



FIT2093 INTRODUCTION TO CYBERSECURITY

Assignment Project Exam Help

<https://tutorcs.com>

Week 11 Lecture

WeChat: cstutorcs

Machine Learning / AI for Cybersecurity

S1 2022



Outline

Machine Learning for Cybersecurity

- Machine learning: What is it?
 - types of **learning**
 - types of machine learning **problems**
 - generative vs discriminative <https://cstutorcs.com>
 - e.g. *pattern recognition, feature extraction*
- Assignment Project Exam Help
- <https://cstutorcs.com>
- WeChat: cstutorcs



Outline

Machine Learning for Cybersecurity

- Machine learning for CyberSecurity problems

- **binary** classification
- **anomaly** detection
- **lie** detection
- **biometrics**
- **spam** classification
- **Tasks** in classification
- **Classification Metrics**
- **cryptography**

Assignment Project Exam Help

<https://tutorcs.com>

WeChat: cstutorcs



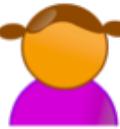
• Question

Machine Learning for Cybersecurity

- **Q:** Give an example scenario where you think Machine Learning/AI might be used to **defend** a system against attacks.

[Assignment Project Exam Help](https://tutorcs.com)
<https://tutorcs.com>

WeChat: cstutorcs



• Question

Machine Learning for Cybersecurity

- **Q:** Give an example scenario where you think Machine Learning/AI might be used to **attack a system**

[Assignment Project Exam Help](https://tutorcs.com)
<https://tutorcs.com>

WeChat: cstutorcs

Activity (5 mins)

- 1) Click the latest link in the Zoom chat
- 2) Add your question response to the Ed forum
- 3) Add your “hearts” to your favourite responses

Machine Learning: What is it?

Assignment Project Exam Help

<https://tutorcs.com>

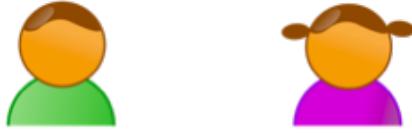
WeChat: cstutorcs



Machine Learning: What is it?

Machine Learning for Cybersecurity

- Machine learning
 - learn from experience, improve over time
 - **learn** to do what? <https://tutorcs.com>
 - how to **improve**?
 - with **time**, what **experience** is enhanced?
WeChat: cstutorcs
- Includes but is not just optimization
 - max or min some objective function s.t. some constraints

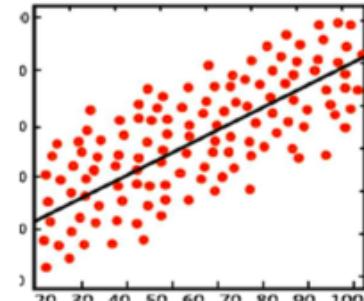
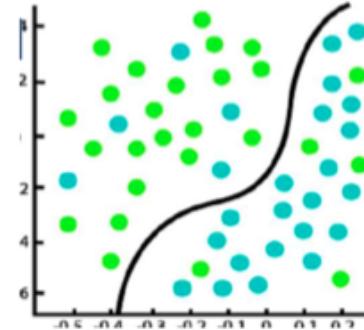


Machine Learning: What is it?

Machine Learning for Cybersecurity

- Machine *learning*: learn from *experience*, *improve* over *time*
 - learn to do
 - **classification**: choose 1 of N labels/categories/classes
<https://tutores.com>
 - **regression**: predict numerical values

WeChat: cstutorcs





Machine Learning: What is it?

Machine Learning for Cybersecurity

- Machine *learning*: learn from *experience*, *improve* over *time*

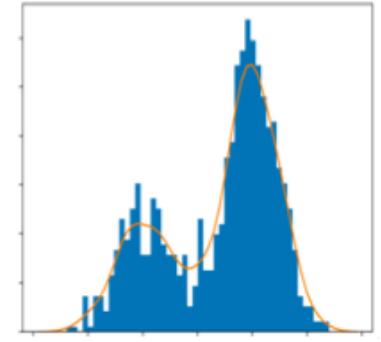
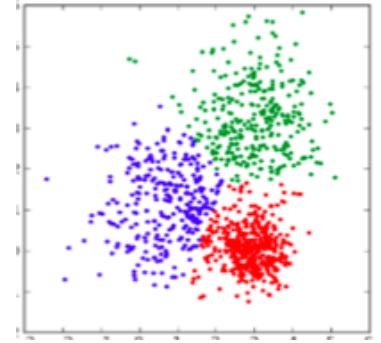
- learn to do

Assignment Project Exam Help

- **clustering**: group similar samples

WeChat: cstutorcs

- **density estimation**: construct the probability distribution of observed samples

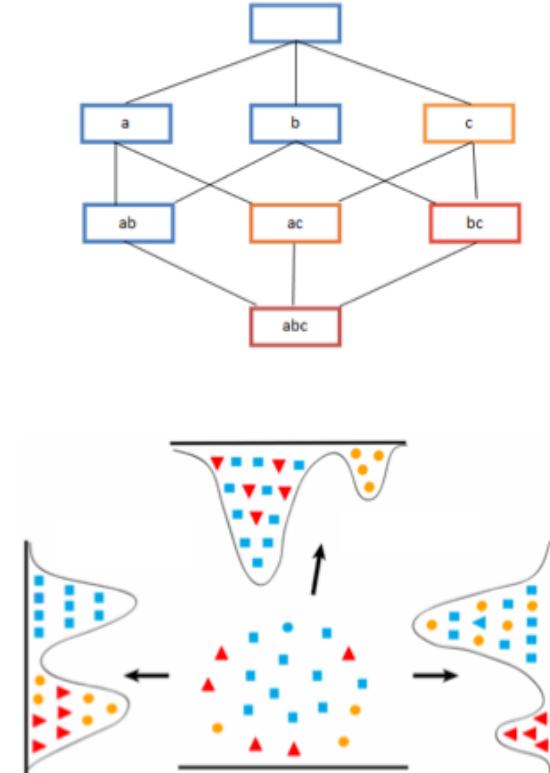




Machine Learning: What is it?

Machine Learning for Cybersecurity

- Machine *learning*: learn from *experience*, *improve over time*
 - learn to do
 - **association rule learning**: discover relations between samples
<https://tutorcs.com>
WeChat: cstutorcs
 - **dimensionality reduction / feature learning / representation learning**: reduce the dimensions of features in the dataset of samples

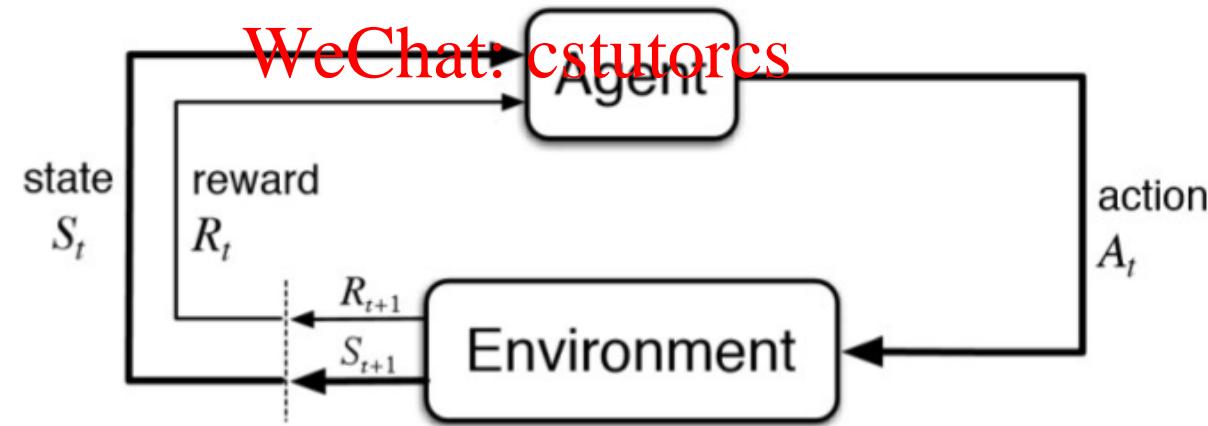




Machine Learning: What is it?

Machine Learning for Cybersecurity

- Machine learning: learn from *experience*, *improve over time*
 - learn to do
 - reinforcement learning: game-like training to max the reward vs penalty
<https://cstutorcs.com>





Machine Learning: What is it?

Machine Learning for Cybersecurity

- Machine *learning*: learn from *experience*, *improve* over *time*
 - how to **improve**? **Assignment Project Exam Help**
 - ↑ accuracy of prediction/estimation, <https://tutorcs.com>
 - ↓ difference from desired
WeChat: cstutorcs
 - with **time**, what **experience** is enhanced?
 - see more (labelled) samples e.g. (un)supervised learning
 - more interactions (reinforcement learning)



Machine Learning: What is it?

Machine Learning for Cybersecurity

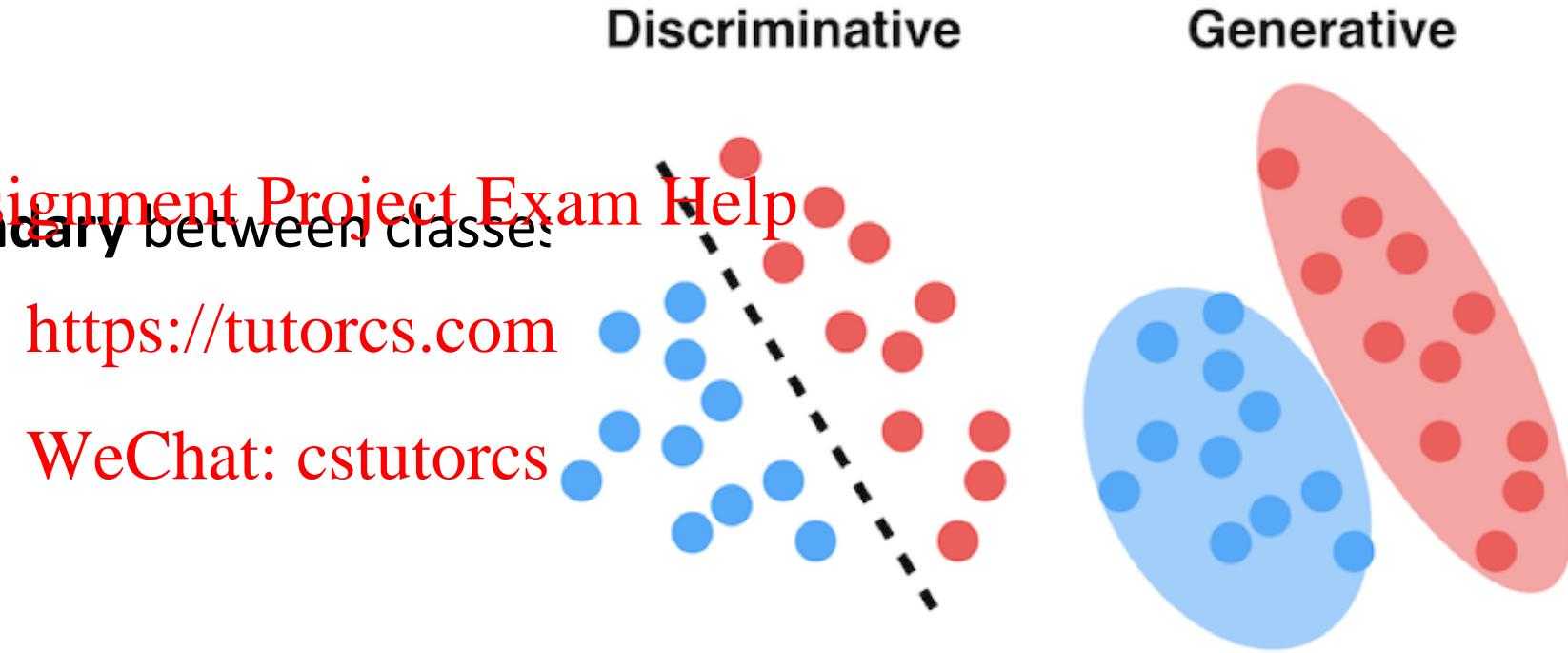
- **supervised** learning: *we'll focus on this, most common*
 - ground truth samples available with class labels
- **unsupervised** learning <https://tutorcs.com>
 - samples without labels, no ground truth
- **semi-supervised**: samples with & without labels
- **reinforcement** learning
 - no samples, learn from interactions & corresponding reward/payoff



Machine Learning: What is it?

Machine Learning for Cybersecurity

- Discriminative Model
 - learn the decision **boundary** between classes among samples
 - **predict the class label**
- Generative Model
 - model the underlying **distribution** of the samples
 - can **generate new samples** from same distribution





Pattern Recognition < Machine Learning

Machine Learning for Cybersecurity

- **pattern recognition = feature learning then classification**

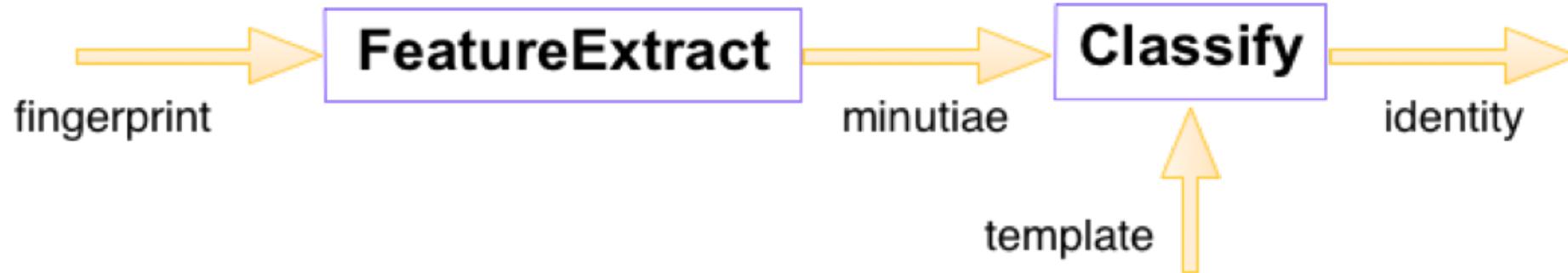
- can treat each stage as a black box, various options

<https://tutorcs.com>

- e.g. for *fingerprint biometrics*

- *feature = minutiae*

WeChat: cstutorcs





Fingerprint Minutiae

Machine Learning for Cybersecurity



Pattern Recognition

Machine Learning for Cybersecurity



- biometrics = pattern recognition applied to a security problem
Assignment Project Exam Help
- Q: Why need to do feature extraction? Why not directly classify?
<https://tutorcs.com>

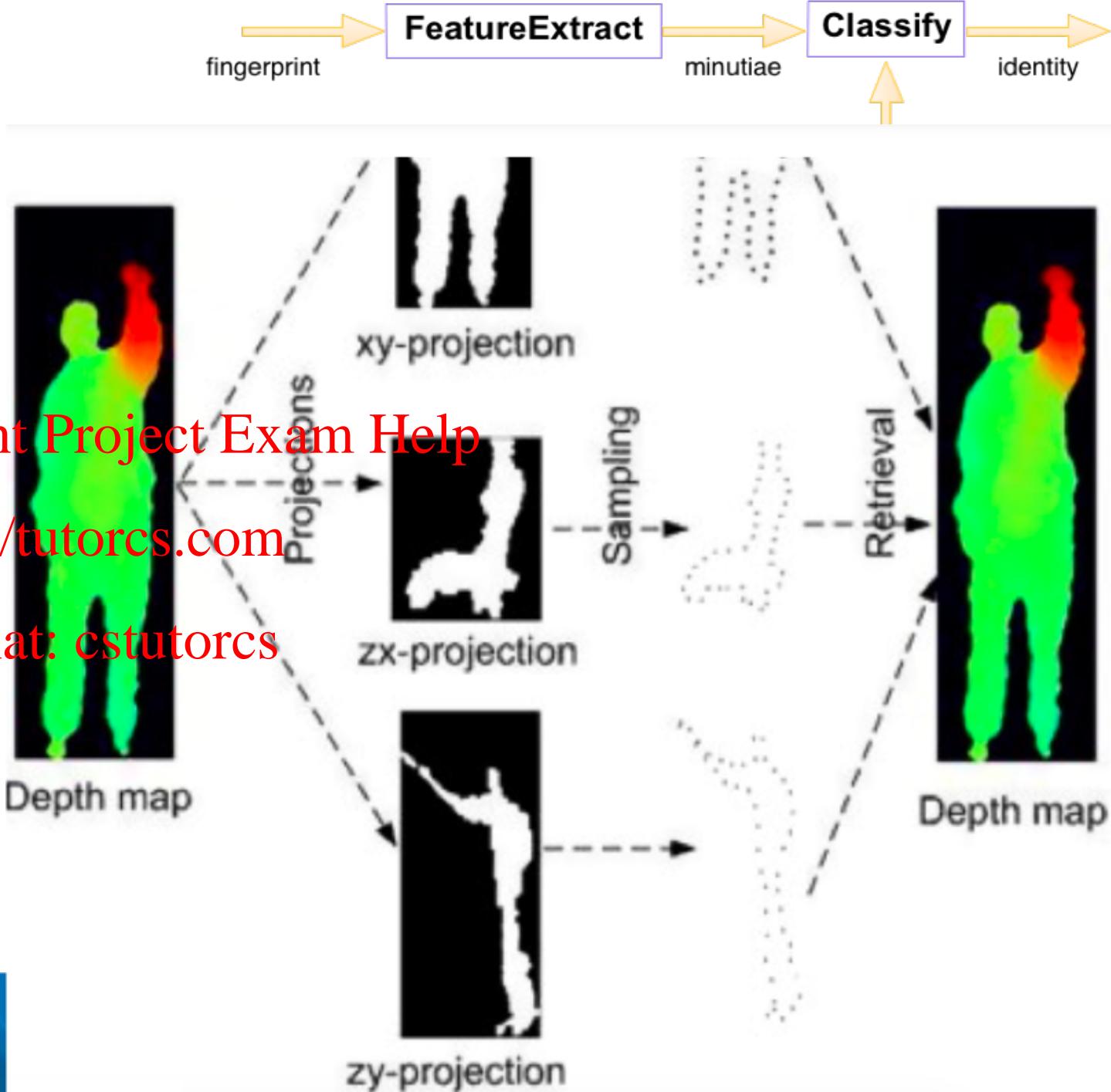
WeChat: cstutorcs

Feature Extraction

Machine Learning for Cybersecurity

- Q: Why need to do feature extraction? Why not directly classify?
- A: this allows to obtain various projection views, features from various viewpoints are suited for different classification problems e.g. top-down view of a human does not help to predict his/her height, side view is better

18



Machine Learning for Assignment Project Exam Help CyberSecurity Problems

<https://tutorcs.com>

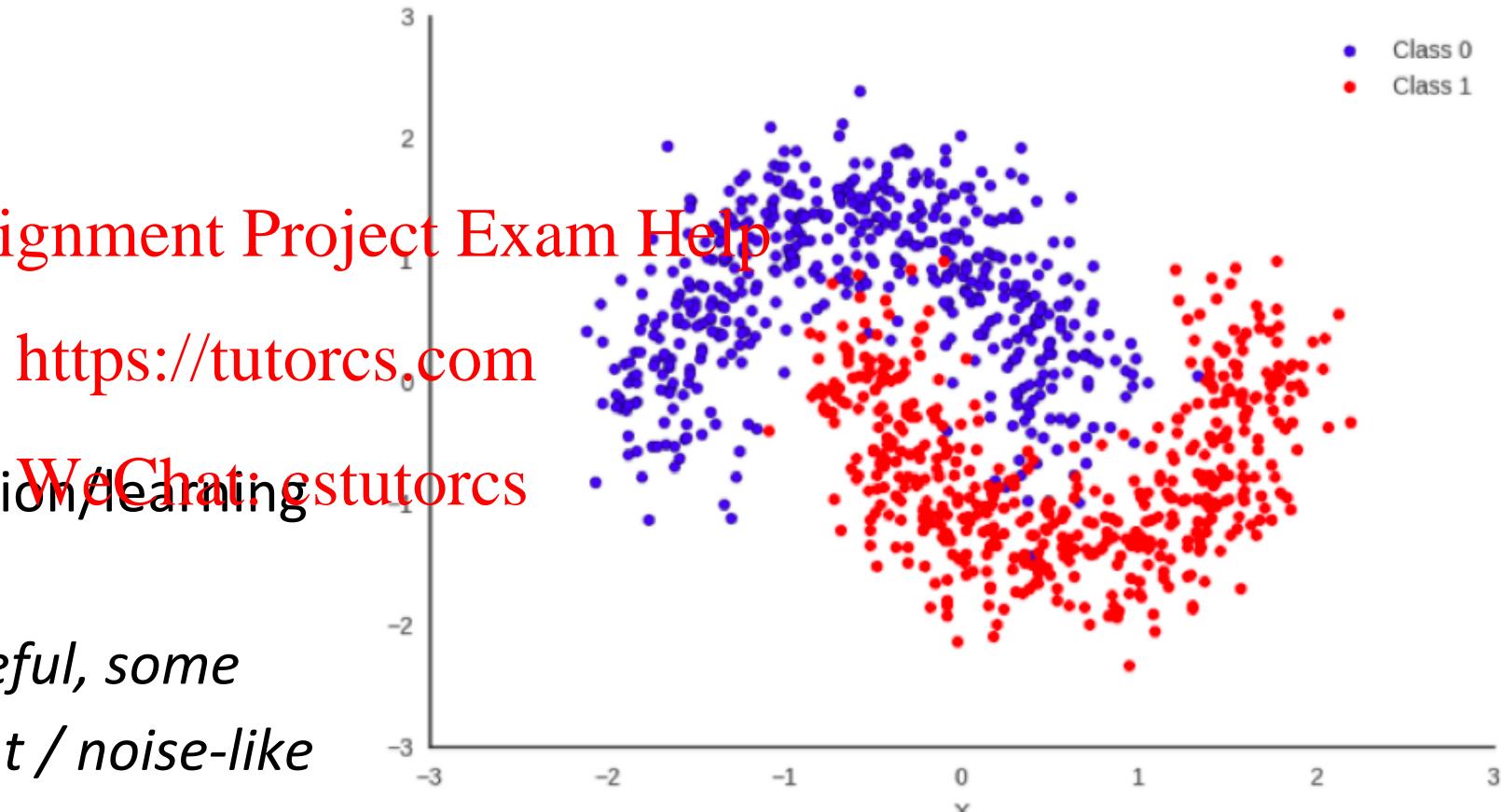
WeChat: cstutorcs



Security Problems as Classification Problems

Machine Learning for Cybersecurity

- binary classification
 - 2 classes/categories
 - yes / no
 - involves feature extraction/learning before classification
 - *not everything is useful, some redundant / noise-like*
 - *we only need to focus on some key aspects*





Anomaly Detection as Binary Classification

Machine Learning for Cybersecurity

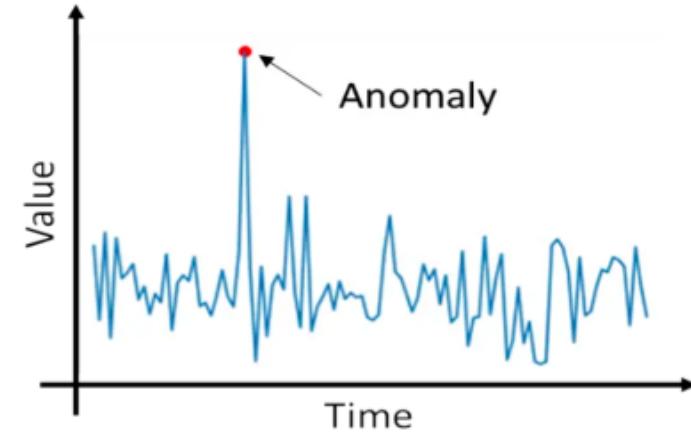
- binary classification: 2 classes
 - normal / abnormal **Assignment Project Exam Help**

○ benign / malicious traffic/behaviour e.g. *intrusion/malware detection*

- (I) detect **presence** of known **malicious** signatures
- (II) detect **absence** of known **normal** behaviour

WeChat: cstutorcs

- **Q:** when is approach (II) better than (I)?
- **Q:** when will both approaches not work?
- **Q:** how could an attacker cause both approaches to fail?

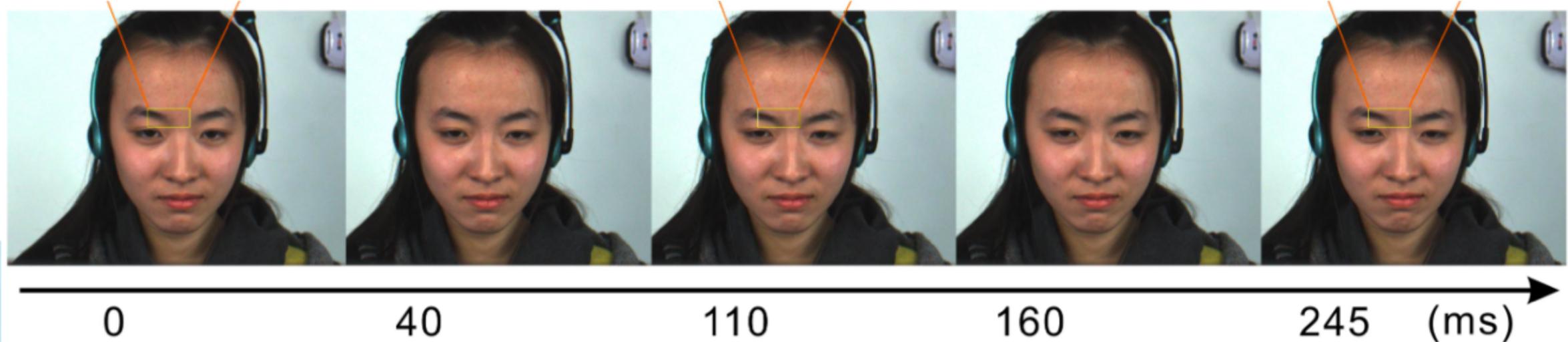




Lie Detection as Binary Classification

Machine Learning for Cybersecurity

- lie detection:
 - innocent vs suspicious
 - e.g. presence of micro-expressions <https://tutores.com>
 - subtle facial movements, fraction of a second
 - involuntary, sub-conscious reaction to emotion





Biometrics as Multiclass Classification

Machine Learning for Cybersecurity

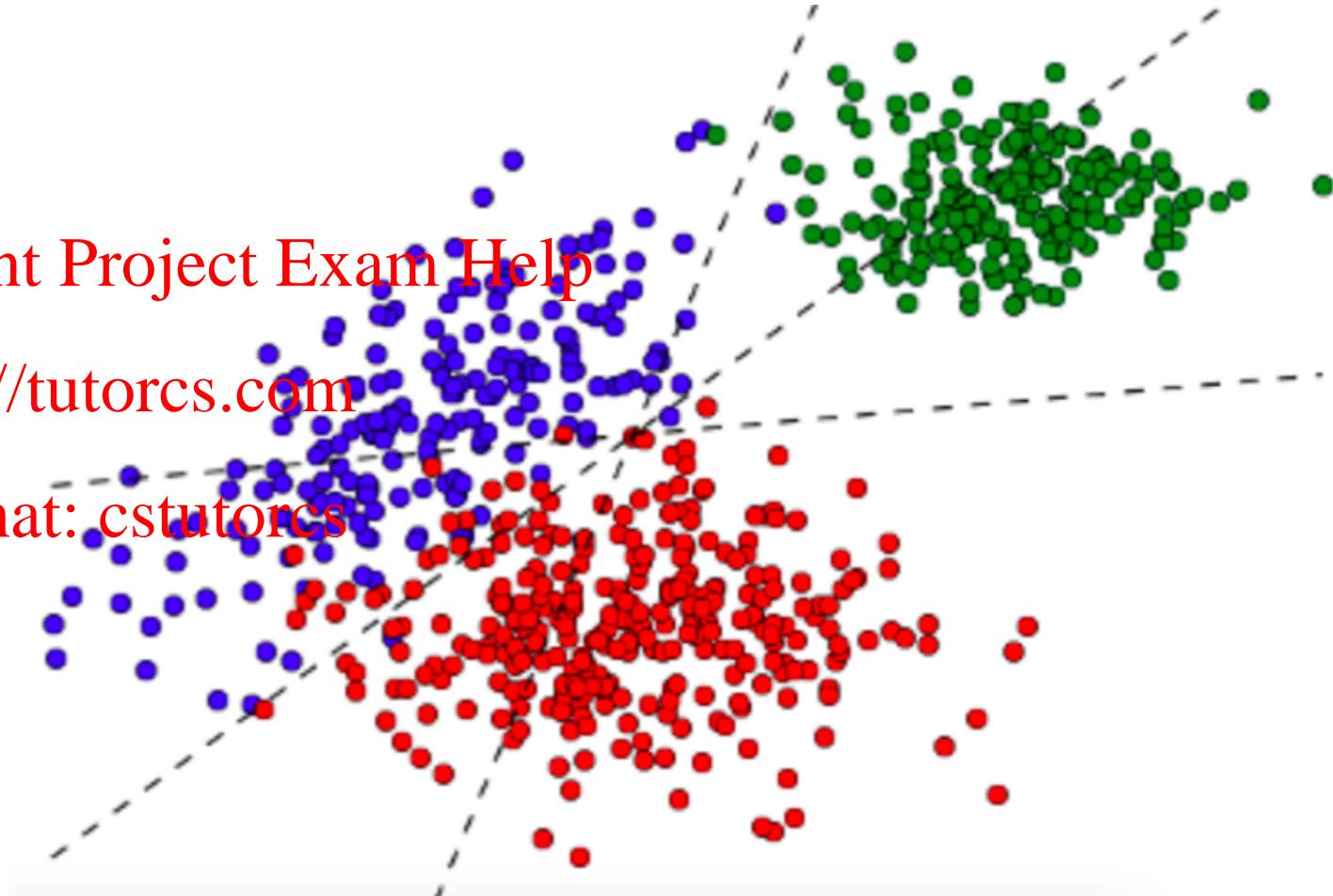
- Biometrics identification:

- each ID is a class label
- multi-class classification

Assignment Project Exam Help

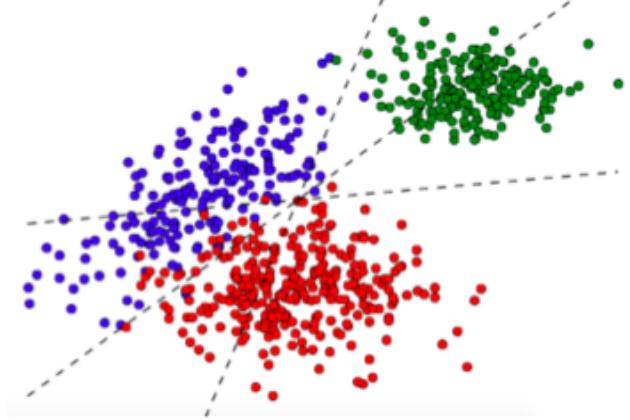
<https://tutorcs.com>

WeChat: cstutorcs



Biometrics as Multiclass Classification

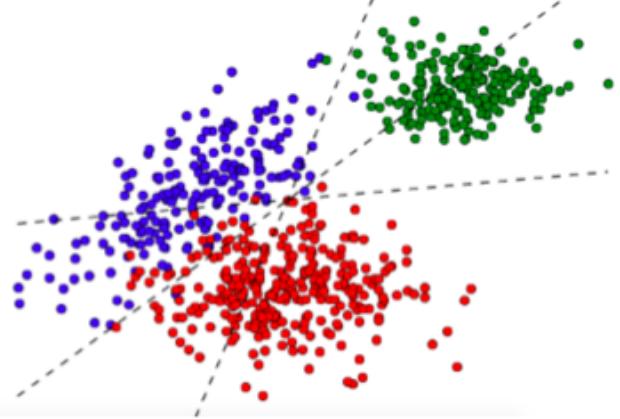
Machine Learning for Cybersecurity



- *Recap:*
 - biometrics ~~forgery~~ ~~Assignment Project Exam Help~~
 - biometrics considered by most to be highly unforgeable ~~<https://tutorcs.com>~~
 - thus much reliance & false sense of security due to this assumption
 - thus consequences ~~WeChat~~ ~~tutorcs~~ forged
 - anti-forgery solution: *liveness detection*
 - real vs forgery/fake
 - assumption: real has liveness, fake does not
 - *e.g. sweat, pulsation, capacitance, ...*

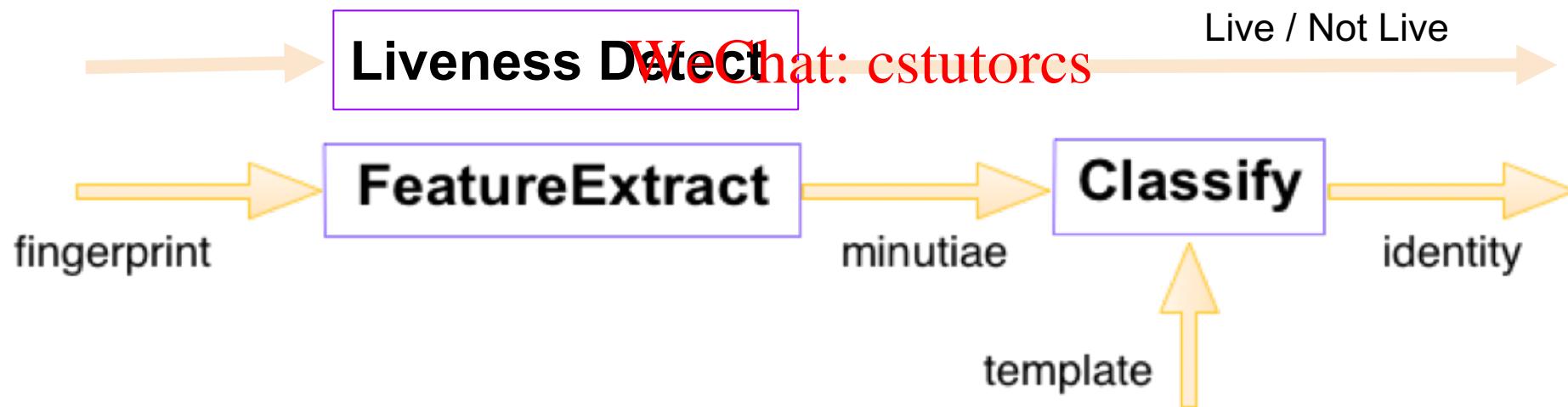
Biometrics as Multiclass Classification

Machine Learning for Cybersecurity



- counter anti-forgery attacks e.g. [Bowden-Peters et al. 2012]
 - Q: why liveness detection not work?
Assignment Project Exam Help
 - liveness separately detected from feature extraction

<https://tutorcs.com>



Spam Detection as Binary Classification

Machine Learning for Cybersecurity

- *Spam detection:*
 - *Spam vs Non-Spam* Assignment Project Exam Help
 - *datasets comprise the entries (label, text)*
 • where $label \in \{\text{non-spam}, \text{spam}\}$
 - *could use the Natural Language Processing (NLP)*
 • *breaks language into shorter pieces*

Classification: Empirical Tasks

Machine Learning for Cybersecurity

- separate dataset into **training set & test set**

Assignment Project Exam Help

- **train**/build the classifier

<https://tutorcs.com>

- let it learn from seeing labelled samples from the training set,
learn the decision boundary

WeChat: cstutorcs

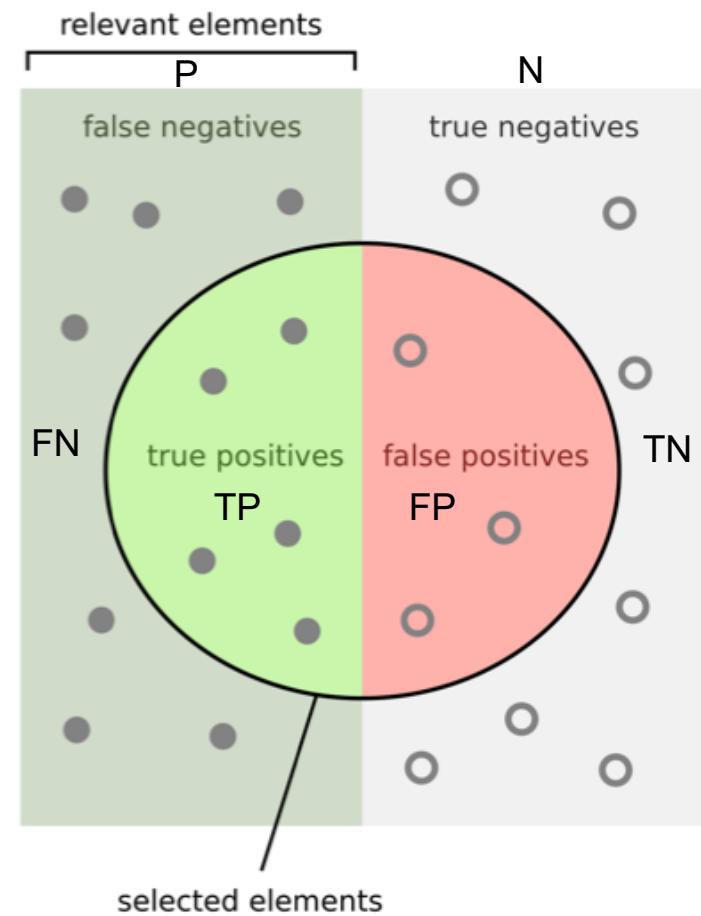
- **test** the classifier

- ...

Classification: Metrics

Machine Learning for Cybersecurity

- **test** the classifier
 - give it the samples from the test set, without the labels, it should **predict** what the label is for each sample
 - compute the **metrics**:
 - **accuracy** = #correct / #totalSamples
= $(TP+TN) / (P+N)$
 - **precision** = TP / all P predictions = TP / (TP+FP)
 - **recall** = TP / all samples in that class = TP / (TP+FN)
(a.k.a. TP rate or sensitivity)



How many selected items are relevant?

How many relevant items are selected?

$\text{Precision} = \frac{\text{How many selected items are relevant?}}{\text{How many relevant items are selected?}}$

$\text{Recall} = \frac{\text{How many selected items are relevant?}}{\text{How many relevant items are selected?}}$

Classification: Metrics

Machine Learning for Cybersecurity

- test the classifier: compute the **metrics**

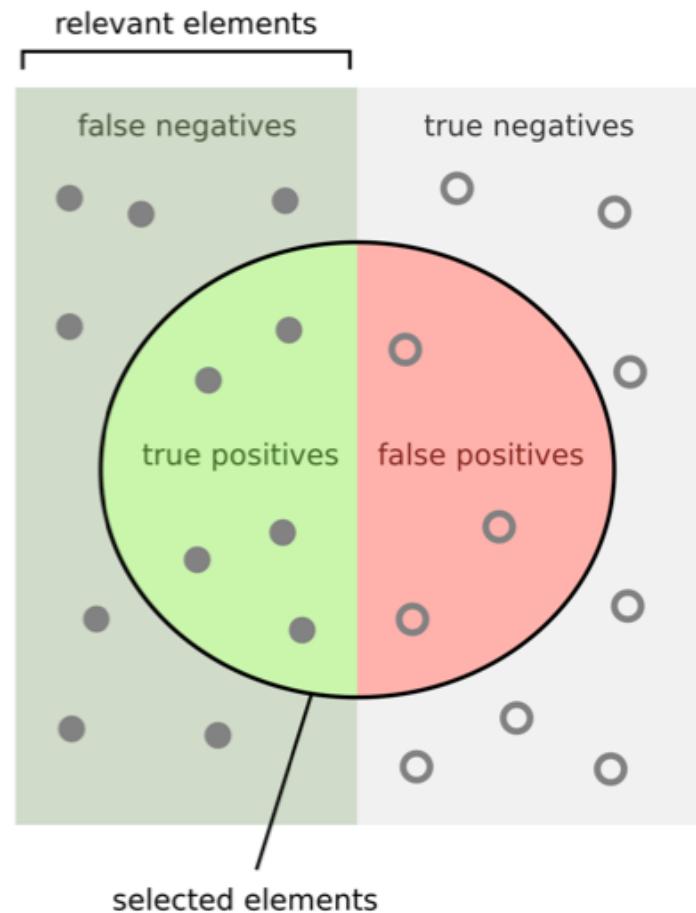
- accuracy
- precision
- recall
- **F1 score** = harmonic mean of precision & recall

Assignment Project Exam Help
<https://tutorcs.com>

$$H = \frac{n}{\frac{1}{x_1} + \frac{1}{x_2} + \frac{1}{x_3} + \dots + \frac{1}{x_n}}$$

WeChat: cstutorcs

$$\frac{2}{\frac{1}{\text{recall}} \times \frac{1}{\text{precision}}} = 2 \times \frac{\text{precision} \times \text{recall}}{\text{precision} + \text{recall}}$$



How many selected items are relevant?

$$\text{Precision} = \frac{\text{green}}{\text{green} + \text{red}}$$

How many relevant items are selected?

$$\text{Recall} = \frac{\text{green}}{\text{green} + \text{grey}}$$

Classification: Metrics

Machine Learning for Cybersecurity

- Recall last lecture's alternative classification metrics
- (commonly used in biometrics)

- False Accept Rate (FAR) = $FP / \text{all } N \text{ samples} = FP / (TN + FP)$

- compare with (1-precision) = $FP / (TP + FP)$

Assignment Project Exam Help

- False Reject Rate (FRR) = $FN / \text{all } P \text{ samples} = FN / (FN + TP)$

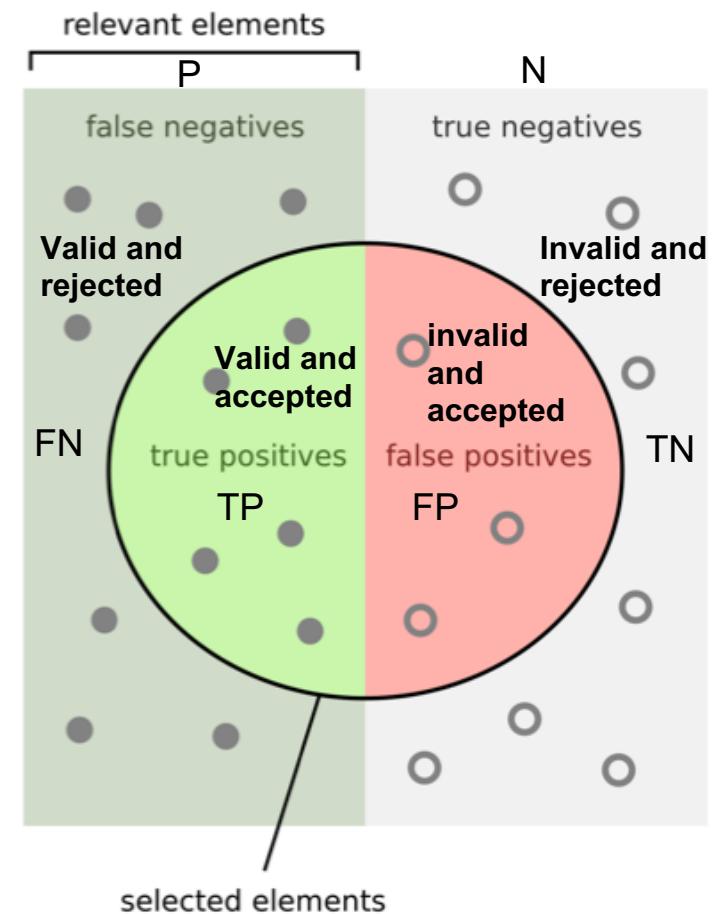
- compare with (1-recall) = $FN / (FN + TP)$

WeChat: cstutorcs

Note: above definitions assume

P = "positive" = "valid", N = "negative" = "invalid" (see next slide)

QUESTION: Which metric do you think is more relevant for security classification (assume P = "valid", N = "invalid"):
FAR or (1-precision)? Why?



How many selected items are relevant?

$$\text{Precision} = \frac{\text{TP}}{\text{TP} + \text{FP}}$$

How many relevant items are selected?

$$\text{Recall} = \frac{\text{TP}}{\text{TP} + \text{FN}}$$

A note on definitions

- Definition of “positive” (P) / “negative” (N) can depend on the scenario / context
 - can be defined explicitly in each scenario to avoid confusion.
- However, ~~Assignment Project Exam Help~~:
 - In security scenarios, the following natural convention for the meaning of Accept and Reject is usually followed:
 - Accept = gets through system = (if no errors) honest/valid/good
 - Reject = gets blocked = (if no errors) malicious/invalid/bad
 - Therefore FAR and FRR error rates can be defined generally as:
 - **False Accept Rate, FAR** = (# of malicious samples falsely classified as accept)/(total # of malicious samples)
 - **False Reject Rate, FRR** = (# of honest samples falsely classified as reject)/(total # of honest samples)

Spam Detection as Binary Classification

Machine Learning for Cybersecurity

- feature learning

- bag of words (BoW)
 - word & frequency

e.g.

(1) John likes to watch movies. Mary likes movies too.

Assignment Project Exam Help

BoW =

{"John":1,"likes":2,"to":1,"watch":1,"movies":2,"Mary":1,"too":1};

<https://tutorcs.com>

WeChat: cstutorcs

(source: Wikipedia)



Spam Detection as Binary Classification

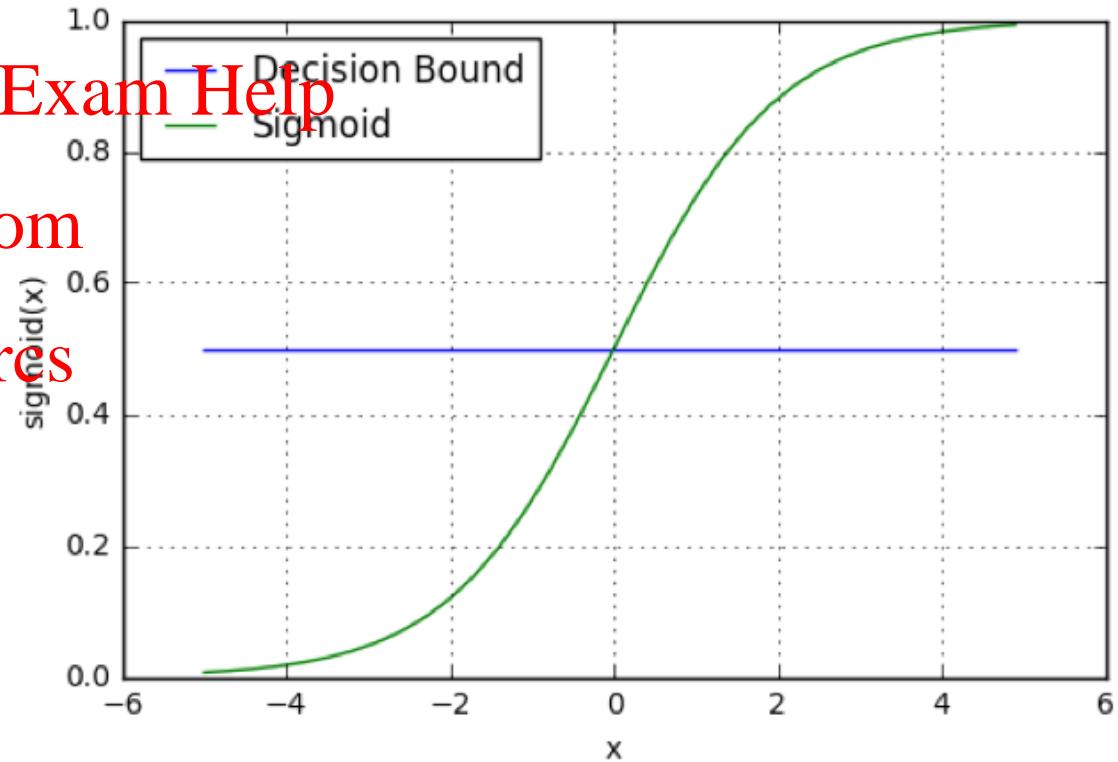
Machine Learning for Cybersecurity

- classification
 - logistic regression classifier
 - features (in this case, BoW, i.e. frequency count) are input to sigmoid function
 - $\sigma(x) = \frac{1}{1 + e^{-x}}$
 - outputting a value between 0 and 1 i.e. probability p
 - $p > 0.5 \rightarrow \text{class} = 1, \text{else class} = 0$

Assignment Project Exam Help

<https://tutorcs.com>

WeChat: cstutordes





Cryptanalysis as Binary Classification

Machine Learning for Cybersecurity

- *Recap:*
 - CONFidentiality problem aims to safeguard secrecy of message m
<https://tutorcs.com>
- How? ensure only specific parties can access
 - lock it up: unlock only if have key
 - transform it: reverse transform needs a key
 - encryption



Encryption for CONF?

Machine Learning for Cybersecurity

- Encrypt? $c = E(m, k)$
 - m = message a.k.a. plaintext p
 - k = secret key
 - c = output ciphertext
- Q: *since m needs to be CONF, how to ensure m remains CONF although c can be seen by anyone?*





Encryption for CONF?

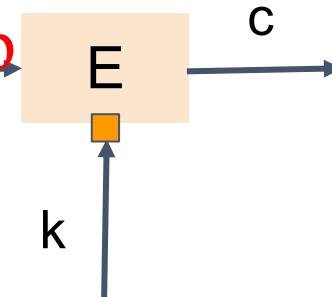
Machine Learning for Cybersecurity

- Encrypt? $c = E(m, k)$
 - m = message a.k.a. plaintext p
 - k = secret key
 - c = output ciphertext
- Q: *since m needs to be CONF, how to ensure m remains CONF although c can be seen by anyone?*
- A: reversing E needs the k , only knowing k can E be reversed
- the end goal = CONF of m , the means = CONF of k

Assignment Project Exam Help

<https://tutorcs.com>

WeChat: cstutorcs

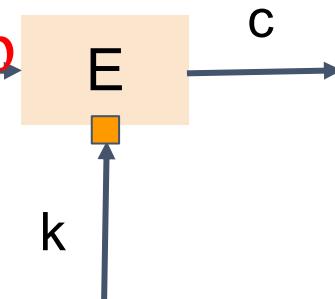




Breaking Encryption CONF

Machine Learning for Cybersecurity

- the end goal = CONF of m ,
the means = CONF of k
Assignment Project Exam Help
- guessing k leads to breaking CONF of m
https://tutorcs.com
- cryptanalysis = security analysis of crypto techniques e.g. encryption
WeChat: cstutorcs
- basic cryptanalysis of encryption = binary classification problem
 - we'll discuss a few variations of cryptanalysis problems, any of these would indicate some weakness in the encryption design

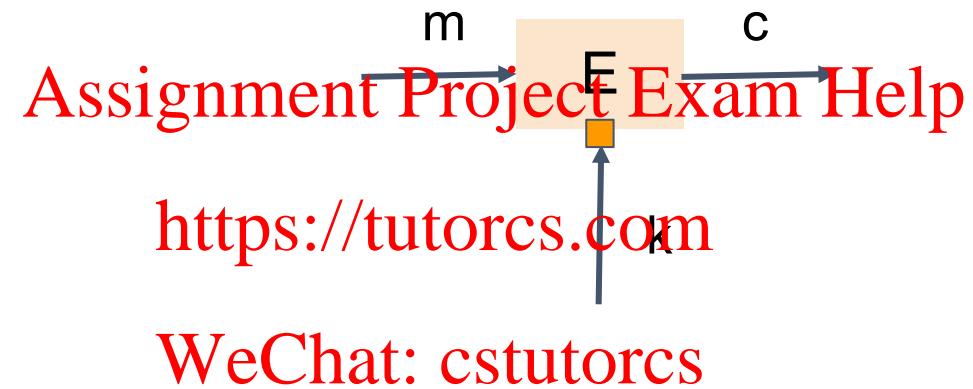




Cryptanalysis of Encryption: Case I

Machine Learning for Cybersecurity

- User
 - chooses secret k



- Attacker
 - sees c
 - wants to guess k

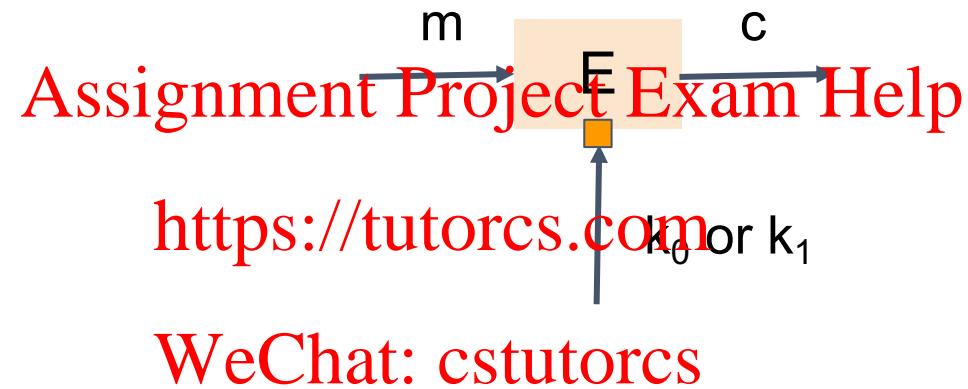
- Case I: **key-recovery** attack



Cryptanalysis of Encryption: Case II

Machine Learning for Cybersecurity

- User
 - flips a coin {0,1}
 - if 0: use k_0
 - if 1: use k_1
- Case II: **key-distinguishing** attack
- Q: *which is easier for attacker to do? Case I or II?*



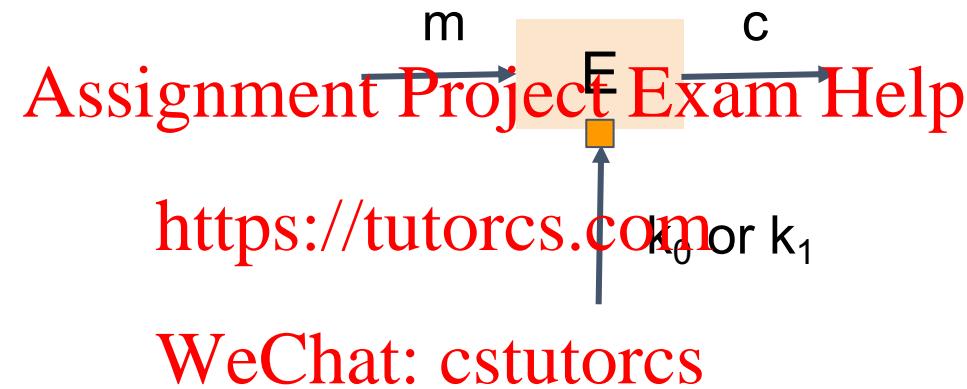
- Attacker
 - sees c
 - wants to guess if k_0 or k_1



Cryptanalysis of Encryption: Case II

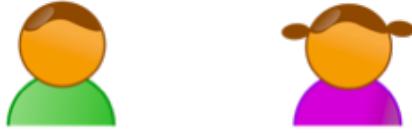
Machine Learning for Cybersecurity

- User
 - flips a coin {0,1}
 - if 0: use k_0
 - if 1: use k_1



- Q: *which is easier for attacker to do? Case I or II?*
- A: Case II is easier, just guess 1 of 2 options, vs Case I i.e. guess all bits of k

- Attacker
 - sees c
 - wants to guess if k_0 or k_1

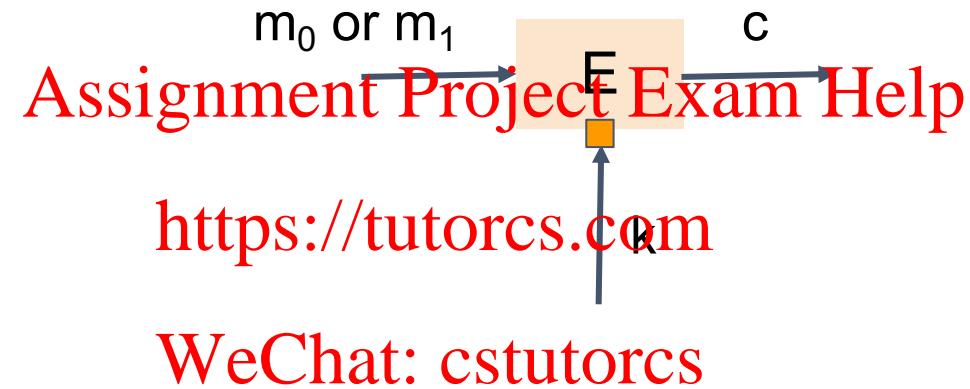


Cryptanalysis of Encryption: Case III

Machine Learning for Cybersecurity

- User

- flips a coin {0,1}
 - if 0: use m_0
 - if 1: use m_1



- Case III: **plaintext-distinguishing** attack

- Case II and Case III would be similarly difficult for the attacker since it is guessing 1 of 2 options

- Attacker

- sees c
 - wants to guess if m_0 or m_1

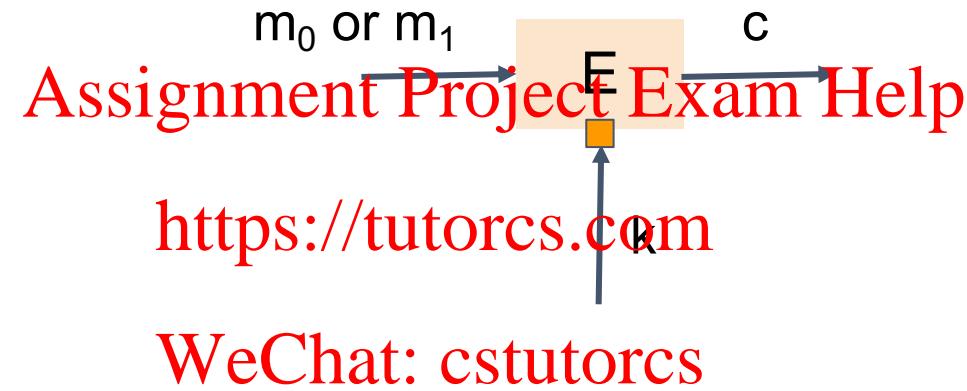


Cryptanalysis of Encryption: Case III

Machine Learning for Cybersecurity

- User

- flips a coin {0,1}
 - if 0: use m_0
 - if 1: use m_1



- Case III: **plaintext-distinguishing** attack

- Q: *what weakness or problem could cause an attacker to be able guess if it is m_0 or m_1 ?*

- Attacker

- sees c
 - wants to guess if m_0 or m_1

Further Reading

- A. Polyakov, "Machine Learning for Cybersecurity", article available at:
<https://towardsdatascience.com/machine-learning-for-cybersecurity-101-7822b802790b>
Assignment Project Exam Help
- This article discusses most of the topics we discussed, and provides links to further references. <https://tutorcs.com>

Note: Our focus in this week's lecture is on a high-level overview on the types of ML techniques and their potential applications in cyber security. We do not expect students to learn the details of how ML algorithms work – this is outside the scope of this unit.