# Malware Defense II
## TDDD17 – Information Security, Second Course

Alireza Mohammadinodooshan

Department of Computer and Information Science
Linköping University

**LiU** LINKÖPING
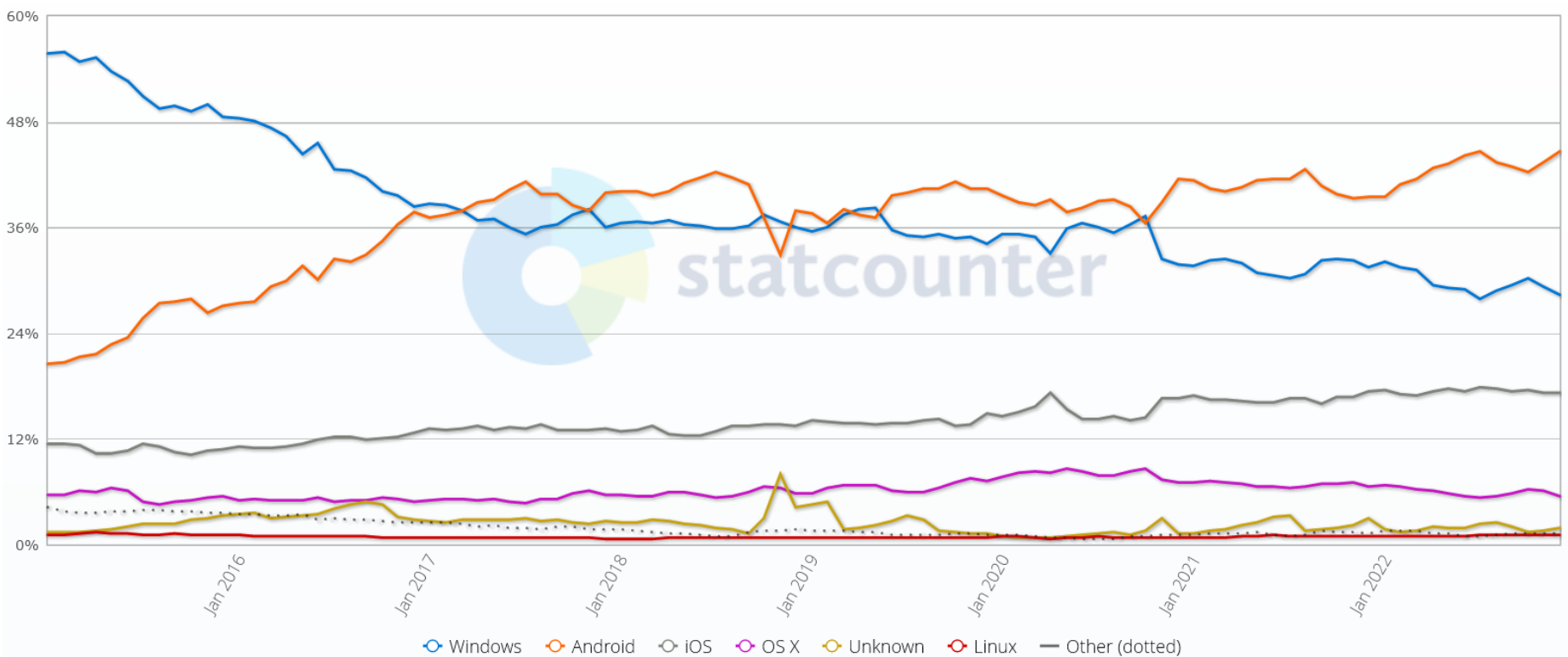UNIVERSITY

# What Has Been Covered …

- Malware basics
  - Different types of functionality
  - Different infection Methods
- AV cat and mouse game
  - Signatures based detection
  - More complex signatures and static heuristics
  - Static unpacking and emulation
  - Cloud-based detection
  - Machine learning detection

LINKÖPING UNIVERSITY

# Agenda

- Mobile malware

    - Specific challenges

    - Specific risks

    - Security models and their effect on malware detection

        - iOS

        - Android

    - Detection countermeasures

- Machine learning for malware detection

    - Motivation

    - Terminology

    - Learning types

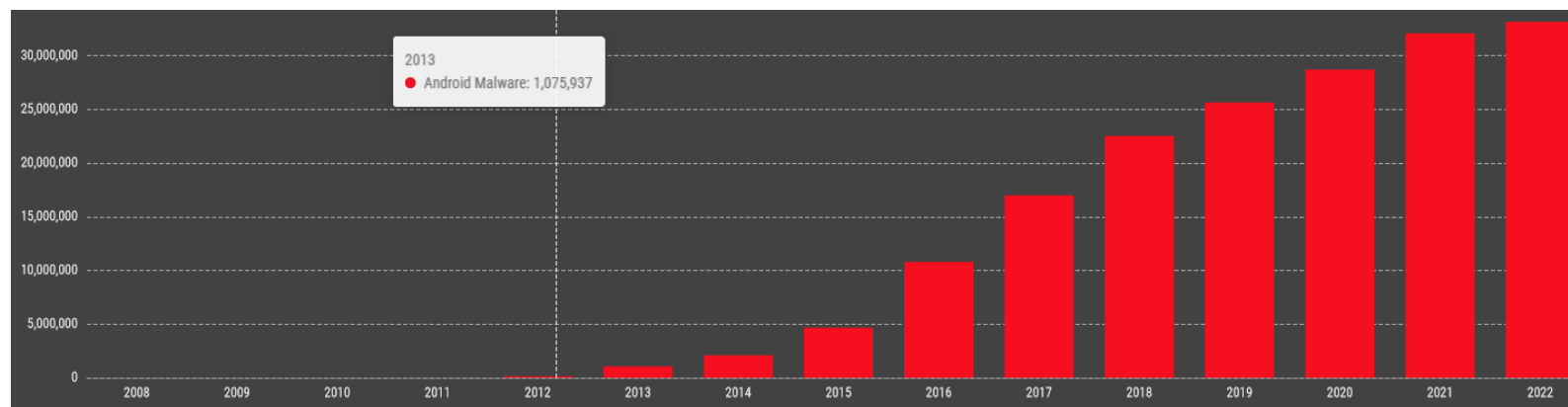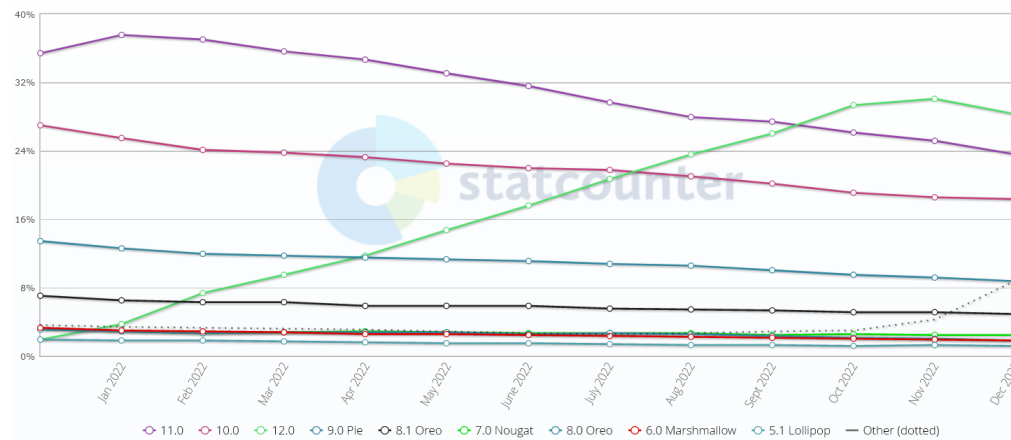    - Machine learning-based malware detection challenges

**LINKÖPING UNIVERSITY**

# Motivation

- 5.48 billion smartphone users in the world in 2022

# Motivation

- Using old versions of android

- It is not surprising that Mobile platform became an appealing target for malware authors.





https://portal.av-atlas.org/malware/statistics

https://gs.statcounter.com/android-version-market-share/mobile-tablet/worldwide

# Mobile Malware Definition

- Malicious software designed to attack mobile devices
  - Phone
  - Tablet
  - Watch
  - TV

LINKÖPING
UNIVERSITY

# Samples of Mobile Malware

- iOS stock
  - PawnStorm.A
    - Able to upload GPS location, contact list, photos to a remote server.
  - YiSpecter
    - Able to download, install and launch arbitrary apps
- Android
  - Android/Filecoder.C
    - Able to spread via text messages and contains a malicious link. Encrypts all of your local files in exchange for a ransom between $94 and $188.
  - Plankton
    - Communicates with a remote server, downloads and install other applications and sends premium SMS messages

# Mobile Malware Specific Challenges

1. Lots of users
   - Botnets

2. More personalized and privacy concerns
   - Banking info
   - Personal Photos
   - Contact info

3. Widespread access to networks
   - 4G
   - Wifi
   - Bluetooth

# Mobile Malware Specific Challenges

4. Less computation power

   – Limited capabilities for on-device detection


5. Almost exclusively trojans

   – Repackaged apps

     • It is easier to reverse engineer Android apps

     • A very simple technique is to replace the advertisement logic and re-bundle and publish the app

   – Fake apps also exist!

LINKÖPING UNIVERSITY

# Mobile Malware Specific Challenges

6. Due to limited computation power, most of the trust in apps is moved to app stores to analyze the apps

   - While for the 3rd party stores and somehow even for the google play store, this is a mistrust (we will elaborate on this …)

   - Attackers also have the motivation to deliver their malware through stores (official or third party)

7. As the Android's kernel is open source, attackers have a better understanding of its vulnerabilities if they exist

LINKÖPING UNIVERSITY

# Mobile Malware Specific Challenges

8. Harder to detect with 3$^{rd}$ party AV on the device compared to PC malware due to stronger isolation between apps

   – Memory isolation

   – User isolation

   - Each app is treated as a separate user

   - Applications cannot interact with each other, and they have limited access to the system as well as other apps resources

# Mobile Malware Risks

- System damage
  - Battery draining
  - Cryptocurrecy mining
  - Disabling system functions
    - Block calling functionality

- Economic
  - Sending SMS or MMS messages to premium numbers
  - Dialing premium numbers
  - Deleting important data

Peng, S., Yu, S., & Yang, A. (2013). Smartphone malware and its propagation modeling: A survey. IEEE Communications Surveys & Tutorials, 16(2), 925-941

# Mobile Malware Risks

- Information leakage
  - Privacy
  - Stealing bank account information

- Disturbing mobile networks
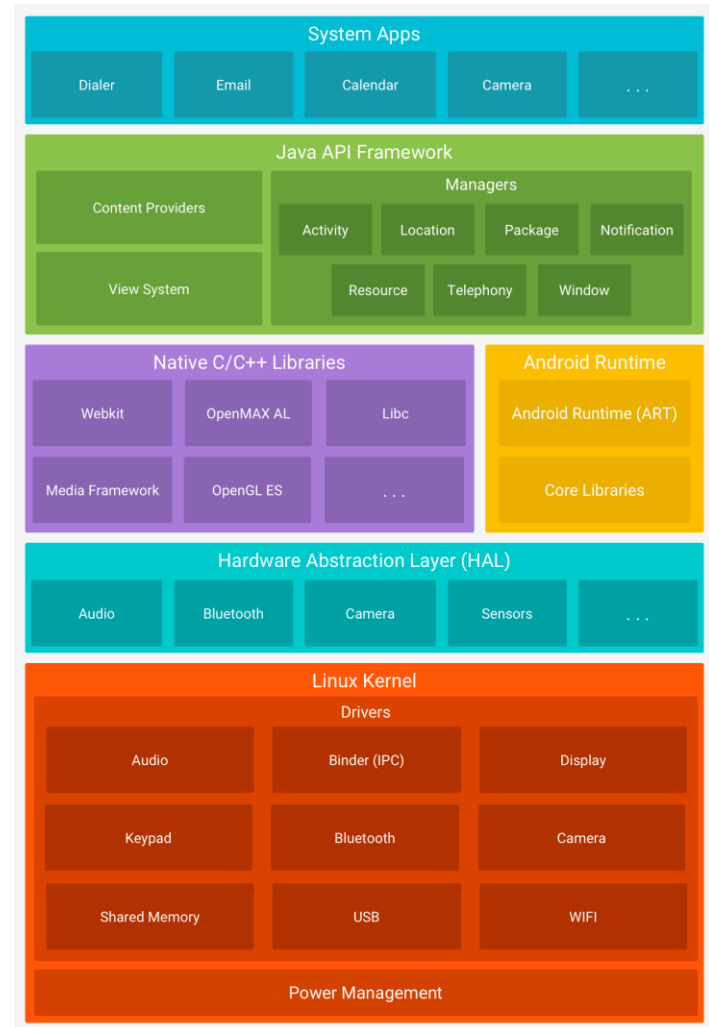  - Denial-of-service (DoS)

# iOS Security Model

- ## System Security

  - Startup and updates are authorized

- ## Data security

  - File-level data protection uses strong encryption keys derived from the user's unique passcode.

- ## App security

  - Application run in their sandboxes.

  - More important than this …

https://developer.apple.com/app-store/review/
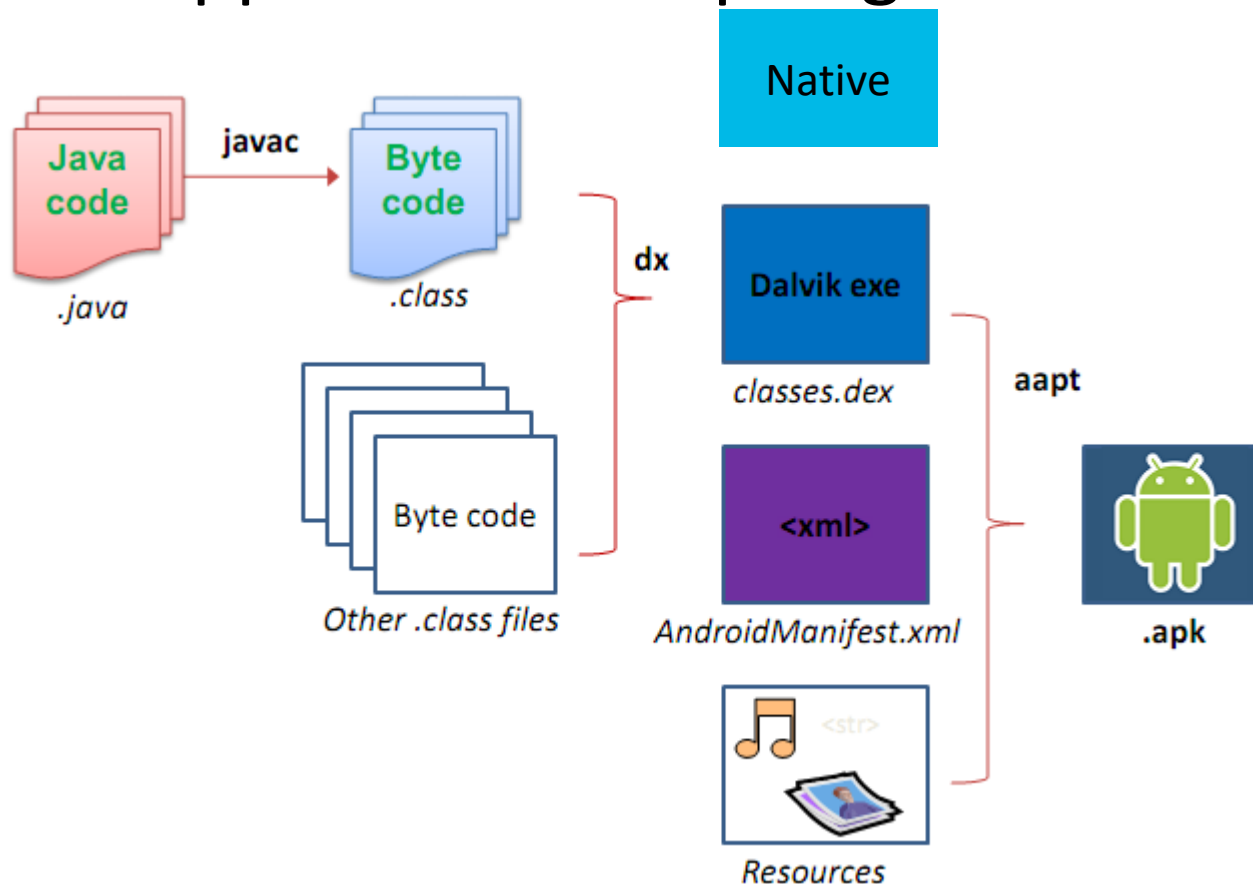
# iOS Security Model

- Before releasing on store, they go through a strict vetting process

  – Manual testing

  – Static analysis

  – Apps can not do actions outside of what they claim

LINKÖPING
UNIVERSITY

# Android Architecture

- ## Hardware Abstraction Layer (HAL)

    – provides standard interfaces that make the
    device hardware capabilities available to the
    higher-level Java API framework.

- ## Android Runtime

    – For new Android devices, each app runs in its
    own process and with its own instance of the
    Android Runtime (ART).  Before ART, the
    Dalvik VM has been used

- ## Native C/C++ Libraries

    – It is possible to have compiled c/c++ code
    packaged with an Apk which can be called
    through Java Native Interface (JNI)



https://developer.android.com/guide/platform

# Android Application Compiling

https://justamomentgoose.wordpress.com/2013/06/04/android-started-note-2-android-file-apk-decompile/

# Androidmanifest.xml

- Provides the essential information to the Android system regarding this app
  - Minimum android API
  - Linked libraries
  - Components, activities, services, …
  - Required permissions

LINKÖPING UNIVERSITY

# Android  Security Model

- Application Sandboxing

  – Android automatically assigns a unique UID to each app at installation

  – App is allowed to access :

    - Own files

    - World-accessible resources

  – More access :

    - Managed through defining in the androidmanifest.xml

      – <uses-permission android:name="android.permission.READ_PHONE_STATE " />

LINKÖPING UNIVERSITY

# Android  Security Model

- Vetting process

  - Does not require an exhaustive app vetting process

    - More lenient comparing to iOS

  - Apps are dynamically tested with a Google security service known as Bouncer.

    - The results are combined with the output coming from the google reputation system

  - Researchers have shown the feasibility of fingerprinting Bouncer

    - Android ID.

    - phone number

    - ….

    - Check Jon Oberheide and Charlie Miller's work

  - Malware may be able to bypass Bouncer

    - They have the motivation to bypass it because they want not to be detected by it. So if they detect that they are running in bouncer they do not show their actual behavior

LINKÖPING
UNIVERSITY

# Mobile Malware Detection

- Static Code Analysis

  - Signature-Based Technique

    - Specific strings or patterns in the byte code

    - Extracting the strings is straightforward

  - Permission-Based Technique

    - Analysing the requested permissions to identify the potential malware samples

  - Dalvik Bytecode-Based Technique

    - Analysing the byte code to identify the malicious Android samples(API calls, data flows,…)

LINKÖPING UNIVERSITY

# Mobile Malware Detection

- Dynamic Behavior Analysis
  - Sequence of system calls
  - Accessed files

- Hybrid Analysis

# Malware Detection Countermeasures

- Static
  - Obfuscation
    - Making the byte code hard to understand
    - Making signature or even some static heuristics-based analysis harder
  - Packing
- Dynamic
  - Sandbox detection
    - Many of the sandboxes still do not have real device behaviors
      - E.g. do not support GPS or do not have a real GPS accuracy

LINKÖPING UNIVERSITY

# Obfuscation

- Identifier renaming
  - garble the key identifiers used in their source code. e.g., 'a', 'b', 'aa', 'ab', 'ac'

- String encryption
  - Replacing the constant strings in the dex file with their encrypted form and adding the code to decrypt them on the fly

- Control flow obfuscation: changing the logical flow of the program

  - Injecting dead code

  - Re-ordering statements

  - Inserting opaque predicates.

    - It is always true or false

    - Malware author knows this

    - But it is hard for the analyst to follow and find the value

```
obj = benign()
var1 = 10
var2 = [var1 for i in range(10)]
if var1 == var2[0]:
    obj = malware()
obj.load()
```

LINKÖPING
UNIVERSITY

# Packing

# Machine Learning for Malware Analysis

# Why?

- Creating detection rules (signatures) manually couldn't keep up with the emerging flow of malware.

  – Zero-day malware

- Need a more reliable method when we know that the relation between the app features is hard to find for the human

- Sometimes we need a triage method

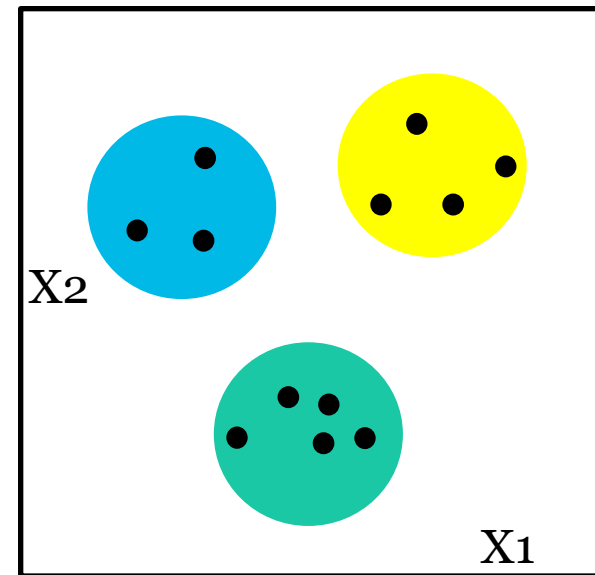  – A procedure we use to prioritize the apps that should be examined

# Machine Learning

- Machine learning is a set of methods that gives computers the ability to learn without being explicitly programmed
  - Learning from the data
  - It is used when we want to (explicitly or implicitly) learn the relation using some available data (known as training data)

# Terminology

- (Predictive) Model: The hidden relation

- Training data: Data based on which we make the model

- Testing data: Data based on which we evaluate the model

- (Hidden relation) Learning types:
  - Unsupervised
  - Supervised

# Unsupervised learning

- Given X(X1 and X2 in the following figure)

- The goal is to discover the structure of the data

  - Clustering: splitting a data set into groups of similar objects

    - Application example

      » Malware family detection

# Supervised Learning

- Having both X (X1 and X2 in the following fig)  and y (the colors in the following fig)

  – we try to find the relation between them(X and y)

- Application

  – Malware detection

    • X: features of malware and benign apps

    • Y: malware or benign label

- Can be classified to

  – Discriminative methods

    • Need samples from both classes

  – Anomaly detection

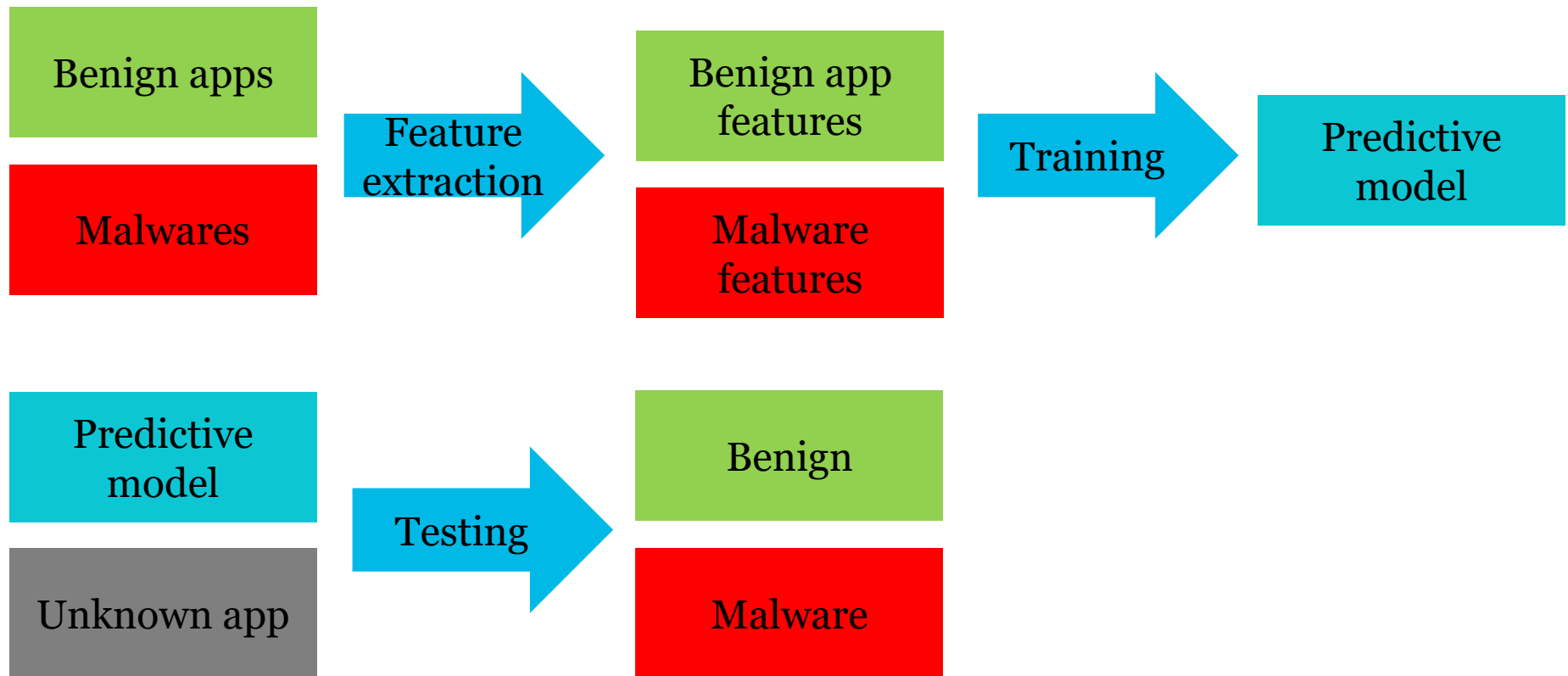    • Need samples from one class (The benign class)

    • To detect deviation from the "normal" class

# Classification vs. Regression

- Classification

  – When y is a discrete variable

    - Malware or benign


- Regression

  – When Y is a continuous variable

    - Probability of belonging to a specific malware family (e.g., can be used for triaging the app)

# ML-based Malware Detection procedure

- Collecting training data

- Extracting features from training data

- Training the model: finding the model

- Testing (Evaluating) the model

# ML-based Malware Detection Workflow

# Collect Training data

- Dataset should be representative of real-world malware
  - Example of bad practice
    - Suppose that we collected some benign and malware apps but all the malware samples we collected have a size between 1-2 megabytes which is not representative of real malware
    - The model overfits to this unrealistic pattern
    - The model will have a large false positive
      - Can have disastrous effects as we have seen in the previous lecture

# Extracting Features

- The extracted features should be relevant.

- Usually, a domain knowledge helps a lot here
  - Examples
    - PC
      - The header values
    - Mobile
      - Set of Privileges
    - Both
      - Obfuscation status
  - Feature selection methods can be used to limit the number of features
    - For example, low variance features can be removed

LINKÖPING
UNIVERSITY

# Training

- Most of the models have some parameters which during the training phase are optimized using the training data

  – This optimization happens based on a particular metric.

  – This particular metric is usually the classification or regression error
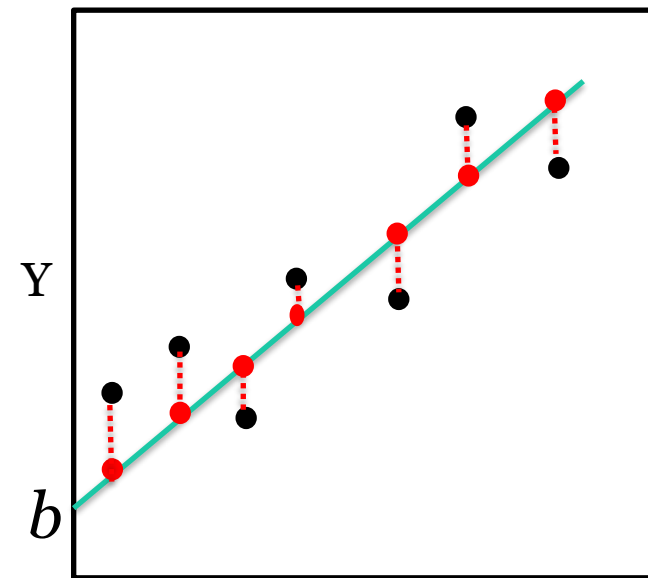
# Training(Example)

- Linear regression
  - We have a set of $(X_i, Y_i)$ *training points*
  - We want to find the regression line
    - Which with the least error estimates the points
    - $F = aX + b$
      - $a_{opt\ ?}$
      - $b_{opt\ ?}$

# Training(Example)

- Learning workflow
  - For each point $X_i$ compute the response $F_i$
    - $F_i = aX_i + b$
  - Compute $ERR_{tot} = SUM((F_i - Y_i)^2)$
  - Now we can compute $a_{opt}$ and $b_{opt}$
    - Which minimizes $ERR_{tot}$
      - Closed form
      - Optimization
- This was a regression example
  - For the classification, for example
    - We can find the discriminative line or hyperplane between the points

# Testing

- By finding the optimal values of parameters ( in this case $a$ and $b$) we test it on testing data.

  - To see whether it can generalize to unseen data

  - Or it has just memorized the training data

- In this case (Testing) we will also have some error

  - We train a model by minimizing its error on the training data

  - The training error is different from the testing error

  - This testing error value is computed on test data.

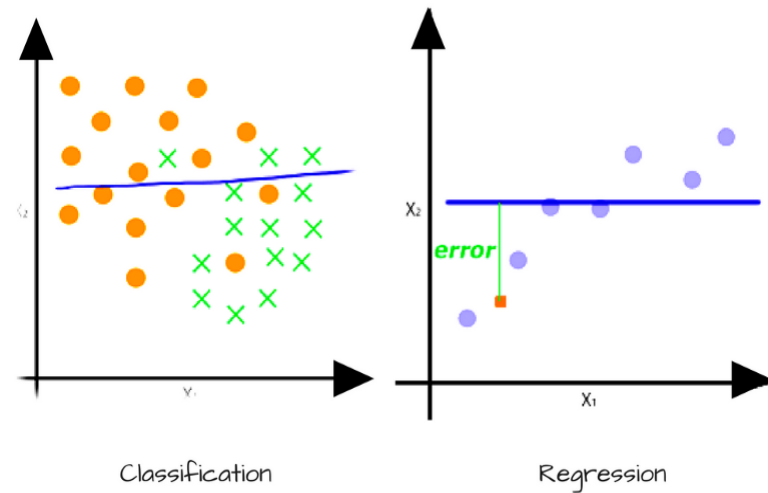# Machine Learning-based Malware Detection Challenges

- Under- and Over-fitting

- Imbalanced datasets

- Performance evaluation measures
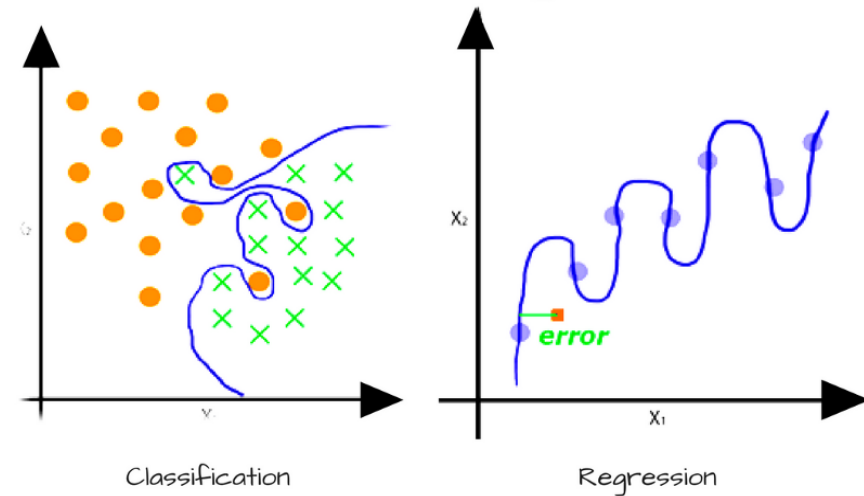
- Dataset quality

# Underfitting and overfitting

- Underfitting
  - The model is unable to obtain a low error even on the training set.


- Overfitting (Memorization)
  - The training error is small enough but not the testing error.
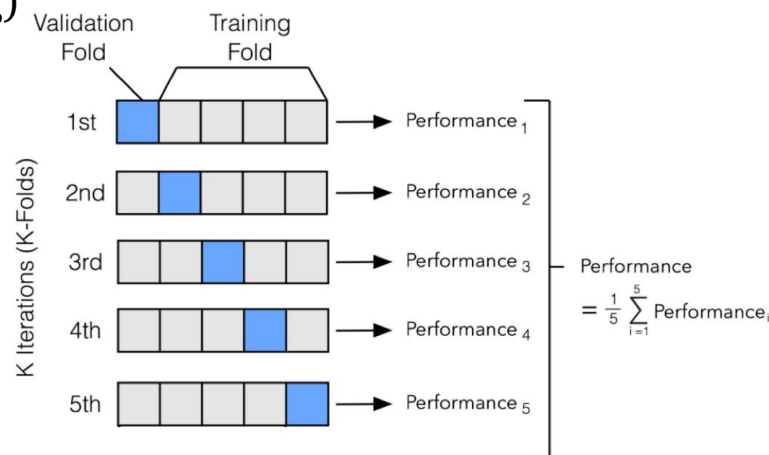
# Underfitting and overfitting

# Solution : Cross-Validation

- Basic idea
  - Each observation in our dataset has the opportunity of being tested

- Procedure
  - we divide the dataset into k sets
  - For k rounds, we go over the dataset, and in each round:
    - one part is used for validation (Testing)
    - Remaining parts used for training
    - Based on the total performance value we can select the optimal model

# The problem of imbalanced datasets

- Malware datasets are usually imbalanced

- Suppose that we have a dataset in which 99 percent of samples are benign

  – Now a naïve malware detection classifier which classifies all the samples as being benign reaches an accuracy of 99 percent

  – Probably no other model can reach this optimal accuracy

  – But is accuracy a good metric to train the model on ?

  – Evidently not. This model can not detect any malware!

- We need to focus on some other performance measures!

LINKÖPING UNIVERSITY

# Performance Measures



- Accuracy

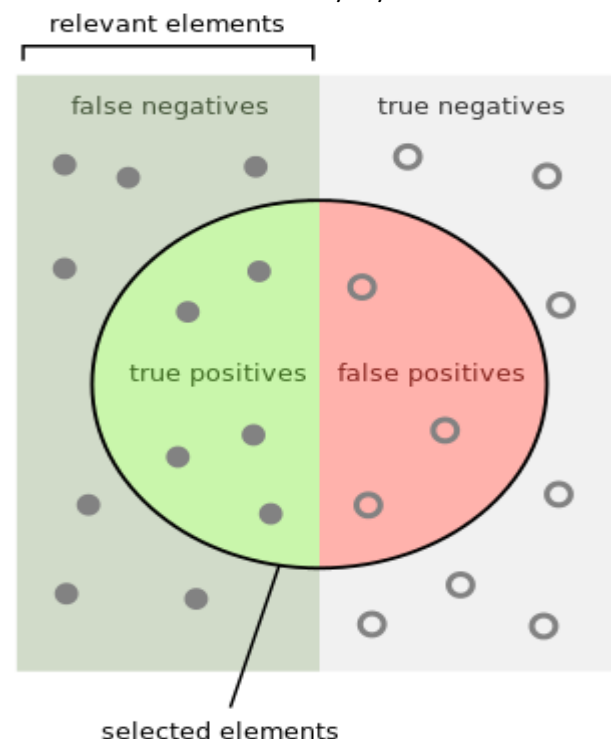    $$\frac{TP+TN}{TP+FP+TN+FN}$$

- Recall (Sensitivity)

    $$\frac{TP}{TP+FN}$$

- Precision

    $$\frac{TP}{TP+FP}$$

- F-score : F-Score is the weighted average of Precision and Recall.

    $$\frac{2 * precision * recall}{precision+recall}$$

https://en.wikipedia.org/wiki/Precision_and_recall

# Dataset quality

- Having a representative dataset is crucial for machine learning methods.

    – Recall the bad practice for data collection

- It is not possible to train the models on the end points

    – We cannot collect representative data there!

- The training is done on the cloud

# Summary

- We motivated the need for mobile malware detection

- We discussed Mobile malware specific challenges

    – Lots of users, privacy concerns, Widespread access to networks, ...

- Mobile malware risks were reviewed

    – System damage, Economic risks, ....

- We reviewed the security model of iOS and Android

    – We discussed the differences between iOS and android vetting processes

    – We have seen how it effects the malware prevalence in each platform

- We have reviewed different techniques for mobile malware detection

    – Static, dynamic, hybrid

- Obfuscation techniques were reviewed

# Summary

- The role of machine learning in malware detection

- Different learning types

    – Supervised

    – Unsupervised

- ML-based Malware Detection procedure

    – Collecting training data

    – Extracting features from training data

    – Training  the model

    – Validating the model

- Machine Learning-based Malware Detection Challenges

    – Under- and Over-fitting

    – Imbalanced datasets

    – Performance evaluation measures

    – Dataset quality

**II.U** LINKÖPING
UNIVERSITY