代写代做 CS 编程辅导

- Week 9: Adversarial Machine Learning – Vulnerabilities

- Definition + examples
- Classification
- Evasion attacks
 - Gradient-descent based approaches
 - Automatic differentiation
 - Real-world example
- Poisoning attacks
- Transferability



WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

程序代写代做 CS 编程辅导

$$\arg \min_{\delta \in [-\epsilon, \epsilon]^d} \|\delta\| - c \cdot f_{true}(x + \delta)$$

Indiscriminate

$$\arg \max_{\delta \in [-\epsilon, \epsilon]^d} + c \cdot f_{target}(x + \delta)$$

Targeted

WeChat: cstutorcs

- Start with the initial input x_0
- Repeat $x_i \leftarrow x_{i-1} - \alpha \frac{\partial J}{\partial x_{i-1}}, i > 0$
- Until $C(x_i) \neq C(x_0)$ (or $C(x_i) = l_{target}$)

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

How to calculate the partial derivatives?
<https://tutorcs.com>

程序代写代做 CS编程辅导

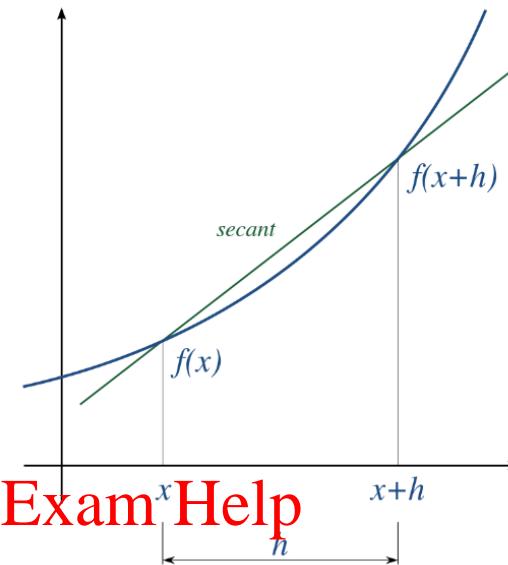
Derivative

- Definition: $f'(x) = \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h}$
- Numerical differentiation
 - $\frac{f(x+h) - f(x)}{h}, \frac{f(x+h) - f(x-h)}{2h}$
 - Significant round-off errors



WeChat: cstutorcs

Assignment Project Exam Help



Email: tutorcs@163.com

- Symbolic differentiation: apply chain rules to symbolic expressions
 - Exponentially-long results

<https://tutorcs.com>

程序代写代做 CS编程辅导

Automatic differentiation

- A set of techniques to automatically evaluate the derivative of a function specified by a computer program – Wikipedia
- Any complicated function can be rewritten as the composition of a sequence of primitive functions:

$$f = f_0 \circ f_1 \circ f_2 \circ \dots \circ f_n$$

- Apply the chain rule [Assignment](#) [Project](#) [Exam](#) [Help](#)

- Forward mode: $\frac{\partial f}{\partial x} = \frac{\partial f_0}{\partial x} \left(\frac{\partial f_1}{\partial f_0} \left(\dots \left(\frac{\partial f_n}{\partial f_{n-1}} \frac{\partial f_n}{\partial x} \right) \right) \right)$
Email: tutorcs@163.com
QQ: 749389476

- Reverse mode: $\frac{\partial f}{\partial x} = \left(\left(\left(\frac{\partial f_0}{\partial f_1} \frac{\partial f_0}{\partial f_2} \right) \dots \right) \frac{\partial f_n}{\partial f_{n-1}} \right) \frac{\partial f_n}{\partial x}$



程序代写代做 CS编程辅导

- Given $y = f(x_1, x_2) = \ln(x_1) + x_1 x_2 - \sin(x_2)$, calculate $\frac{\partial y}{\partial x_1}$ at (2,5)
- Forward mode [5]



Output variables:

$$y_{m-i} = v_{l-i}, i = m - 1, \dots, 0$$

Assignment Project Exam Help

Working variables:

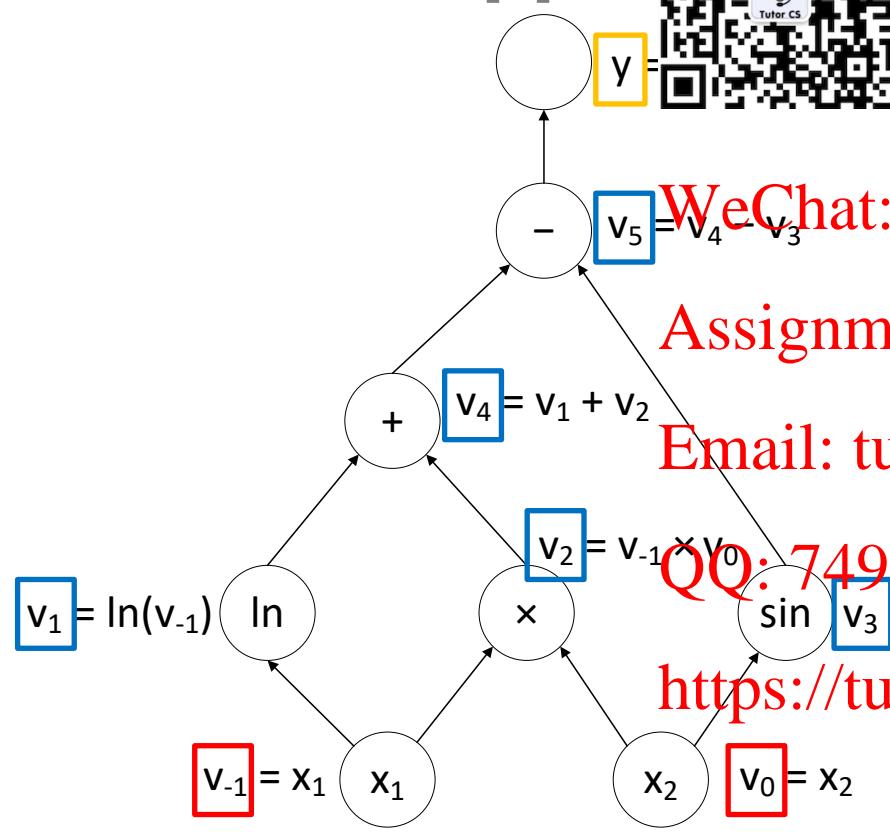
Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

$$v_0 = x_2$$

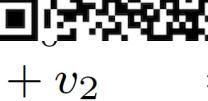
Computational graph

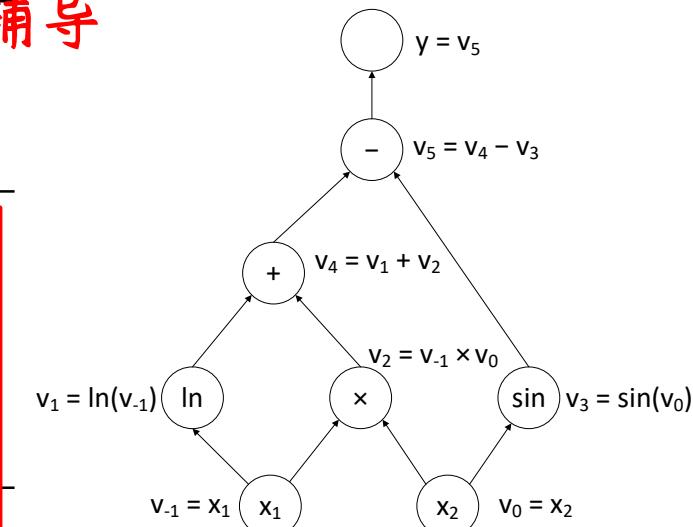


$$v_{i-n} = x_i, i = 1, \dots, n$$

Evasion attacks (automatic differentiation)

程序代写代做 CS编程辅导

Forward evaluation trace		
$v_{-1} = x_1$	= 2	
$v_0 = x_2$	= 5	
$v_1 = \ln v_{-1}$	= 	= ln 2
$v_2 = v_{-1} \times v_0$	= 	= 2 × 5
$v_3 = \sin v_0$	= 	= sin 5
$v_4 = v_1 + v_2$	= 0.6931 + 10	
$v_5 = v_4 + v_3$	= 10.6931 + 0.9589	
$y = v_5$	= 11.6521	



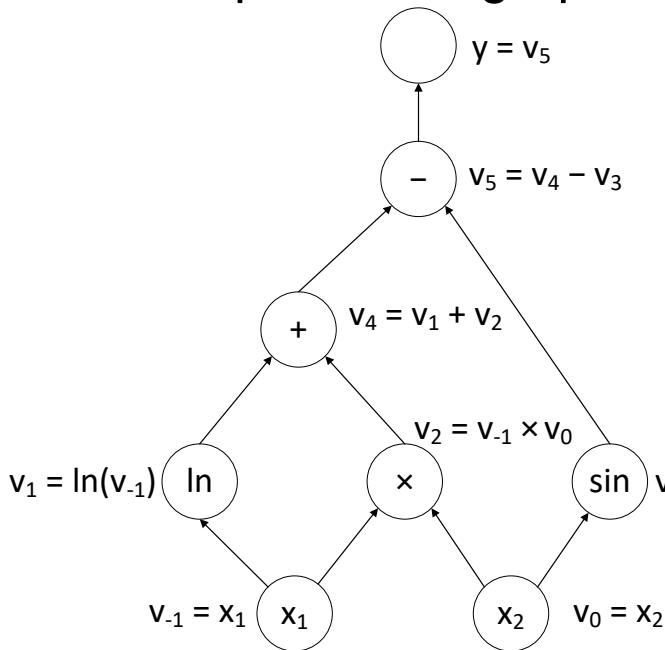
Assignment Project Exam Help

Forward derivative trace

$\dot{v}_i = \frac{\partial v_i}{\partial x_1}$	$\dot{v}_{-1} = \dot{x}_1$	= 1	Email: tutorcs@163.com
	$\dot{v}_0 = \dot{x}_2$	= 0	
	$\dot{v}_1 = \dot{v}_{-1}/v_{-1}$	= 1/2	
	$\dot{v}_2 = \dot{v}_{-1} \times v_0 + v_{-1} \times \dot{v}_0$	= 1 × 5 + 2 × 0	
	$\dot{v}_3 = \cos v_0 \times \dot{v}_0$	= cos 5 × 0	
	$\dot{v}_4 = \dot{v}_1 + \dot{v}_2$	= 0.5 + 5	
	$\dot{v}_5 = \dot{v}_4 - \dot{v}_3$	= 5.5 - 0	
	$\dot{y} = \dot{v}_5$	= 5.5	2-cos5

- Reverse mode [5]

Computational graph



程序代写代做CS编程辅导

Forward evaluation trace



WeChat: cstutorcs

Assignment Project Exam Help

Reverse adjoint trace

Email: tutorcs@163.com

Adjoint

$\bar{v}_i = \frac{\partial y}{\partial v_i}$

QQ: 749389476

<https://tutorcs.com>

$$\begin{array}{lll} v_{-1} & = x_1 & = 2 \\ v_0 & = x_2 & = 5 \end{array}$$

$$\begin{array}{lll} v_1 & = \ln v_{-1} & = \ln 2 \\ v_2 & = v_{-1} \times v_0 & = 2 \times 5 \\ v_3 & = \sin v_0 & = \sin 5 \\ v_4 & = v_1 + v_2 & = 0.6931 + 10 \\ v_5 & = v_4 - v_3 & = 10.6931 + 0.9589 \\ y & = v_5 & = 11.6521 \end{array}$$

Reverse adjoint trace

$$\begin{array}{lll} \bar{x}_1 & = \bar{v}_{-1} & = 5.5 \\ \bar{x}_2 & = \bar{v}_0 & = 1.7163 \\ \bar{v}_{-1} & = \bar{v}_{-1} + \bar{v}_1 (\partial v_1 / \partial v_{-1}) & = \bar{v}_{-1} + \bar{v}_1 / v_{-1} \\ \bar{v}_0 & = \bar{v}_0 + \bar{v}_2 (\partial v_2 / \partial v_0) & = \bar{v}_0 + \bar{v}_2 \times v_{-1} \\ \bar{v}_1 & = \bar{v}_2 (\partial v_2 / \partial v_{-1}) & = \bar{v}_2 \times v_0 \\ \bar{v}_2 & = \bar{v}_4 (\partial v_4 / \partial v_2) & = \bar{v}_3 \times \cos v_0 \\ \bar{v}_3 & = \bar{v}_5 (\partial v_5 / \partial v_3) & = \bar{v}_4 \times 1 \\ \bar{v}_4 & = \bar{v}_5 (\partial v_5 / \partial v_4) & = \bar{v}_4 \times 1 \\ \bar{v}_5 & = \bar{y} & = \bar{v}_5 \times (-1) \\ & & = \bar{v}_5 \times 1 \end{array}$$

程序代写代做 CS 编程辅导

- Example 1: $y = \ln(x_1) + x_1 x_2 - \sin(x_2)$

- Calculate $\left(\frac{\partial y}{\partial x_1}, \frac{\partial y}{\partial x_2} \right)$

- Forward mode:



- Reverse mode: time(s)

- Example 2: $y_1 = \ln(x) + x, y_2 = x - \sin(x)$

- Calculate $\left(\frac{\partial y_1}{\partial x}, \frac{\partial y_2}{\partial x} \right)$

Email: tutorcs@163.com

- Forward mode: time(s)
QQ: 749389476

- Reverse mode: time(s)
<https://tutorcs.com>

程序代写代做 CS编程辅导

Function $f: R^n \rightarrow R^m$

- n independent x_i as inputs, m dependent y_j as outputs
- Forward mode: $m \gg n$, forward run can calculate $\frac{\partial y}{\partial x_i}$
- Reverse mode: $n \gg m$, reverse run can calculate $\frac{\partial y_j}{\partial x}$



#http://laid.delanover.com/gradients-in-tensorflow/
Assignment Project Exam Help
Import tensorflow as tf

Tensorflow example

Email: tutorcs@163.com

~~y = tf.Variable(2.)~~

~~z = tf.subtract(2*x, y)~~

~~grad = tf.gradients(z, [x, y])~~

~~https://tutorcs.com~~

```
sess = tf.Session()
sess.run(tf.global_variables_initializer())
print(sess.run(grad)) # [2.0, -1.0]
```

程序代写代做 CS 编程辅导

- Week 9: Adversarial Machine Learning – Vulnerabilities

- Definition + examples

- Classification

- Evasion attacks

- Gradient-descent based approaches

- Automatic differentiation

- Real-world example

- Poisoning attacks

- Transferability



WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

Evasion attacks (real-world example)

- Robust Physical-World Attacks on Deep Learning Visual Classification [6]
 - Stop sign, Right Turn speed Limit 45
 - Drive-By (Field) Test
 - Start from 250 ft away
 - Classify every 10th frame

程序代写代做 CS编程辅导



Table 1: Sample of physical adversarial examples against LISA-CNN and GTSRB-CNN.

WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

40' 0°



Targeted-Attack Success

100%

73.33%

66.67%

100%

80%

程序代写代做 CS 编程辅导

- Week 9: Adversarial Machine Learning – Vulnerabilities

- Definition + examples
- Classification
- Evasion attacks
 - Gradient-descent based approaches
 - Automatic differentiation
 - Real-world example
- Poisoning attacks
- Transferability



WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

Poisoning attacks

- 程序代写代做 CS编程辅导
- Insert extra points to maximally decrease the accuracy [8]



Poisoning attacks

程序代写代做 CS编程辅导

- Attacker's aim: maximise the hinge loss over the validation data
 $D_{val} = \{x_i, y_i\}_{i=1}^m$
- Optimisation problem

$$\arg \max_{x_c} \sum_{i=1}^m \left(1 - y_i f_{x_c}(x_i) \right)$$

WeChat: cstutorcs

To find the optimal poisoning data x_c :

Assignment Project Exam Help

- Random initial attack point x_c
- Update: re-compute the SVM;

Email: tutorcs@163.com

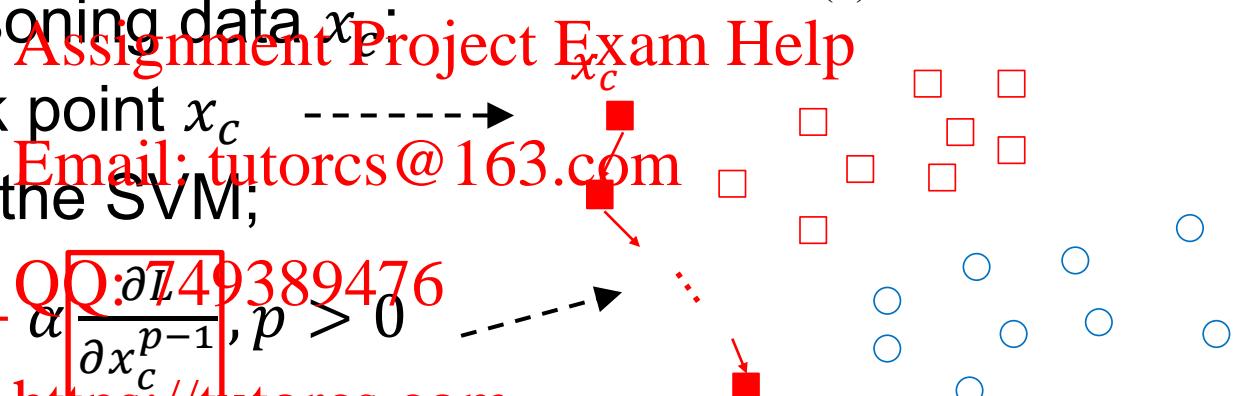
$$x_c^p \leftarrow x_c^{p-1} + \alpha \frac{\partial L}{\partial x_c^{p-1}}, p > 0$$

- Until $L(x_c^p) - L(x_c^{p-1}) < \varepsilon$

<https://tutorcs.com>

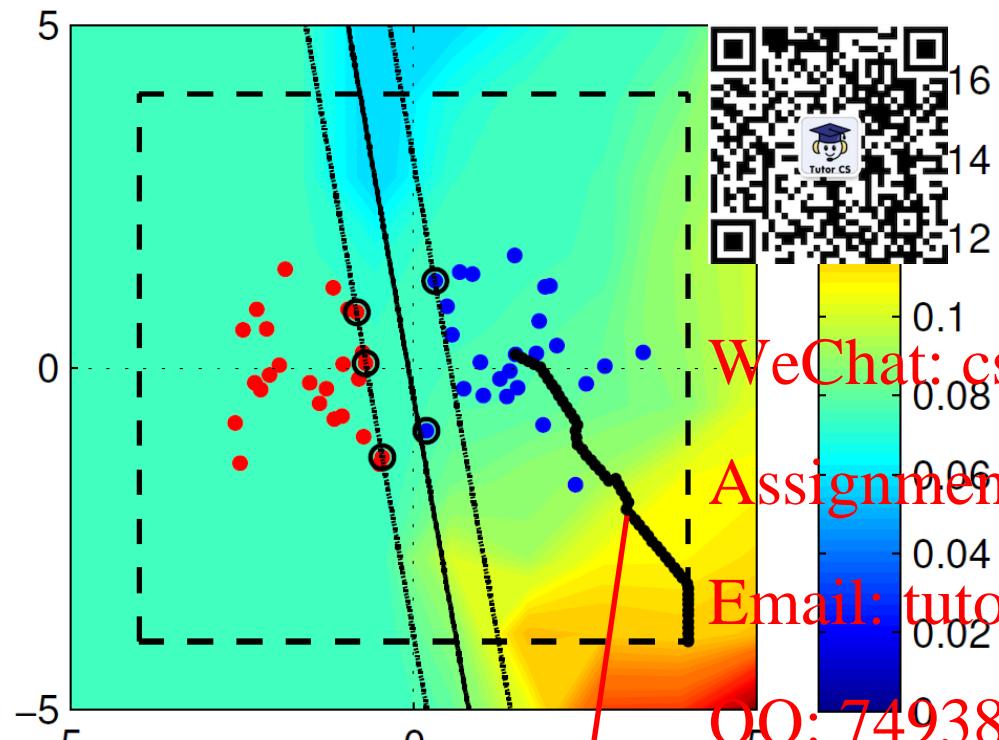
Class 1 (+)

Class 2 (-)



Poisoning attacks

mean $\sum_i \xi_i$ (hinge loss)



程序代写代做 CS编程辅导

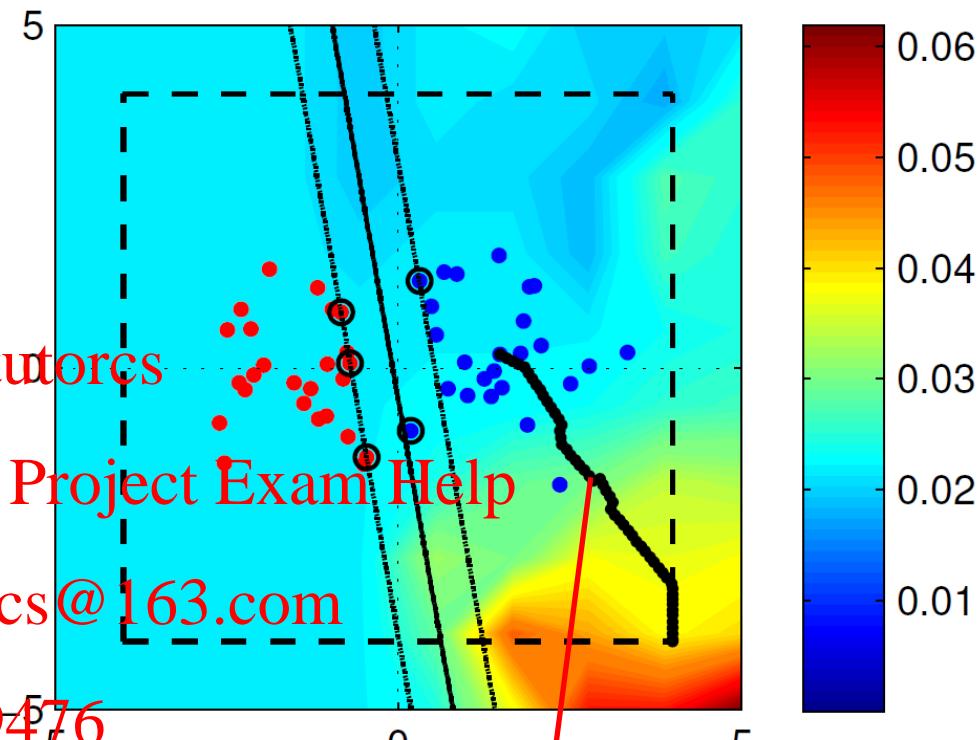
classification error

WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476



As the attack point x_c moves towards a local maximum, both the hinge loss and the classification error increase.

Poisoning attacks

程序代写代做 CS编程辅导

- Poison frog attacks [10]

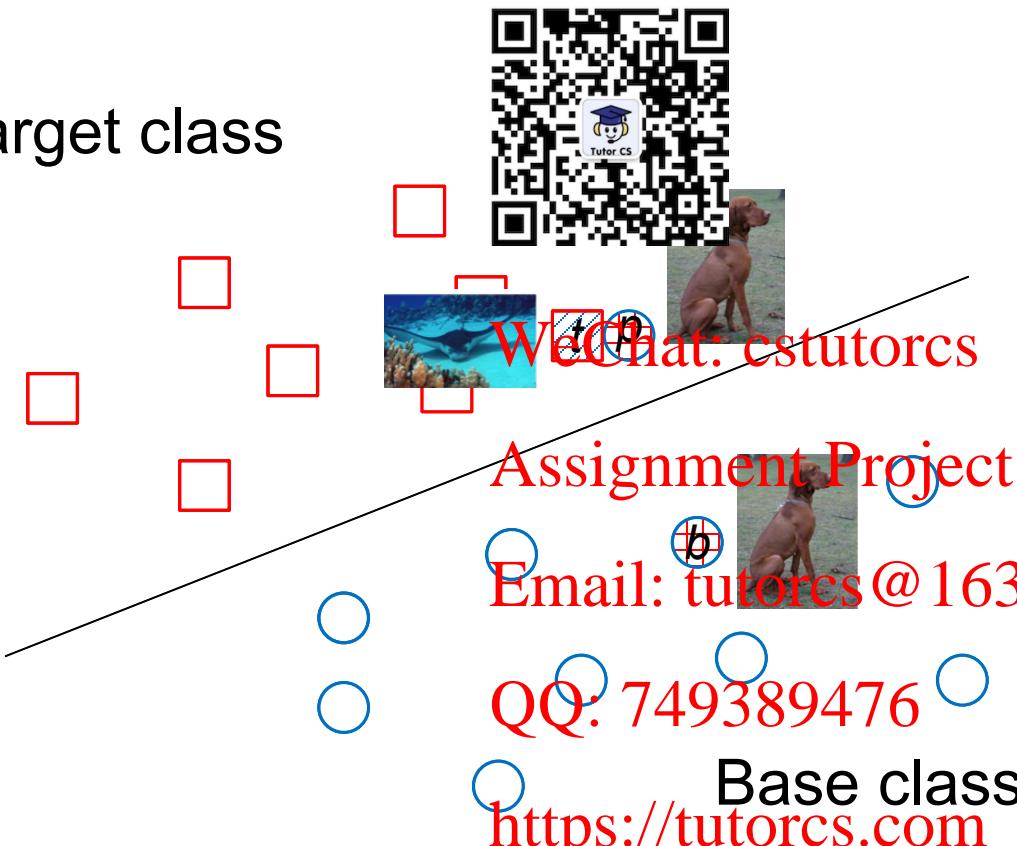
- E.g., add a seemingly normal image (that is properly labeled) to a training set, and control the quantity of a chosen image at test time



Poisoning attacks

程序代写代做 CS编程辅导

Target class



Step 1: choose an instance from the target class – t (target instance)

Step 2: sample an instance from the base class – b (base instance)

Step 3: perturb b to create a poison instance – p

Step 4: inject p into the training dataset

The model is then re-trained.
The attack succeeds if the poisoned model labels t as the base class

程序代写代做 CS编程辅导

- Generate poison data p
 - Optimisation problem: $\underset{x}{\operatorname{argmin}} \|f(x) - f(t)\|_2^2 + \beta \|x - b\|_2^2$
 - $f(x)$: output of the second last layer of the neural network
 - $\|f(x) - f(t)\|_2^2$: makes p move toward the target instance in **feature space** and get embedded in the target class distribution
 - $\beta \|x - b\|_2^2$: makes p appear like a base class instance to a human labeller

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

Tutor CS

WeChat: cstutorcs

程序代写代做 CS编程辅导

- Forward-backward-splitting iterative procedure [11]
 - Forward step: gradient descent update to minimise the L2 distance to the target instance in feature space
 - Backward step: projection update that minimises the Euclidean distance from the base instance in input space

WeChat: cstutorcs

Algorithm 1 Poisoning Example Generation Assignment Project Exam Help

Input: target instance t , base instance b , learning rate λ

Initialize x : $x_0 \leftarrow b$

Email: tutorcs@163.com

Define: $L_p(x) = \|f(\mathbf{x}) - f(\mathbf{t})\|^2$

for $i = 1$ to $maxIters$ do

QQ: 749389476

Forward step: $\hat{x}_i = x_{i-1} - \lambda \nabla_x L_p(x_{i-1})$

Backward step: $x_i = (\hat{x}_i + \lambda \beta b) / (1 + \beta \lambda)$

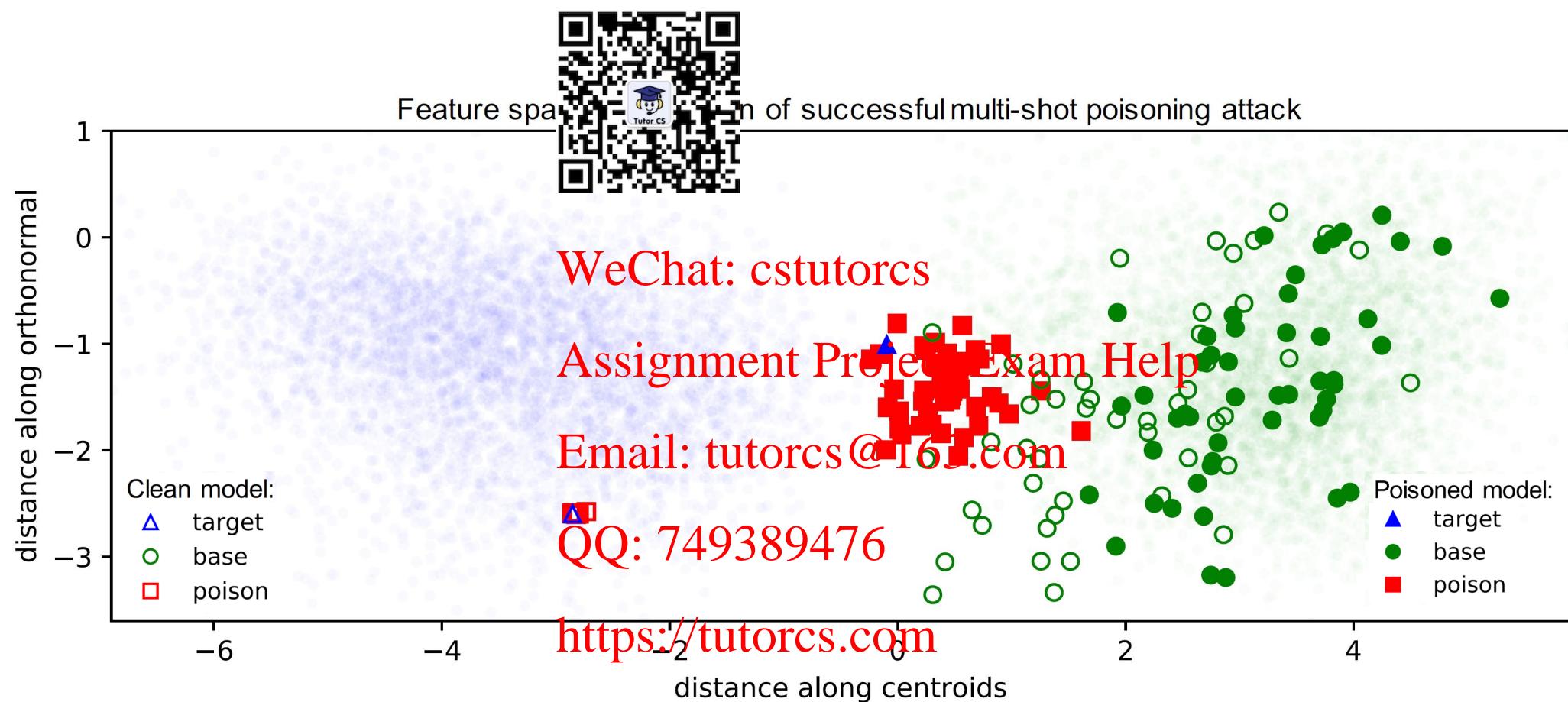
end for

<https://tutorcs.com>

Poisoning attacks

程序代写代做 CS编程辅导

- Results



程序代写代做 CS编程辅导

Using Machine Teaching to Identify Optimal Training-Set Attacks on
Machine Learners [9]

- Attacker's objective : $\hat{\theta}_D = \|\hat{\theta}_D - \theta^*\| + \|D - D_0\|_2$
 - $\hat{\theta}_D$: parameters of the poisoned model after the attack
 - θ^* : parameters of the attacker's target model, i.e., model that the attacker aims to obtain
 - D : poisoned training data
 - D_0 : original training data

<https://tutorcs.com>

程序代写代做 CS编程辅导

- Week 9: Adversarial Machine Learning – Vulnerabilities
 - Definition + examples
 - Classification
 - Evasion attacks
 - Gradient-descent based approaches
 - Automatic differentiation
 - Real-world example
 - Poisoning attacks
 - Transferability



WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

程序代写代做 CS编程辅导

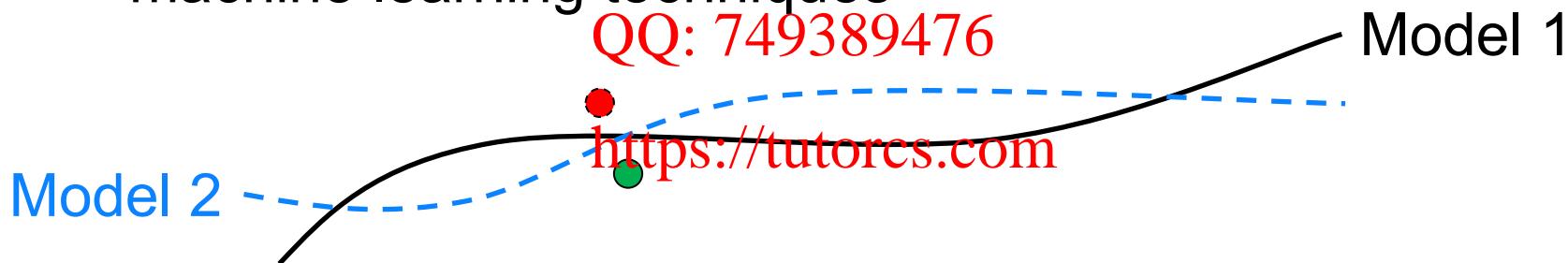
- Implicit assumption: full knowledge of the target model
- What if the target model is unknown to the attacker?
- Transferability: for two models that perform the same task, trained on different datasets, adversarial samples generated against one model can often fool the other model as well [12][13]
 - Intra-technique: both the target and surrogate model use the same machine learning technique
 - Inter-technique: the target and surrogate model use different machine learning techniques

WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476



程序代写代做 CS编程辅导

- Verification on the MNIST dataset of handwritten digits

- Grey-scale, 0-255
- Size: 28px * 28px



0
1
2
3
4
5
6
7
8
9 9

WeChat: cstutorcs

<https://upload.wikimedia.org/wikipedia/commons/0/27/MnistExamples.png>

- DNN, SVM, LR, DT, kNN
- Black-box attack

- Step 1: adversary trains their own model – surrogate/source
- Step 2: generate adversarial samples against the surrogate
- Step 3: apply the adversarial samples against the target model

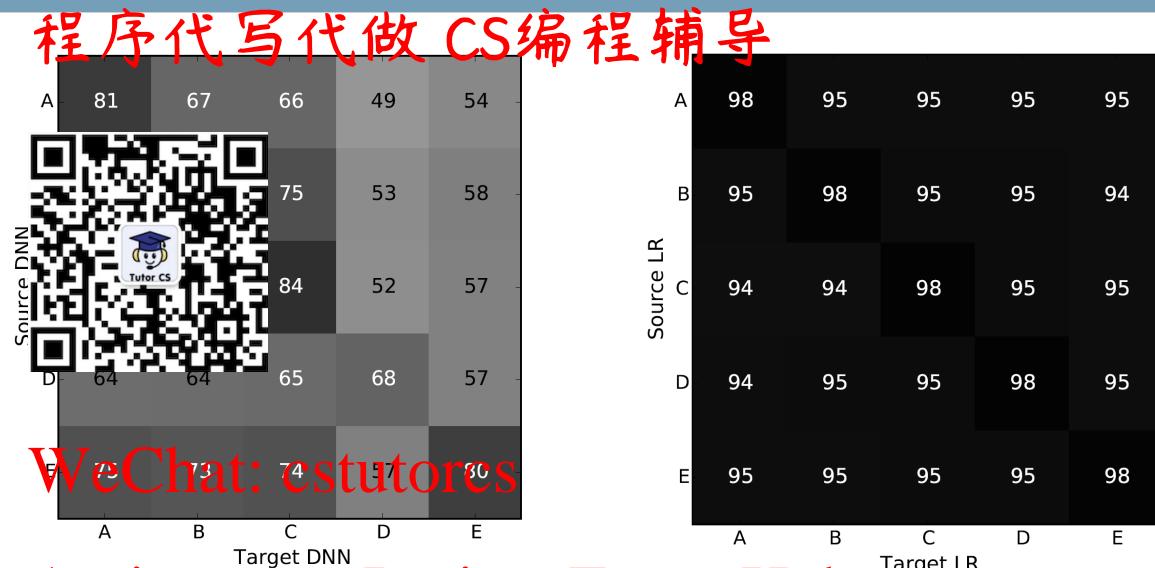
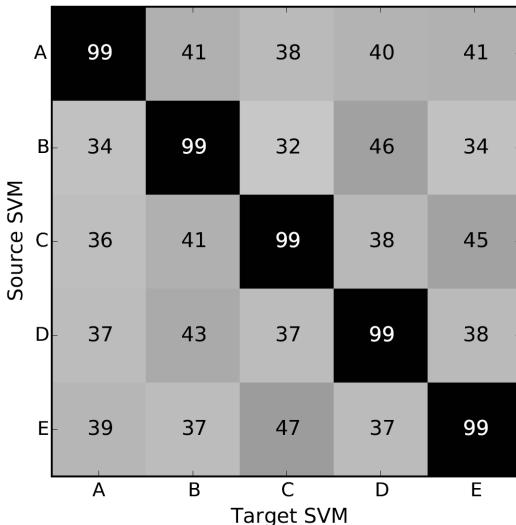
QQ: 749389476

<https://tutorcs.com>

Transferability & Black-box attacks

- Intra-technique

71% adv. samples against the source are also effective against the target



Assignment Project Exam Help

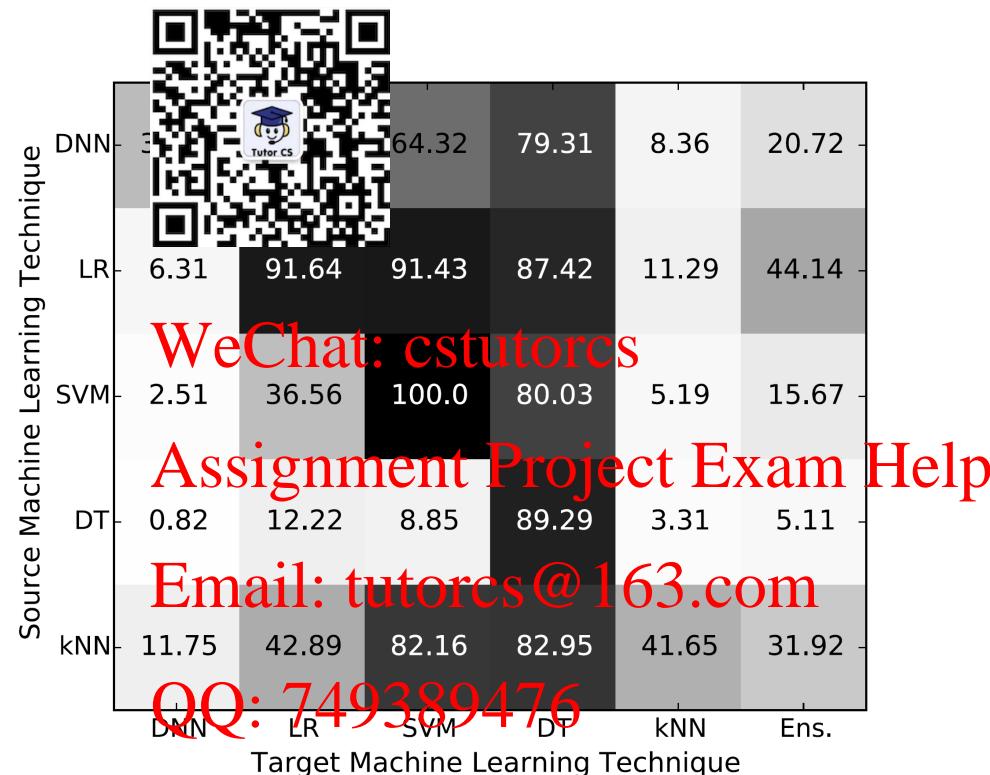
Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

- Inter-technique

程序代写代做 CS编程辅导



WeChat: cstutorcs
Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

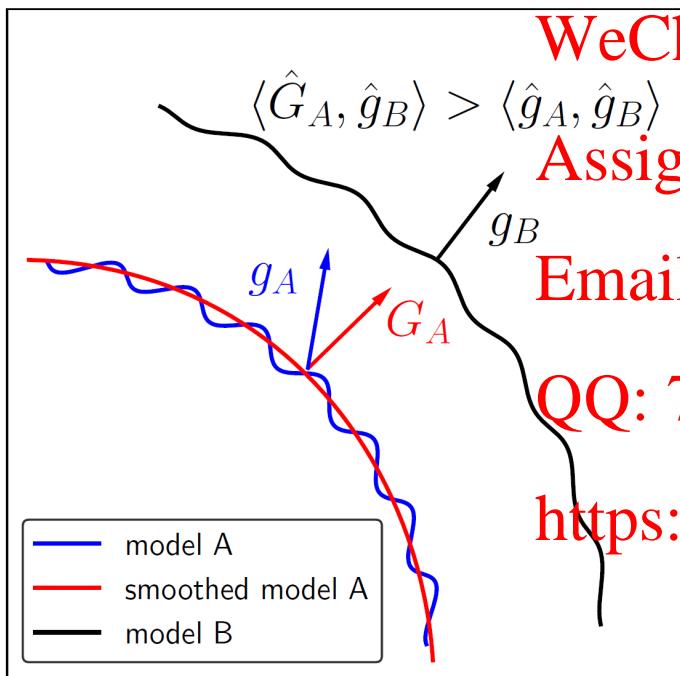
<https://tutorcs.com>

程序代写代做 CS编程辅导

- Non-smoothness can hurt the transferability [7]

- A is the surrogate model for the target model
- Smoothed loss surface contributes to transferability

$$x_i \leftarrow x_{i-1} - \alpha \frac{\partial J(x_{i-1})}{\partial x} \quad \text{vs} \quad x_i \leftarrow x_{i-1} - \alpha \frac{1}{m} \sum_{j=1}^m \frac{\partial J(x_{i-1} + \xi_j)}{\partial x}, \xi_j \sim \mathcal{N}(0, \sigma^2)$$



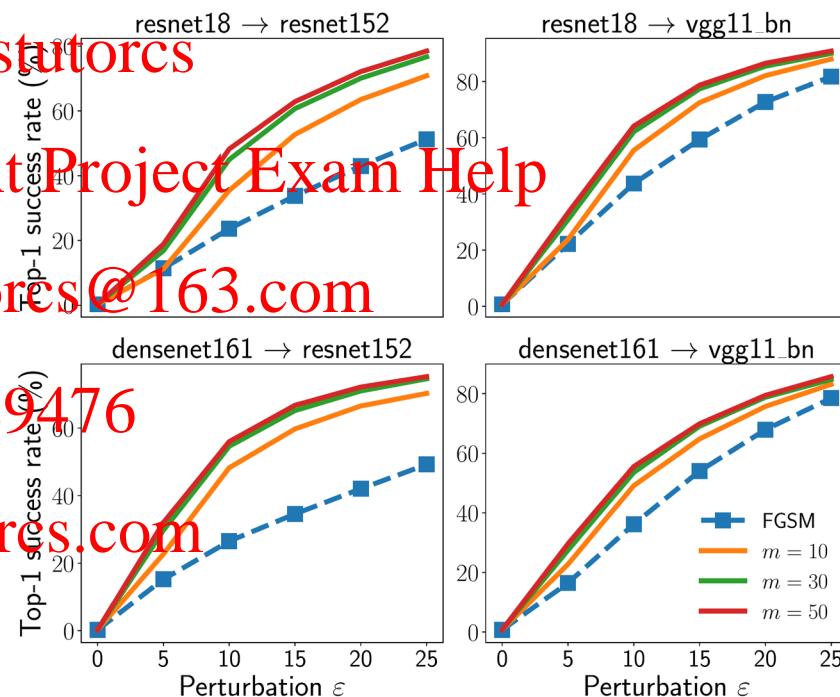
WeChat: cstutorcs

Assignment Project Exam Help

Email: tutores@163.com

QQ: 749389476

<https://tutores.com>



程序代写代做 CS编程辅导

- Input diversity improves transferability [15]
 - Adversarial samples  fit to the surrogate model
 - Data augmentation
 - Random resizing: stretch input image to a random size
 - Random padding: pad zeros around an image in a random manner
 - Diverse Inputs Iterative Fast Gradient Sign Method (DI²-FGSM)

$$x_i \leftarrow x_{i-1} - \alpha \cdot \text{sign} \left(\frac{\partial J(x_{i-1})}{\partial x} \right)$$

WeChat: **tutorcs**
Assignment Project Exam Help

$$x_i \leftarrow x_{i-1} - \alpha \cdot \text{sign} \left(\frac{\partial J(T(x_{i-1}; p))}{\partial x} \right) \quad T(x_{i-1}; p) = \begin{cases} T(x_{i-1}) & \text{with prob. } p \\ x_{i-1} & \text{with prob. } 1 - p \end{cases}$$

- Momentum Diverse Inputs Iterative Fast Gradient Sign Method (M-DI²-FGSM)

$$g_i = \mu \cdot g_{i-1} + \frac{\nabla_x J(x_{i-1})}{\|\nabla_x J(x_{i-1})\|_1} \quad g_i \leftarrow \mu \cdot g_{i-1} + \frac{\nabla_x J(T(x_{i-1}; p))}{\|\nabla_x J(T(x_{i-1}; p))\|_1}$$

QQ: 749389476
<https://tutorcs.com>

- Backpropagation smoothness [16], backpropagation linearity [17]

- Non-linear activation e.g., ReLU, sigmoid
- Non-continuous derivative zero during backpropagation
- Continuous derivatives can improve transferability
- Keep the ReLU function in the forward pass, but during backpropagation approximate the ReLU derivative with a continuous derivative, e.g. using softplus function ($\log(1 + e^x)$)

WeChat: estutorcs

Assignment Project Exam Help



Figure 1. Activation functions (left) and their derivatives (right).

程序代写代做 CS编程辅导

- Evasion attacks

- Indiscriminate: $\arg \min_{\delta \in [0,1]^d} \cdot f_{true}(x + \delta)$
 - Targeted: $\arg \min_{\delta \in [0,1]^d} f_{target}(x + \delta)$

- Poisoning attacks

- Attacker's objective: $Q_4(D, \hat{\theta}_D) = \|\hat{\theta}_D - \theta^*\| + \|D - D_0\|_2$

- $\hat{\theta}_D$: poisoned model after the attack

- θ^* : attacker's target, i.e., model that the attacker aims to obtain
 - D : poisoned training data
 - D_0 : original training data

- Transferability

QQ: 749389476

- Intra, inter-technique
 - Black-box attacks



WeChat: cstutorcs

Assignment Project Exam Help

ning data

Email: tutorcs@163.com

References

- [1] M. Barreno, B. Nelson, A. D. Joseph, and J. D. Tygar, “The Security of Machine Learning,” *Machine Learning*, vol. 81, no. 2, pp. 121–148, Nov. 2010.
- [2] N. Carlini and D. Wagner, “Towards Evaluating the Robustness of Neural Networks,” *eprint arXiv:1608.04645*, 2016.
- [3] I. J. Goodfellow, J. Shlens, and C. Szegedy, “Explaining and Harnessing Adversarial Examples,” *eprint arXiv:1412.6572*, 2014.
- [4] A. Kurakin, I. Goodfellow, and S. Bengio, “Adversarial examples in the physical world,” *arXiv preprint arXiv:1607.02533*, 2016.
- [5] A. G. Baydin and B. A. Pearlmutter, “Automatic Differentiation of Algorithms for Machine Learning,” *arXiv:1404.7456 [cs, stat]*, Apr. 2014.
- [6] I. Evtimov *et al.*, “Robust Physical-World Attacks on Machine Learning Models,” *arXiv preprint arXiv:1707.08945*, 2017.
- [7] Wu, L. and Zhu, Z., “Towards Understanding and Improving the Transferability of Adversarial Examples in Deep Neural Networks.” Proceedings of The 12th Asian Conference on Machine Learning:837–850. Available from <https://proceedings.mlr.press/v129/wu20a.html>.



WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749989476

[Https://tutorcs.com](https://tutorcs.com)

References

- 程序代写代做 CS编程辅导
- Assignment Project Exam Help
- WeChat: cstutors
- Email: tutors@163.com
- QQ: 749389476
- [8] B. Biggio, B. Nelson, and P. Laskov, “Poisoning Attacks against Support Vector Machines,” in *Proceedings of the 29th International Conference on International Conference on Machine Learning*, Edinburgh, Scotland, 2012, pp. 1467–1474.
 - [9] S. Mei and X. Zhu, “Using Machine Teaching to Identify Optimal Training-set Attacks on Machine Learners,” in *Proceedings of the Twenty-Ninth AAAI Conference on Artificial Intelligence*, Austin, Texas, 2015, pp. 2871–2877.
 - **[10] A. Shafahi et al., “Poison Frogs! Targeted Clean-Label Poisoning Attacks on Neural Networks,” arXiv:1804.00792 [cs, stat], Apr. 2018.**
 - [11] T. Goldstein, C. Studer, and R. Baraniuk, “A field guide to forward-backward splitting with a fasta implementation,” arXiv preprint arXiv:1411.3406, 2014
 - **[12] N. Papernot, P. McDaniel, and I. Goodfellow, “Transferability in Machine Learning: from Phenomena to Black-Box Attacks using Adversarial Samples,” eprint arXiv:1605.07720, 2016.**
 - [13] N. Papernot, P. McDaniel, I. Goodfellow, S. Jha, Z. B. Celik, and A. Swami, “Practical Black-Box Attacks against Deep Learning Systems using Adversarial Examples,” eprint arXiv:1602.02697, 2016.
 - [14] Yinpeng Dong, Fangzhou Liao, Tianyu Pang, Hang Su, Jun Zhu, Xiaolin Hu, and Jianguo Li, “Boosting adversarial attacks with momentum,” in CVPR, 2018.

References

程序代写代做 CS编程辅导

- [15] Xie, Cihang and Zhang, Zhishuai and Zhou, Yuyin and Bai, Song and Wang, Jianyu and Ren, Zhou and Tang, Jianxin, “Improving Transferability of Adversarial Examples with Input Diversification,” in CVPR, 2019.
- [16] Chaoning Zhang*, Philip H. S. Torr*, Gyusang Cho*, Adil Karjauv, Soomin Ham, Chan-Hyun Youn, In So Kyung, “Backpropagating Smoothly Improves Transferability of Adversarial Examples,” in IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshop on Adversarial Machine Learning, 2021 (*: Equal Contribution)
- [17] Y. Guo, Q. Li, and H. Chen, “Backpropagating linearly improves transferability of adversarial examples,” arXiv preprint arXiv:2012.03528, 2020

WeChat: **tutorcs**

Assignment Project Exam Help

Email: **tutorcs@163.com**

QQ: **749389476**

<https://tutorcs.com>