# Vulnerability Prediction in Android Apps
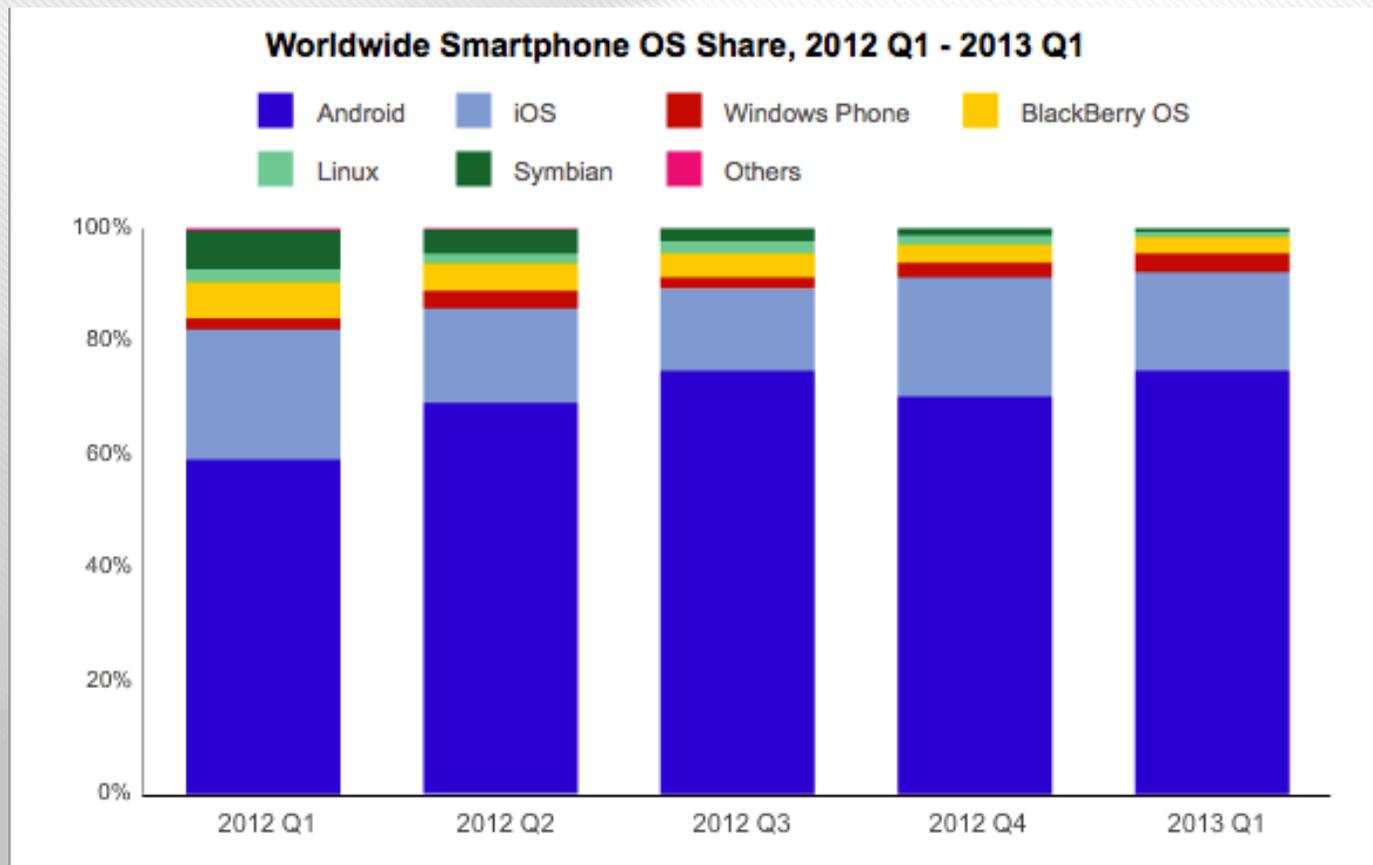
**Aram Hovsepyan[1], Riccardo Scandariato[1], James Walden[2,] Viet Hung Nguyen[3] Wouter Joosen[1]**

[1] iMinds-DistriNet, Katholieke Universiteit Leuven
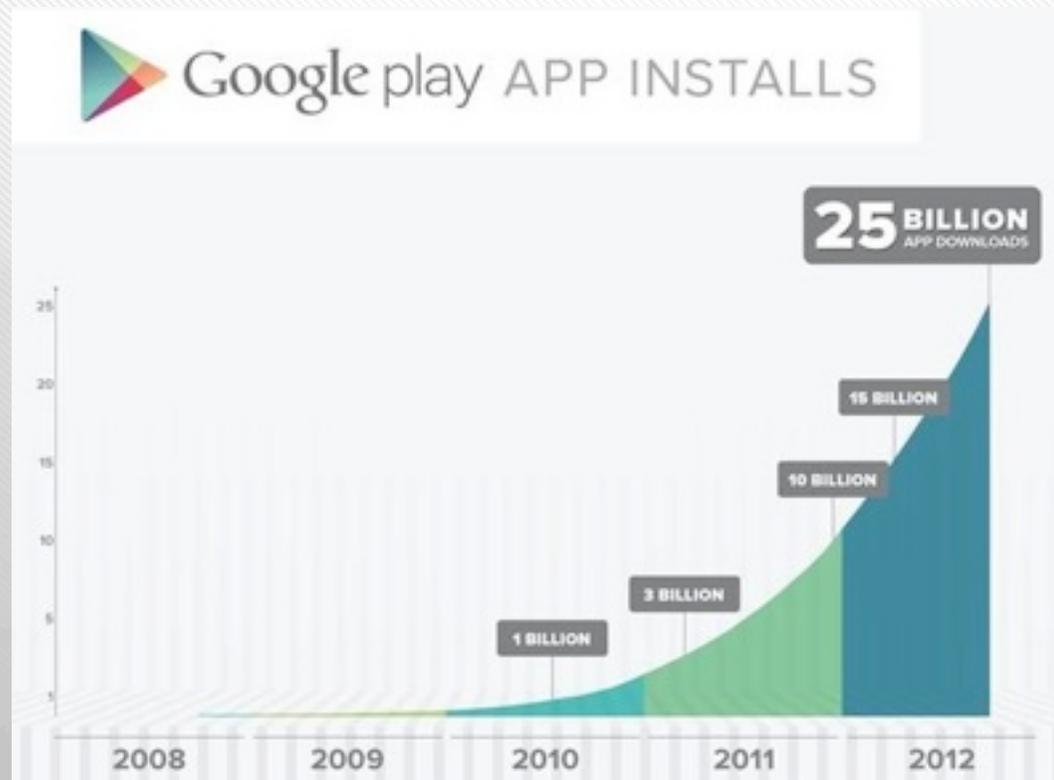[2] Northern Kentucky University, [3] University of Trento

Monday 10 June 13 week

# Android apps are an attractive target

## Android has 75% market share as of Q1 2013 [IDC]

**Worldwide Smartphone OS Share, 2012 Q1 - 2013 Q1**

Legend:
- Android
- iOS
- Windows Phone
- BlackBerry OS
- Linux
- Symbian
- Others

Monday 10 June 13 week

# Android apps are an attractive target

Google play has over 775K apps and over 48B total installs [IDC, Google I/O keynote]

Monday 10 June 13 week

# Android apps are an attractive target

App security is not guaranteed by the platform provider

- ➲ Apps that are well intended, but not exploit free

A single vulnerability could affect a massive number of users

Not yet much explored

- ➲ Focused on Mozilla Firefox / RHEL

Monday 10 June 13 week

# How to find vulnerabilities?

Monday 10 June 13 week

# How to find vulnerabilities?

## Code inspection

- ⮎ Manual verification is not feasible

- ⮎ Not all apps can afford security experts

- ⮎ Even security experts cannot analyze every line of code

Monday 10 June 13 week

# How to find vulnerabilities?

## Code inspection

- ⮌ Manual verification is not feasible
- ⮌ Not all apps can afford security experts
- ⮌ Even security experts cannot analyze every line of code

## Penetration testing / security testing

Monday 10 June 13 week

# How to find vulnerabilities?

## Code inspection

- ⮌ Manual verification is not feasible
- ⮌ Not all apps can afford security experts
- ⮌ Even security experts cannot analyze every line of code

## Penetration testing / security testing

## Static code analysis

Monday 10 June 13 week

# How to find vulnerabilities?

## Code inspection

- ⮑ Manual verification is not feasible

- ⮑ Not all apps can afford security experts

- ⮑ Even security experts cannot analyze every line of code

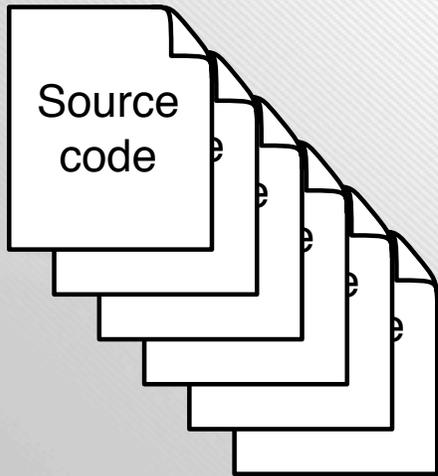## Penetration testing / security testing

## Static code analysis

## Magic

- ⮑ Vulnerability prediction models

Monday 10 June 13 week

# Vulnerability prediction model

Monday 10 June 13 week

# Vulnerability prediction model

Source
code

Monday 10 June 13 week

# Vulnerability prediction model

Monday 10 June 13 week

# Vulnerability prediction model

Source
code

Machine learning

Source
code

Monday 10 June 13 week

# Outline

**Existing tools and techniques**

➲ Vulnerability prediction models

Our approach

Results

Conclusions and future research

Monday 10 June 13 week

# Vulnerability prediction models

Monday 10 June 13 week

# Vulnerability prediction models

Monday 10 June 13 week

# Vulnerability prediction models

Monday 10 June 13 week

# Vulnerability prediction models

## Start from a hunch = feature

⮕ e.g., larger components are more likely to be vulnerable

Monday 10 June 13 week

# Vulnerability prediction models

## Start from a hunch = feature

 ⮎ e.g., larger components are more likely to be vulnerable

## Fetch the features from the components

 ⮎ e.g., calculate the size for each component

Monday 10 June 13 week

# Vulnerability prediction models

## Start from a hunch = feature

⮑ e.g., larger components are more likely to be vulnerable

## Fetch the features from the components

⮑ e.g., calculate the size for each component

## Determine the vulnerabilities

⮑ e.g., National Vulnerability Database, MFSA

## Investigate the correlation

⮑ Use machine learning techniques

Monday 10 June 13 week

# Vulnerability prediction models

Monday 10 June 13 week

# Vulnerability prediction models

## Typical "hunches"

➲ Use size and complexity metrics

➲ Leverage developer activity metrics

➲ Leverage code churn metrics

➲ Leverage design churn metrics

➲ Number of import statements

Monday 10 June 13 week

# Vulnerability prediction models

## Typical "hunches"

➲ Use size and complexity metrics

➲ Leverage developer activity metrics

➲ Leverage code churn metrics

➲ Leverage design churn metrics

➲ Number of import statements

## Inspired on the defect prediction work

➲ Vulnerabilities are actually defects, but much more scarce ("needle in a haystack")

Monday 10 June 13 week

# Outline

## Existing tools and techniques

- Static code analysis
- Vulnerability prediction using metrics

## Our approach

## Results

## Conclusions and future research

Monday 10 June 13 week

# Our approach

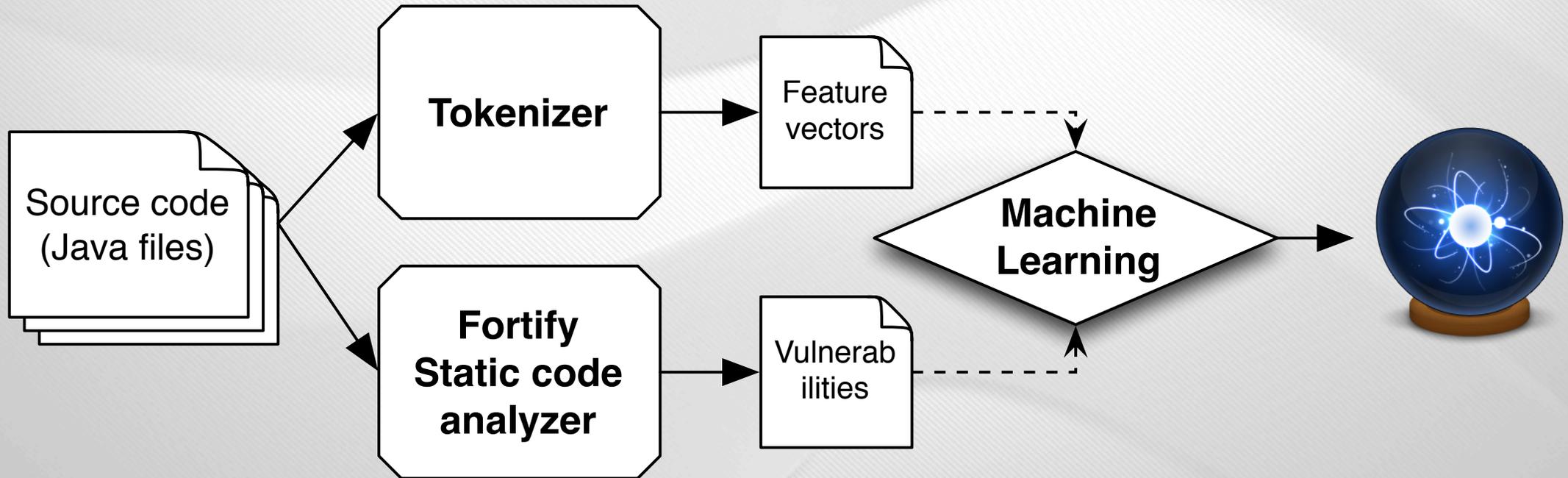Use the source code itself in a tokenized form

Use the token frequency as features

- ➲ Simplicity
- ➲ No explicit assumptions regarding the code characteristics
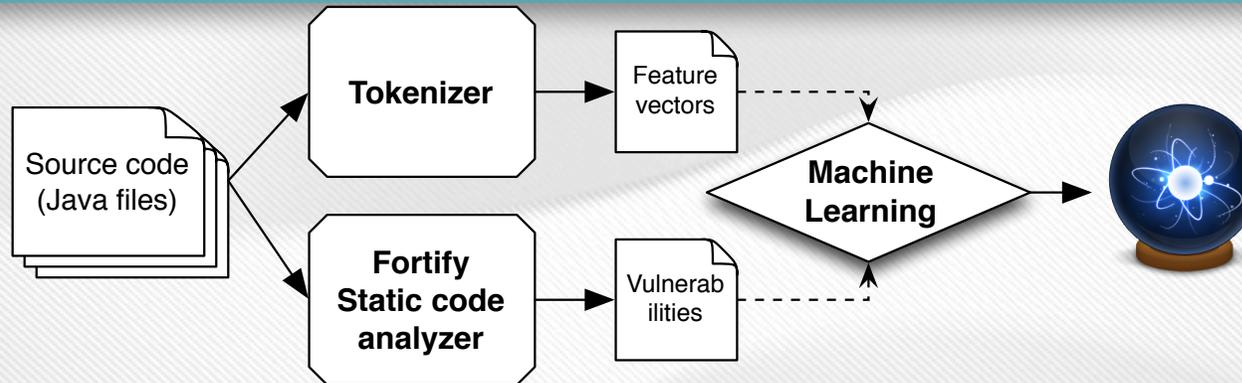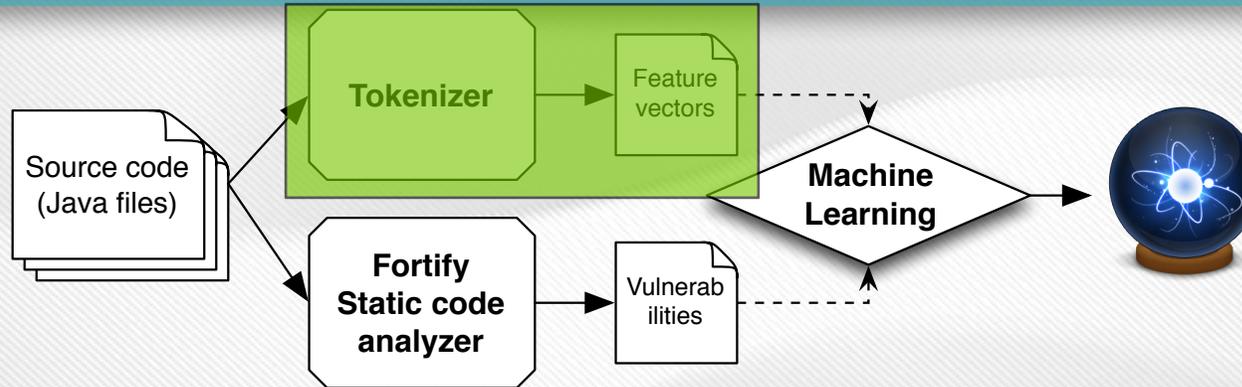
#text analysis

#SPAM filtering

#machine learning

Monday 10 June 13 week
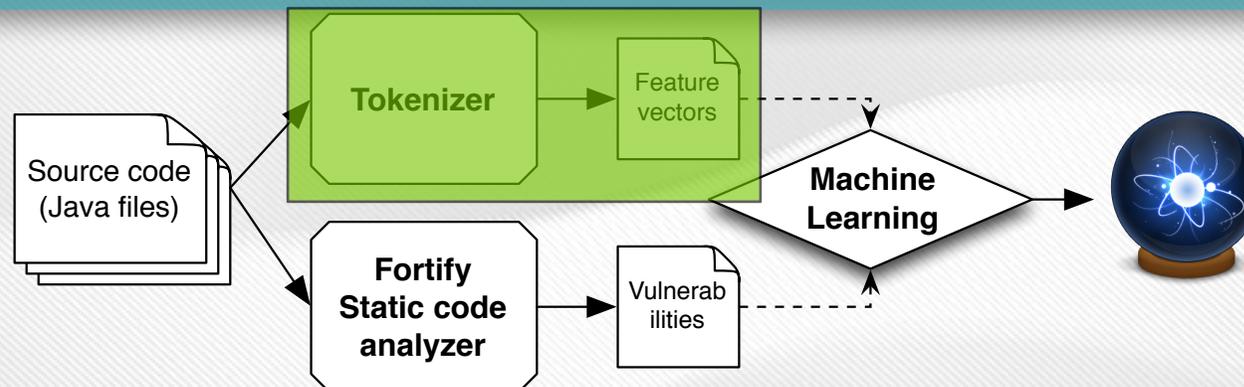
# Our approach

Monday 10 June 13 week

# Tokenizer

Monday 10 June 13 week

# Tokenizer

Monday 10 June 13 week

# Tokenizer

## Transform each source code token into a feature vector

- each token ("monogram") is a feature

- tokenize by delimiters, mathematical and logical operations

  . , ! < > [ ] = + - ^ * / etc.

- each feature has a count assigned to it

Monday 10 June 13 week

# Feature vector

Monday 10 June 13 week

# Feature vector

```java
package com.fsck.k9;
import android.text.util.Rfc822Tokenizer;
import android.widget.AutoCompleteTextView.Validator;
public class EmailAddressValidator implements Validator
{

    public CharSequence fixText(CharSequence invalidText)
    {
        return "";
    }
    public boolean isValid(CharSequence text)
    {
        return Rfc822Tokenizer.tokenize(text).length > 0;
    }
}
```

Monday 10 June 13 week

# Feature vector

```
package com.fsck.k9;
import android.text.util.Rfc822Tokenizer;
import android.widget.AutoCompleteTextView.Validator;
public class EmailAddressValidator implements Validator
{
    public CharSequence fixText(CharSequence invalidText)
    {
        return "";
    }
    public boolean isValid(CharSequence text)
    {
        return Rfc822Tokenizer.tokenize(text).length > 0;
    }
}
```

package: 1

Monday 10 June 13 week

# Feature vector

```
package com.fsck.k9;
import android.text.util.Rfc822Tokenizer;
import android.widget.AutoCompleteTextView.Validator;
public class EmailAddressValidator implements Validator
{
    public CharSequence fixText(CharSequence invalidText)
    {
        return "";
    }
    public boolean isValid(CharSequence text)
    {
        return Rfc822Tokenizer.tokenize(text).length > 0;
    }
}
```
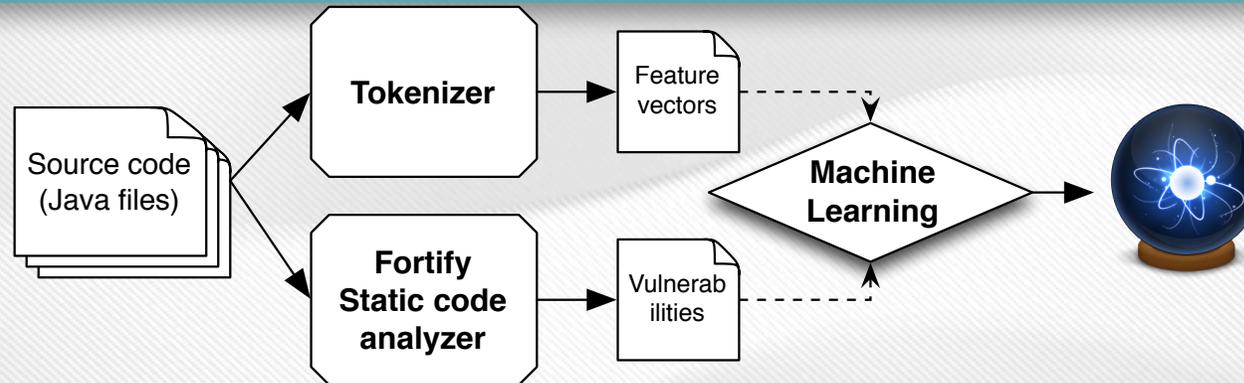
package: 1, com: 1

Monday 10 June 13 week
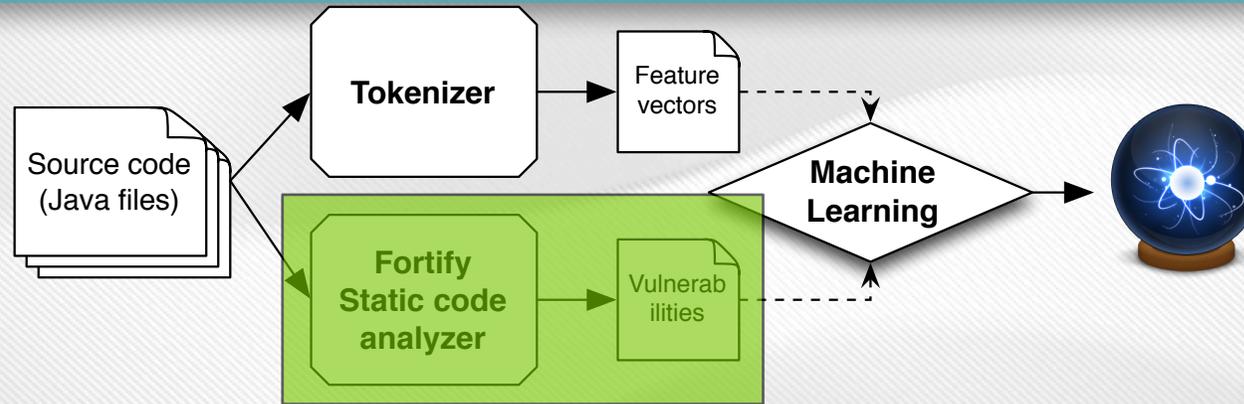
# Feature vector

```
package com.fsck.k9;
import android.text.util.Rfc822Tokenizer;
import android.widget.AutoCompleteTextView.Validator;
public class EmailAddressValidator implements Validator
{
    public CharSequence fixText(CharSequence invalidText)
    {
        return "";
    }
    public boolean isValid(CharSequence text)
    {
        return Rfc822Tokenizer.tokenize(text).length > 0;
    }
}
```

package: 1, com: 1, fsck: 1, k9: 1, import: 2, android: 2, text: 2, util: 1, Rfc822Tokenizer: 2, widget: 1, AutoCompleteTextView:1, Validator: 2, public: 3, class: 1, EmailAddressValidator: 1, implements: 1, CharSequence: 2, fixText: 1, invalidText: 1, return: 2, tokenize: 1, length: 1
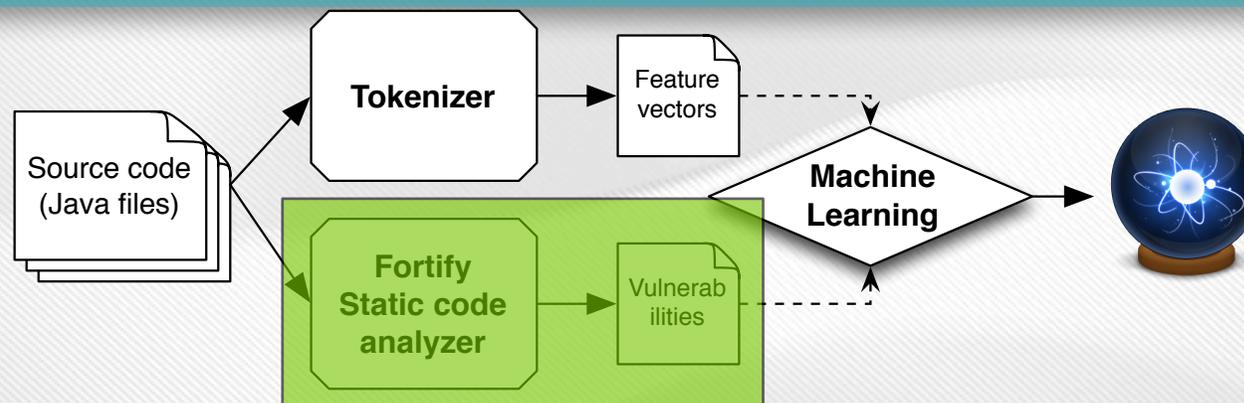
Monday 10 June 13 week

# Vulnerability assignment

Monday 10 June 13 week

# Vulnerability assignment

Monday 10 June 13 week

# Vulnerability assignment

Source code
(Java files)

**Tokenizer**

Feature
vectors

**Fortify
Static code
analyzer**

Vulnerab
ilities

**Machine
Learning**

## Assign vulnerability to each Java file

➲ use Fortify (static code analyzer) for this task

➲ each file is either *vulnerable* or *clean*
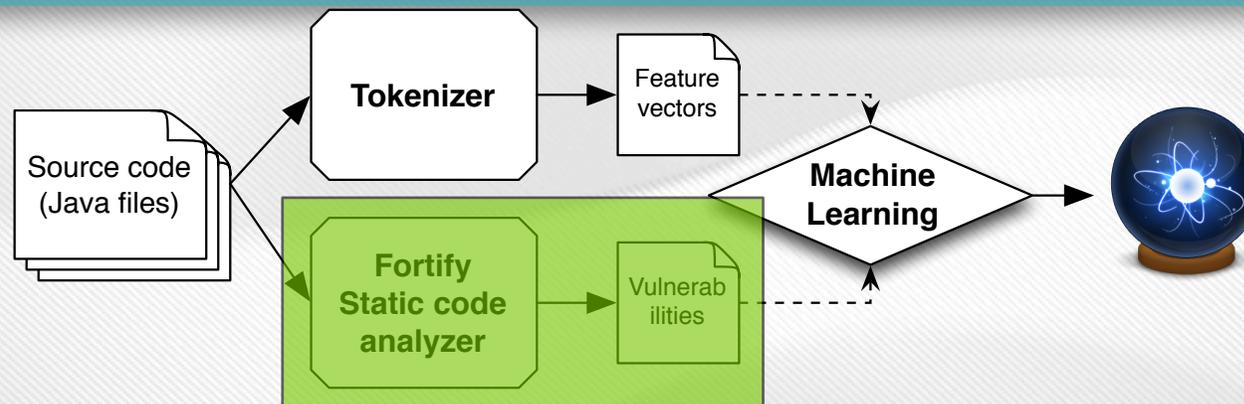
Monday 10 June 13 week
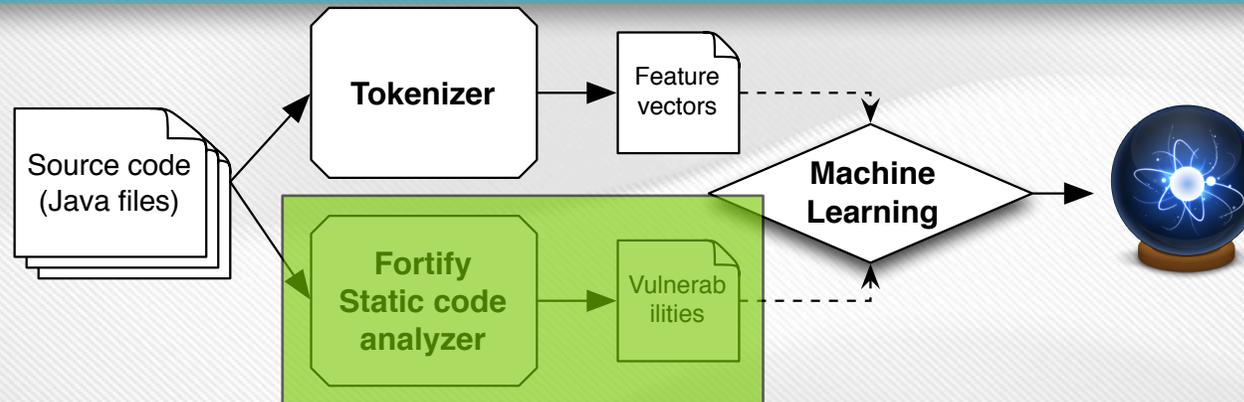
# Vulnerability assignment



## Assign vulnerability to each Java file

- ➲ use Fortify (static code analyzer) for this task

- ➲ each file is either *vulnerable* or *clean*

package: 1, com: 1, fsck: 1, k9: 1, import: 2, android: 2, text: 2, util: 1, Rfc822Tokenizer: 2, widget: 1, AutoCompleteTextView:1, Validator: 2, public: 3, class: 1, EmailAddressValidator: 1, implements: 1, CharSequence: 2, fixText: 1, invalidText: 1, return: 2, tokenize: 1, length: 1
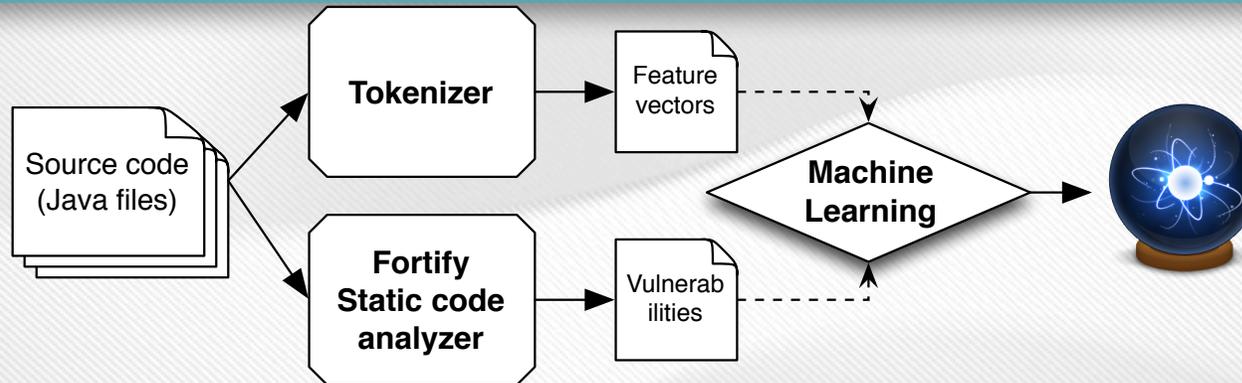
Monday 10 June 13 week

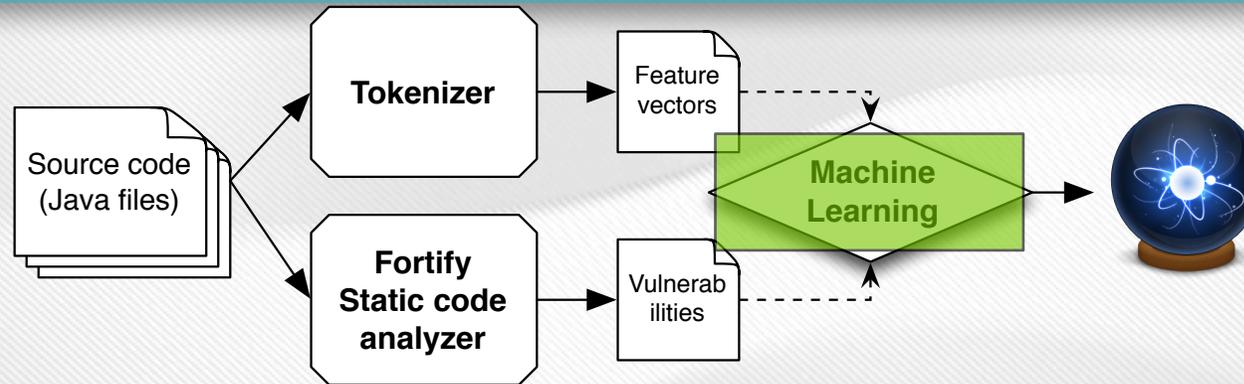# Vulnerability assignment

## Assign vulnerability to each Java file

⮌ use Fortify (static code analyzer) for this task

⮌ each file is either *vulnerable* or *clean*

package: 1, com: 1, fsck: 1, k9: 1, import: 2, android: 2, text: 2, util: 1, Rfc822Tokenizer: 2, widget: 1, AutoCompleteTextView:1, Validator: 2, public: 3, class: 1, EmailAddressValidator: 1, implements: 1, CharSequence: 2, fixText: 1, invalidText: 1, return: 2, tokenize: 1, length: 1, vulnerability: 0
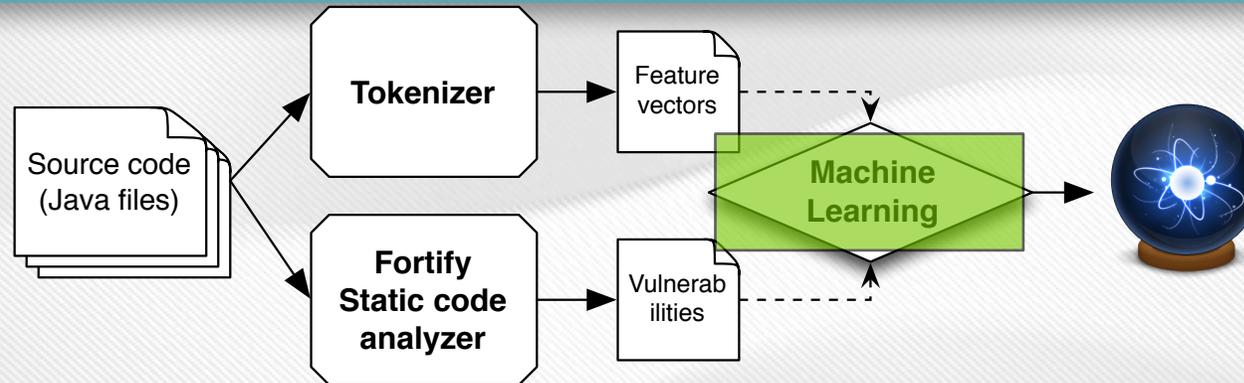
Monday 10 June 13 week

# Machine learning

Monday 10 June 13 week

# Machine learning

Source code (Java files) → Tokenizer → Feature vectors

Source code (Java files) → Fortify Static code analyzer → Vulnerabilities

Feature vectors, Vulnerabilities → Machine Learning → (prediction)

Monday 10 June 13 week

# Machine learning



Leverage machine learning techniques to build a prediction model

- ⊃ Training set -> the data used to train the model
- ⊃ Testing set -> the data used to validate the model

Various techniques available (SVM, Naive Bayes, Random Forest, CART, kNN)

Monday 10 June 13 week

# Experiment 1

Monday 10 June 13 week

# Experiment 1

Can we predict future versions of an app based on its first version?

Monday 10 June 13 week

# Experiment 1

Can we predict future versions of an app based on its first version?

- ➲ **Training set** - the first version (v0) of an app
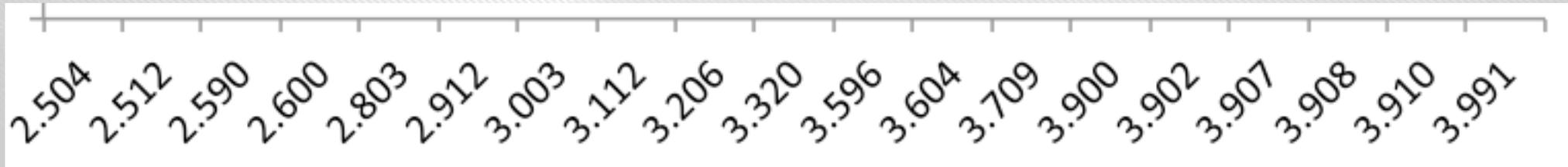
Monday 10 June 13 week

# Experiment 1

Can we predict future versions of an app based on its first version?

- ⮑ **Training set** - the first version (v0) of an app
- ⮑ **Testing set** - all subsequent versions of that app

Monday 10 June 13 week

Can we predict future versions of an app based on its first version?

- **Training set** - the first version (v0) of an app
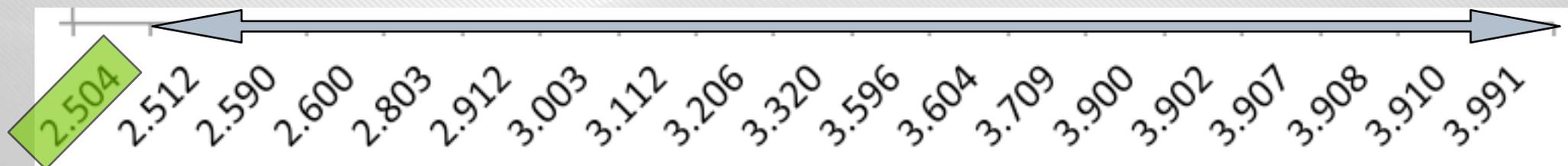- **Testing set** - all subsequent versions of that app

2.504   2.512   2.590   2.600   2.803   2.912   3.003   3.112   3.206   3.320   3.596   3.604   3.709   3.900   3.902   3.907   3.908   3.910   3.991

Monday 10 June 13 week

# Experiment 1

Can we predict future versions of an app based on its first version?

- ➲ **Training set** - the first version (v0) of an app
- ➲ **Testing set** - all subsequent versions of that app

Monday 10 June 13 week

# Experiment 1

## Can we predict future versions of an app based on its first version?

- ➲ **Training set** - the first version (v0) of an app

- ➲ **Testing set** - all subsequent versions of that app

2.504  2.512  2.590  2.600  2.803  2.912  3.003  3.112  3.206  3.320  3.596  3.604  3.709  3.900  3.902  3.907  3.908  3.910  3.991

- ➲ Repeat for all apps

Monday 10 June 13 week

# Experiment 2

Monday 10 June 13 week

# Experiment 2

Can we build a generalized predictor that works on all apps?

- ➲ **Training set** - the first version (v0) of an app
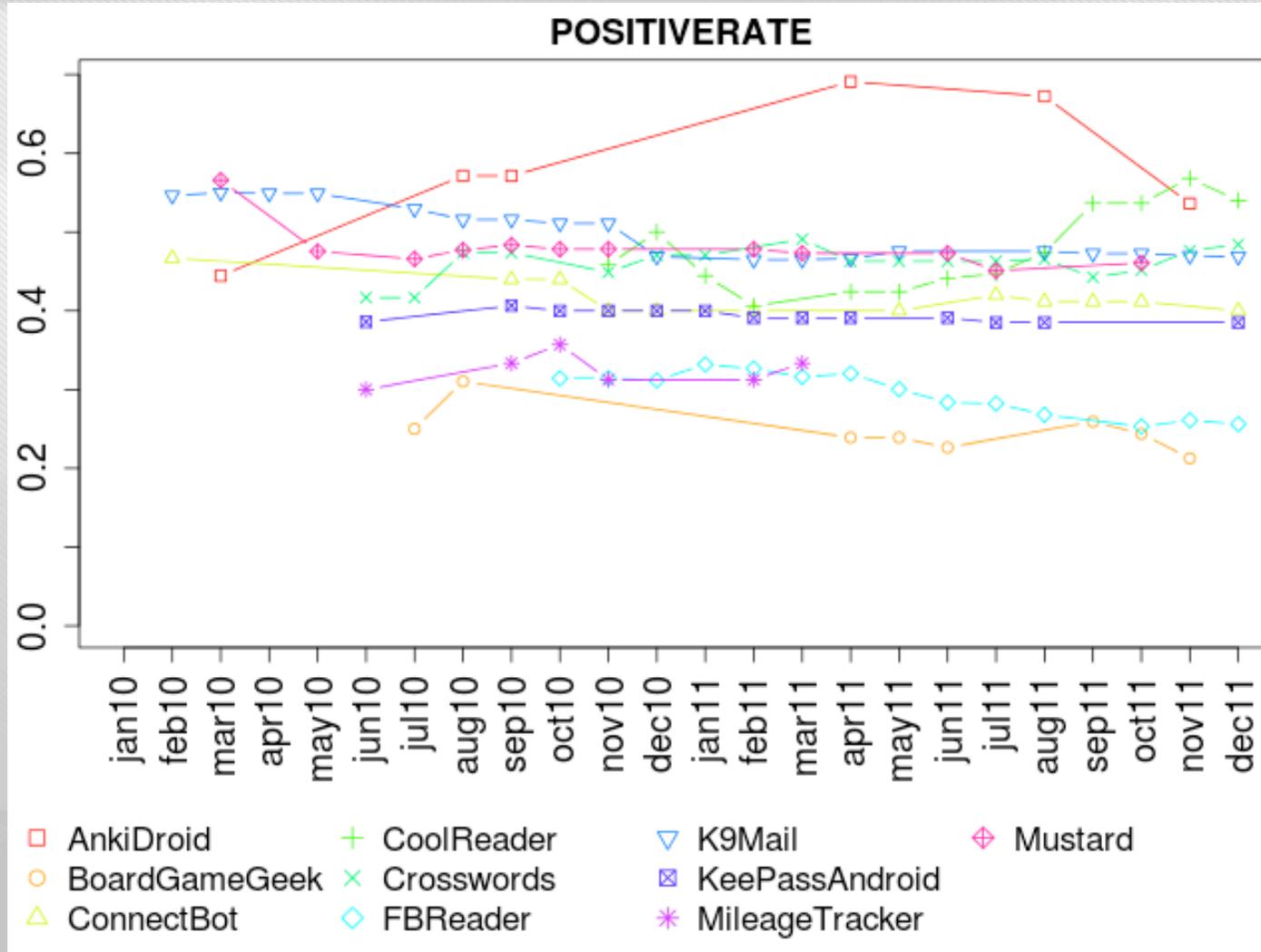- ➲ **Testing set** - first versions of all other apps

Monday 10 June 13 week

# Applications (data from early 2012)

| Application | Category | Downloads | Versions |
|---|---|---|---|
| AnkiDroid | education | 100k - 500k | 8 |
| BoardGameGeek | books | 10k - 50k | 8 |
| Connectbot | communication | 1M - 5M | 12 |
| CoolReader | books | 1M - 5M | 13 |
| Crosswords | brain & puzzle | 5k - 10k | 17 |
| FBReader | books | 1M - 5M | 14 |
| K9 Mail | communication | 1M - 5M | 19 |
| KeePassAndroid | tools | 100k - 500k | 13 |
| MileageTracker | finance | 100k - 500k | 6 |
| Mustard | social | 10k - 50k | 12 |

- F-droid repository: 01/01/2010->31/12/2011

- Selection criteria: open-source, size, number of versions

Monday 10 June 13 week

# Applications: descriptive statistics

Monday 10 June 13 week

Monday 10 June 13 week

# Outline

Existing tools and techniques

- ➲ Static code analysis
- ➲ Vulnerability prediction using metrics

Our approach

**Results**

Conclusions and future research

Monday 10 June 13 week

# Performance indicators

Monday 10 June 13 week

# Performance indicators

Accuracy: percentage of correctly classified files

➲ imagine 90% of the files are clean

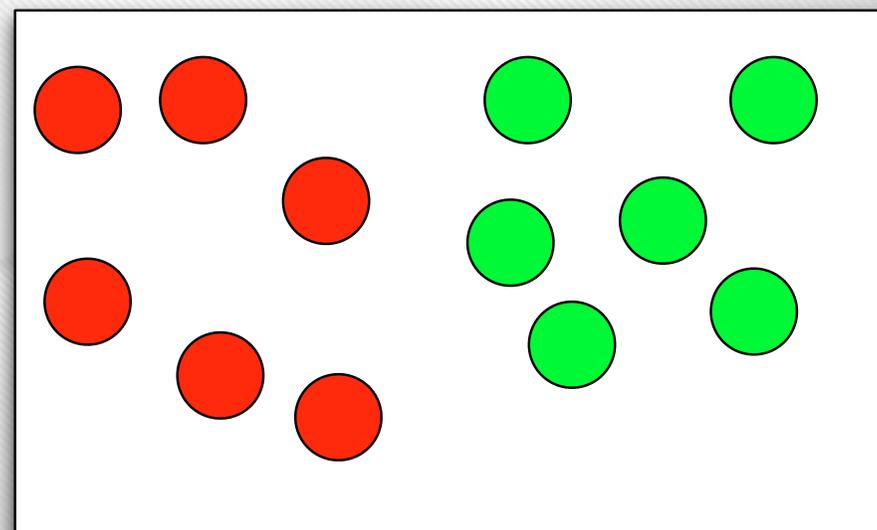➲ saying all files are clean will achieve 90% accuracy

Monday 10 June 13 week

## Accuracy: percentage of correctly classified files

➲ imagine 90% of the files are clean

➲ saying all files are clean will achieve 90% accuracy

Monday 10 June 13 week

# Performance indicators

Accuracy: percentage of correctly classified files

- ⟲ imagine 90% of the files are clean

- ⟲ saying all files are clean will achieve 90% accuracy
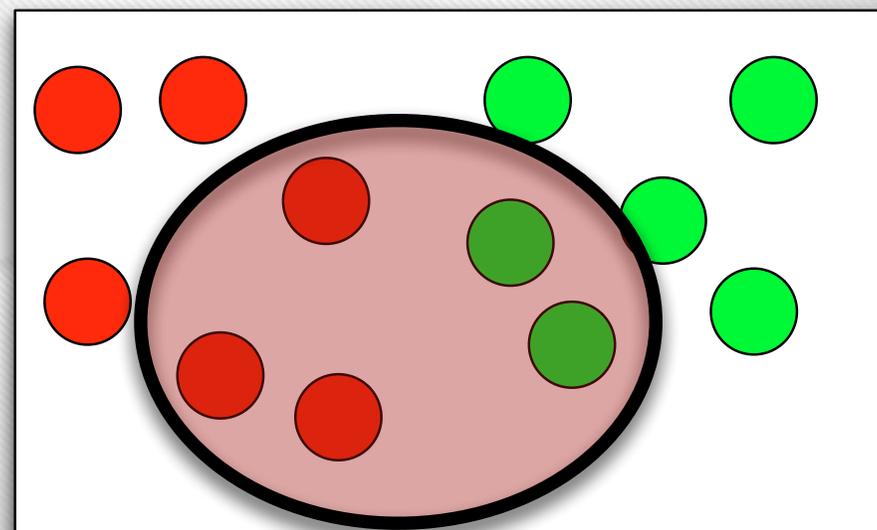
Monday 10 June 13 week

# Performance indicators

Accuracy: percentage of correctly classified files

- ➲ imagine 90% of the files are clean
- ➲ saying all files are clean will achieve 90% accuracy

Prediction vs. reality

- ➲ True positive (TP)
- ➲ True negative (TN)
- ➲ False positive (FP)
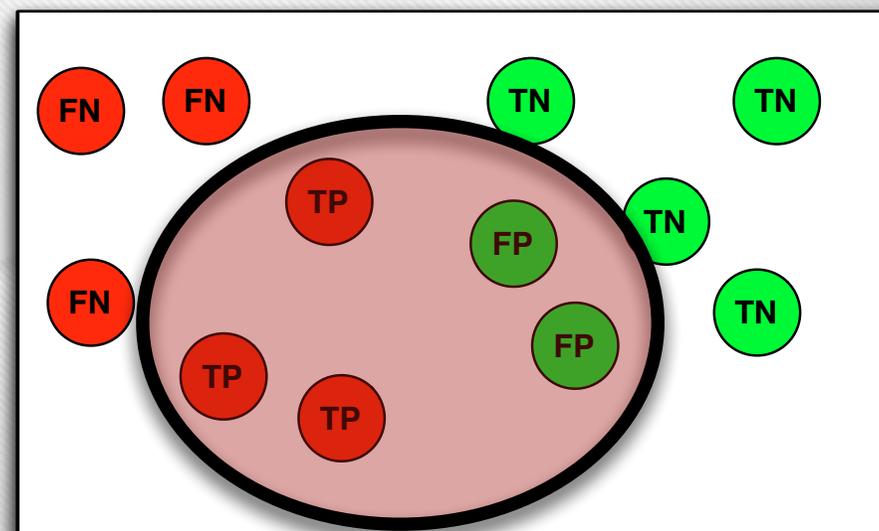- ➲ False negative (FN)

Monday 10 June 13 week

# Performance indicators

## Accuracy: percentage of correctly classified files

➲ imagine 90% of the files are clean

➲ saying all files are clean will achieve 90% accuracy

## Prediction vs. reality

➲ True positive (TP)

➲ True negative (TN)

➲ False positive (FP)

➲ False negative (FN)

Monday 10 June 13 week
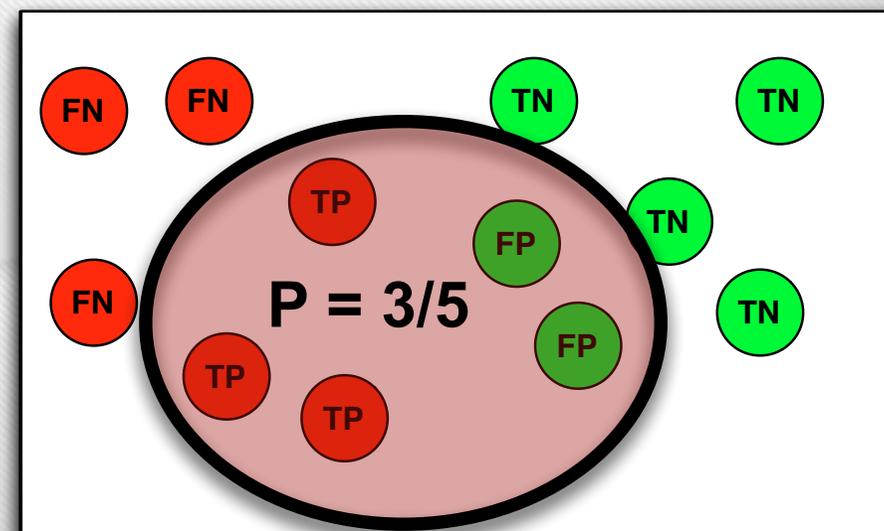
# Performance indicators

**Accuracy: percentage of correctly classified files**

- ➲ imagine 90% of the files are clean
- ➲ saying all files are clean will achieve 90% accuracy

**Prediction vs. reality**

- ➲ True positive (TP)
- ➲ True negative (TN)
- ➲ False positive (FP)
- ➲ False negative (FN)

**Precision: P = TP/(TP+FP)**

Monday 10 June 13 week

# Performance indicators

Accuracy: percentage of correctly classified files

- ➲ imagine 90% of the files are clean
- ➲ saying all files are clean will achieve 90% accuracy

Prediction vs. reality

- ➲ True positive (TP)
- ➲ True negative (TN)
- ➲ False positive (FP)
- ➲ False negative (FN)

Precision: P = TP/(TP+FP)
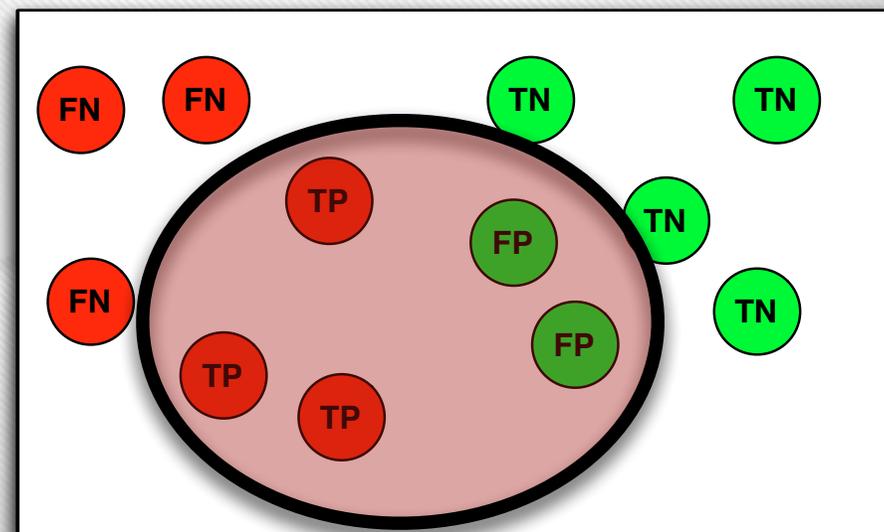
Monday 10 June 13 week

# Performance indicators

Accuracy: percentage of correctly classified files

➲ imagine 90% of the files are clean

➲ saying all files are clean will achieve 90% accuracy

Prediction vs. reality

➲ True positive (TP)

➲ True negative (TN)

➲ False positive (FP)

➲ False negative (FN)

Precision: P = TP/(TP+FP)

Recall: R = TP/(TP+FN)

Monday 10 June 13 week

# Performance indicators
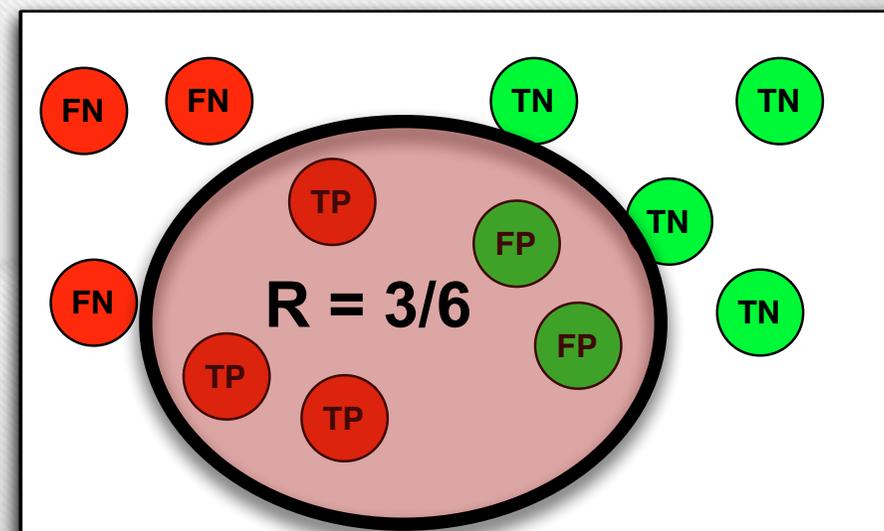
## Accuracy: percentage of correctly classified files
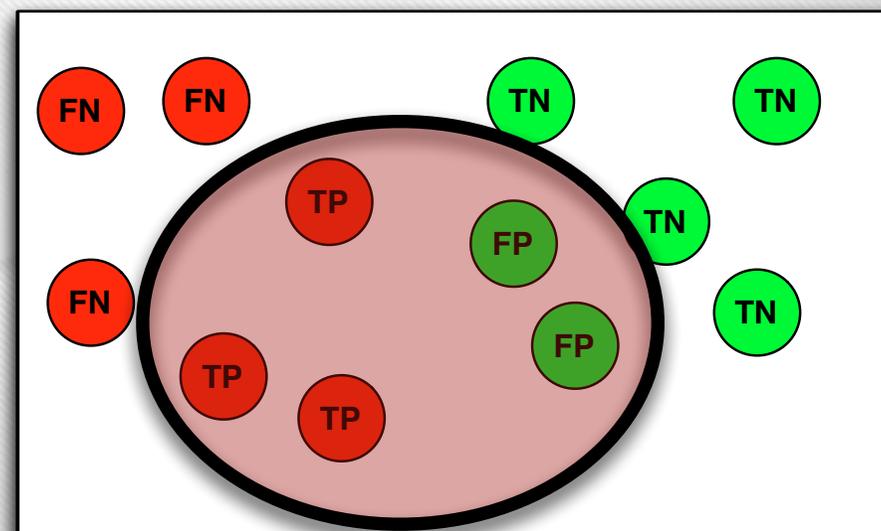
- imagine 90% of the files are clean
- saying all files are clean will achieve 90% accuracy

## Prediction vs. reality

- True positive (TP)
- True negative (TN)
- False positive (FP)
- False negative (FN)

## Precision: P = TP/(TP+FP)

## Recall: R = TP/(TP+FN)

Monday 10 June 13 week

# Performance indicators: K9

Monday 10 June 13 week

# Performance indicators: K9



K9Mail (Random Forest)

Monday 10 June 13 week

# Experiment 1: future predictions

Monday 10 June 13 week

# Experiment 1: future predictions

## When do we need to build a new model?

⮑ Retrain when performance indicators drop with 10%

Monday 10 June 13 week

# Experiment 1: future predictions

## When do we need to build a new model?

➲ Retrain when performance indicators drop with 10%

| Application | Retrain (months) |
|-------------|------------------|
| AnkiDroid | -- |
| BoardGameGeek | 9 |
| ConnectBot | -- |
| CoolReader | 10 |
| Crosswords | 2 |
| FBReader | -- |
| K9Mail | 12 |
| KeePassAndroid | -- |
| MileageTracker | 1 |
| Mustard | -- |

-- no retraining is required

Monday 10 June 13 week

# Results: most influential features

Monday 10 June 13 week

# Results: most influential features

Most influential features

Monday 10 June 13 week

# Results: most influential features

## Most influential features

➲ *e, Exception, try, catch (error handling)*

Monday 10 June 13 week

# Results: most influential features

## Most influential features

- *e, Exception, try, catch (error handling)*
- *if (branching)*

Monday 10 June 13 week

# Results: most influential features

## Most influential features

- ⮌ *e, Exception, try, catch (error handling)*
- ⮌ *if (branching)*
- ⮌ *null (pointer algebra)*

Monday 10 June 13 week

# Results: most influential features

## Most influential features

- ➲ *e, Exception, try, catch (error handling)*
- ➲ *if (branching)*
- ➲ *null (pointer algebra)*
- ➲ java, org (import statements)

Monday 10 June 13 week

# Results: most influential features

## Most influential features

- ➲ *e, Exception, try, catch (error handling)*
- ➲ *if (branching)*
- ➲ *null (pointer algebra)*
- ➲ java, org (import statements)
- ➲ new, Log (others)

Monday 10 June 13 week

# Results: most influential features

## Most influential features

- *e, Exception, try, catch (error handling)*
- *if (branching)*
- *null (pointer algebra)*
- java, org (import statements)
- new, Log (others)

## Produced by InfoGain

Monday 10 June 13 week

# Validity threats

## Use of Fortify tool for vulnerability extraction

- ➲ Some research results have shown that there are strong correlations between static analysis metrics and the quality of reported vulnerabilities

- ➲ **Manual validation seems to confirm our findings (work in progress)!**

- ➲ **We are currently validating the same technique on Mozilla Firefox and the results are slightly better than the existing work**

Monday 10 June 13 week

# Bring your own data

We are looking to validate our technique further

If you have data you are willing to share with us, we would be glad to collaborate

Monday 10 June 13 week