

# DETECTION OF POTENTIALLY DANGEROUS ACTIVITIES FROM LOGS OF MOBILE DEVICES USING MACHINE LEARNING TECHNIQUES

DEPARTMENT OF PARALLEL AND DISTRIBUTED INFORMATION PROCESSING  
INSTITUTE OF INFORMATICS, SLOVAK ACADEMY OF SCIENCES



INSTITUTE OF INFORMATICS  
SLOVAK ACADEMY OF SCIENCES

Tempest

IT makes sense



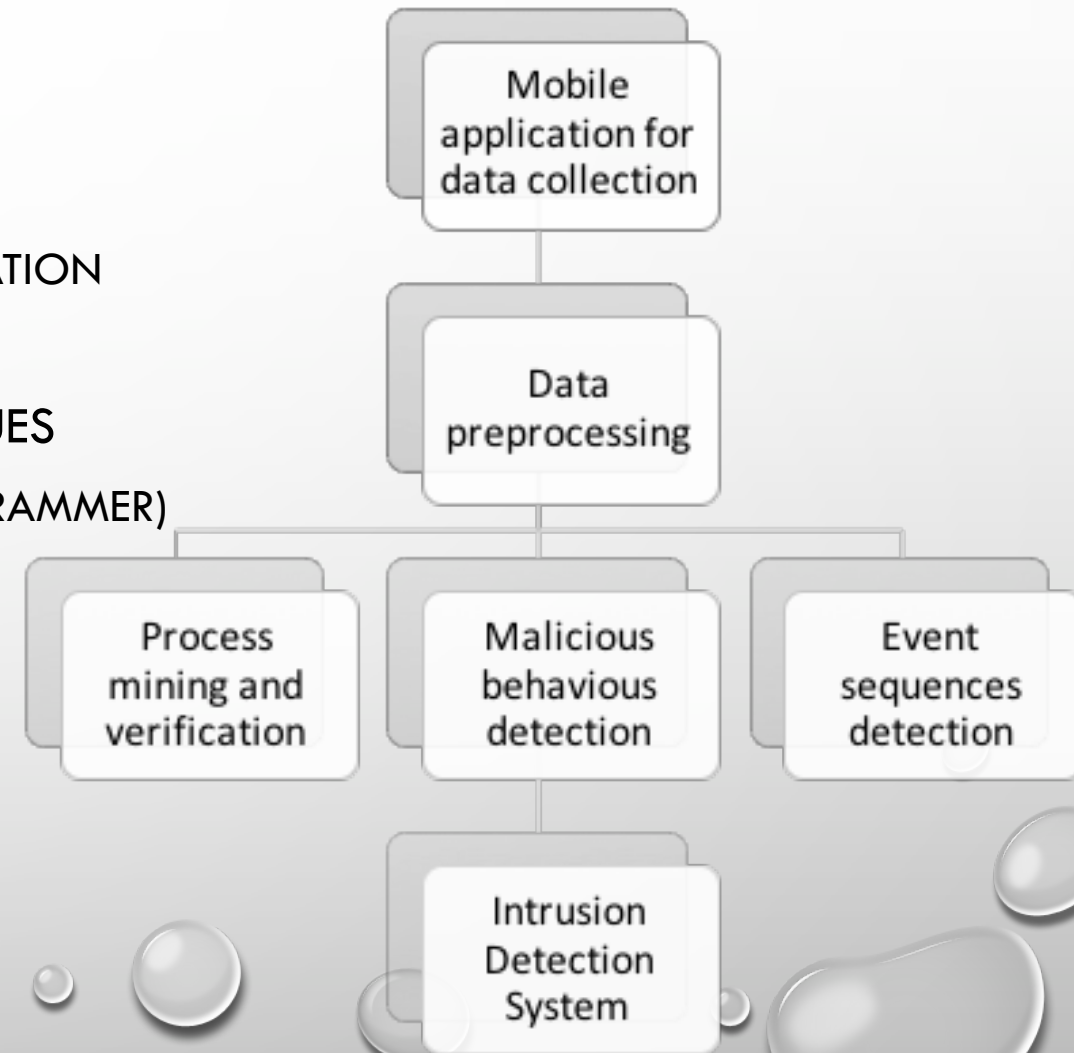
IBM Slovakia

# AIMS AND REQUIREMENTS (IBM)

- **CYBER-SECURITY** RESEARCH FOR MOBILE DEVICES, **BIG DATA TECHNOLOGY ORIENTED**
- RESPONSIBILITIES
  - IBM, TEMPEST A.S.: MOBILE DEVICES, DATA COLLECTIONS, QUEST DEFINITIONS
  - **IISAS TEAM:** DATA MINING USING MACHINE LEARNING TECHNIQUES – **MALWARE DETECTION**
- **SIX MONTHS (6.2015-11.2015)**
  - ALONG WITH DATA COLLECTIONS, INCLUDED COMPLETE DOCUMENTATIONS
  - **PILOT RESEARCH AND DEVELOPMENT WITH RESULT EVALUATIONS**
  - PROOF OF THE CONCEPT, DEFINITIONS, ALGORITHMS, THEOREMS, ANALYZING, EXPERIMENTS ...

# RESEARCH INTERESTS - BUSINESS UNDERSTANDING

- **IISAS TEAM: DOC. ING. LADISLAV HLUCHÝ, PHD.**
  - DEPARTMENT OF PARALLEL AND DISTRIBUTED INFORMATION PROCESSING
- DATA MINING USING MACHINE LEARNING TECHNIQUES
  - **MALICIOUS BEHAVIOR DETECTION** (G. NGUYEN, P. KRAMMER)
  - EVENT SEQUENCES DETECTION (Š. DLUGOLINSKÝ)
  - PROCESS MINING AND VERIFICATION (O. HABALA)
  - INTRUSION DETECTION SYSTEM (V. TRAN)
  - INFRASTRUCTURE (M. ŠELENG)

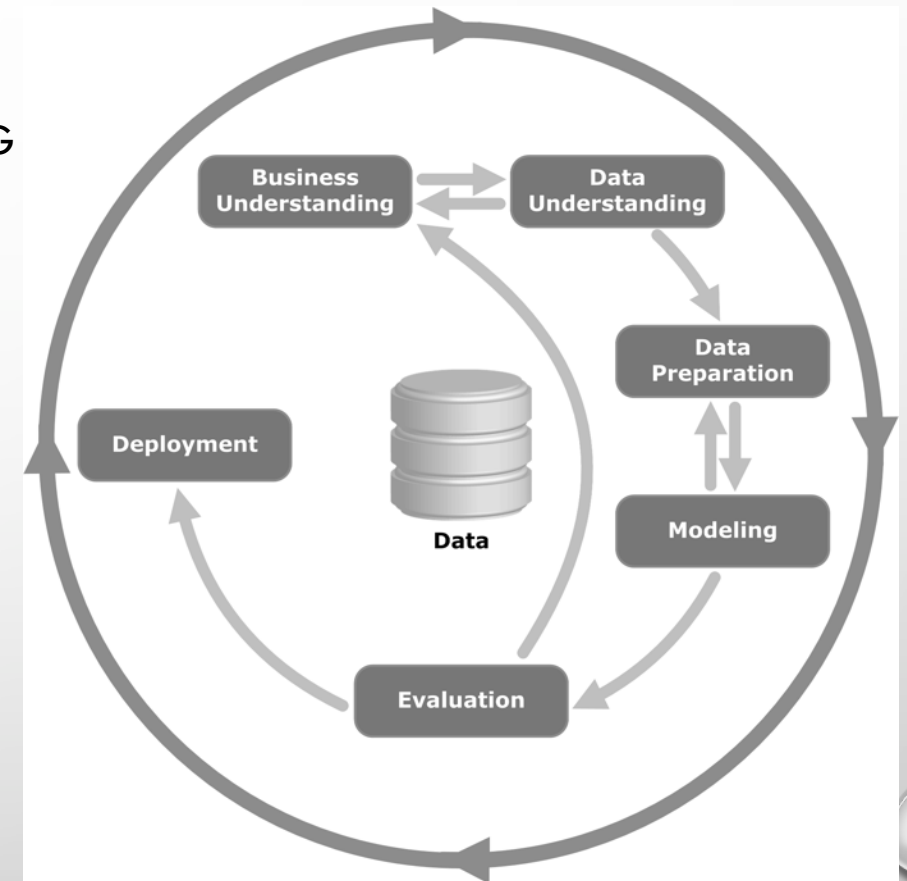


# RAW LOGS

20150908-23:35:55:523\$location\_unavailable\$-9.107543;-0.277727;2.547431\$http://www.hlavnespravdy.sk/\$\$Hlavné správy - Konzervat  
ivny dennik\$20150908-18:12:02:176\$\$(85.248.228.27)\$S80\$Shttp\$S4b09f848cd21186f  
20150908-23:35:55:526\$location\_unavailable\$-9.107543;-0.277727;2.547431\$http://www.pravda.sk/\$\$Pravda.sk\$20150908-21:50:08:17  
1\$\$(217.67.31.48)\$S80\$Shttp\$S4b09f848cd21186f  
20150908-23:35:55:528\$location\_unavailable\$-9.107543;-0.277727;2.547431\$http://spravy.pravda.sk/ekonomika/clanok/366526-regioj  
et-ma-opat-problem-s-licenciou-na-autobusovu-linku/\$\$RegioJet má opäť problém s licenciou na autobusovú linku - Ekonomika - Správ  
y - Pravda.sk\$20150903-17:37:57:394\$\$(217.67.31.48)\$S80\$Shttp\$S4b09f848cd21186f  
20150908-23:35:55:531\$location\_unavailable\$-9.107543;-0.277727;2.547431\$http://spravy.pravda.sk/ekonomika/clanok/366469-sloven  
ske-rajciny-po-cely-rok-domaci-zeleninari-sa-spajaju/\$\$Slovenské rajčiny po\$20150903-17:39:22:851\$\$(217.67.31.48)\$S80\$Shttp\$S4b0  
9f848cd21186f  
20150908-23:35:55:534\$location\_unavailable\$-9.107543;-0.277727;2.547431\$http://spravy.pravda.sk/domace/clanok/366523-kiska-tel  
efonoval-s-ukrajinskym-prezidentom-porosenkom/\$\$Kiska telefonoval s ukrajinským prezidentom Porošenkom - Domáce - Správy - Pravda  
.sk\$20150903-17:43:49:506\$\$(217.67.31.48)\$S80\$Shttp\$S4b09f848cd21186f  
20150908-23:35:55:537\$location\_unavailable\$-9.107543;-0.277727;2.547431\$http://debata.pravda.sk/debata/366523-kiska-telefonova  
l-s-ukrajinskym-prezidentom-porosenkom/\$\$Pravda.sk - Debata - Kiska telefonoval s ukrajinským prezidentom Porošenkom\$20150903-17  
:43:48:572\$\$(217.67.31.59)\$S80\$Shttp\$S4b09f848cd21186f  
20150908-23:35:55:540\$location\_unavailable\$-9.107543;-0.277727;2.547431\$http://debata.pravda.sk/debata/366523-kiska-telefonova  
l-s-ukrajinskym-prezidentom-porosenkom/?view\_mode=vlakna&ordering=od\_najnovsieho&strana=2\$Pravda.sk - Debata - Kiska telefonoval  
s ukrajinským prezidentom Porošenkom\$20150903-17:43:38:565\$\$(217.67.31.59)\$S80\$Shttp\$S4b09f848cd21186f  
20150908-23:35:55:543\$location\_unavailable\$-9.107543;-0.277727;2.547431\$http://spravy.pravda.sk/svet/clanok/366525-orban-sokov  
al-nemeckych-politikov-imigracna-vlna-je-vas-problem-vyhlasil/\$\$Orbán šokoval nemeckých politikov. Migračná vlna je váš problém,  
vyhlásil - Svet - Správy - Pravda.sk\$20150903-17:44:04:015\$\$(217.67.31.48)\$S80\$Shttp\$S4b09f848cd21186f  
20150908-23:35:55:546\$location\_unavailable\$-9.107543;-0.277727;2.547431\$http://debata.pravda.sk/debata/366525-orban-sokoval-ne  
meckych-politikov-imigracna-vlna-je-vas-problem-vyhlasil/\$\$Pravda.sk - Debata - Orbán šokoval nemeckých politikov. Migračná vlna  
je váš problém, vyhlásil\$20150903-17:45:46:794\$\$(217.67.31.59)\$S80\$Shttp\$S4b09f848cd21186f  
20150908-23:35:55:552\$location\_unavailable\$-9.107543;-0.277727;2.547431\$http://www.aktuality.sk/\$\$Aktuálne spravodajstvo na Sl  
ovensku a vo svete | Aktuality.sk\$20150903-17:52:43:075\$\$(91.235.52.35)\$S80\$Shttp\$S4b09f848cd21186f  
20150908-23:35:55:552\$location\_unavailable\$-9.107543;-0.277727;2.547431\$http://www.aktuality.sk/clanok/303370/v-konflikte-na-j  
uhovychode-ukrajiny-zahynulo-uz-vyse-6400-civilistov/\$\$V konflikte na juhovýchode Ukrajiny zahynulo už vyše 6400 civilistov\$2015  
0903-17:50:10:165\$\$(91.235.52.35)\$S80\$Shttp\$S4b09f848cd21186f  
20150908-23:35:55:555\$location\_unavailable\$-9.107543;-0.277727;2.547431\$http://www.aktuality.sk/clanok/303368/narodne-gmo-zaka  
zy-europoslanci-z-polnohospodarskeho-vyboru-odmietaju/\$\$Národné GMO zákazy europoslanci z poľnohospodárskeho výboru odmietajú\$20  
150903-17:50:47:203\$\$(91.235.52.35)\$S80\$Shttp\$S4b09f848cd21186f  
20150908-23:35:55:557\$location\_unavailable\$-9.107543;-0.277727;2.547431\$http://www.cas.sk/\$\$ČAS.sk - Najčítanejší denník\$2015  
0908-18:24:55:167\$\$(91.235.52.131)\$S80\$Shttp\$S4b09f848cd21186f  
20150908-23:35:55:560\$location\_unavailable\$-9.107543;-0.277727;2.547431\$http://www.pluska.sk/\$\$Správy z domova a zahraničia, š  
oubiznis, krimi a šport | Pluska.sk\$20150903-17:59:19:058\$\$(85.248.116.198;85.248.116.194)\$S80\$Shttp\$S4b09f848cd21186f  
20150908-23:35:55:563\$location\_unavailable\$-9.107543;-0.277727;2.547431\$http://195.28.96.161/apps/\$\$\$\$20150908-13:02:08:766\$\$(91.235.52.35)\$S80\$Shttp\$S4b09f848cd21186f  
20150908-23:35:55:566\$location\_unavailable\$-9.107543;-0.277727;2.547431\$http://195.28.96.161/apps/\$\$\$\$20150908-18:15:29:647\$\$(91.235.52.35)\$S80\$Shttp\$S4b09f848cd21186f  
20150908-23:35:55:569\$location\_unavailable\$-9.107543;-0.277727;2.547431\$https://www.google.sk/search?q=facebook&oq=face&aqs=ch  
rome.0.0j69i57j0l2.1988j0j4&client=tablet-android-lenovo&sourceid=chrome-mobile&espv=1&ie=UTF-8\$facebook - Hľadať Google\$20150  
903-18:31:47:034\$\$(173.194.112.127;173.194.112.120;173.194.112.111;173.194.112.119;2a00:1450:4001:803::101f)\$S443\$https\$S4b09f84  
8cd21186f  
20150907-05:33:16:728\$S17.25484552238640;48.09291989707780\$S-0.306458;-0.296881;9.538500\$S754\$Scom.asus.launcher\$S10008\$S0\*\$S0\$S0  
\$S100\$S0\$S17MB\$S4MB\$S18MB\$S10MB\$S1016KB\$S10MB\$S7MB\$S19MB\$S26MB\$S35MB\$S25MB\$S356584067413743  
20150907-05:33:16:851\$S17.25484552238640;48.09291989707780\$S-0.306458;-0.296881;9.538500\$S850\$Scom.google.process.gapps\$S10019\$S0  
\$S893\$S200\$S2\$S13MB\$S5MB\$S13MB\$S4MB\$S1MB\$S4MB\$S5MB\$S13MB\$S7MB\$S23MB\$S19MB\$S356584067413743  
20150907-05:33:16:918\$S17.25484552238640;48.09291989707780\$S-0.306458;-0.296881;9.538500\$S1448\$Scom.asus.contacts\$S10005\$S0\*\$S0\$S0  
\$S130\$S0\$S7MB\$S4MB\$S7MB\$S3MB\$S1MB\$S3MB\$S18MB\$S10MB\$S16MB\$S24MB\$S356584067413743  
20150907-05:33:16:980\$S17.25484552238640;48.09291989707780\$S-0.306458;-0.296881;9.538500\$S838\$Scom.google.android.gms\$S10019\$S0\*\$S  
\$S1\$S0\$S400\$S0\$S10MB\$S4MB\$S10MB\$S6MB\$S1MB\$S6MB\$S9MB\$S13MB\$S11MB\$S25MB\$S19MB\$S356584067413743  
20150907-05:33:17:067\$S17.25484552238640;48.09291989707780\$S-0.306458;-0.296881;9.538500\$S6361\$Sandroid.process.acore\$S10005\$S0\*\$S  
\$S850\$S200\$S1\$S3MB\$S5MB\$S3MB\$S2MB\$S1MB\$S2MB\$S14MB\$S2MB\$S7MB\$S20MB\$S356584067413743  
20150907-05:33:17:115\$S17.25484552238640;48.09291989707780\$S-0.306458;-0.296881;9.538500\$S8624\$Scom.asus.sitd.whatsnext\$S10116\$S0  
\$S1\$S0\$S400\$S0\$S5MB\$S4MB\$S5MB\$S1MB\$S1MB\$S2MB\$S14MB\$S2MB\$S8MB\$S20MB\$S356584067413743  
20150907-05:33:17:126\$S17.25484552238640;48.09291989707780\$S-0.306458;-0.296881;9.538500\$S893\$Scom.google.android.gms.persistent\$  
\$S7MB\$S13MB\$S9MB\$S25MB\$S19MB\$S356584067413743  
20150907-05:33:17:126\$S17.25484552238640;48.09291989707780\$S-0.306458;-0.296881;9.538500\$S12578\$Scom.google.android.apps.fitness\$  
\$S2MB\$S2MB\$S13MB\$S6MB\$S8MB\$S20MB\$S356584067413743  
20150907-05:33:17:126\$S17.25484552238640;48.09291989707780\$S-0.306458;-0.296881;9.538500\$S12234\$Scom.facebook.katana\$S10067\$S0\*\$S  
\$S2MB\$S2MB\$S7MB\$S13MB\$S9MB\$S25MB\$S19MB\$S356584067413743  
20150907-05:33:17:126\$S17.25484552238640;48.09291989707780\$S-0.306458;-0.296881;9.538500\$S192.168.1.110\$S33135\$S216.58.209.202\$S4  
13743  
20150907-05:33:17:126\$S17.25484552238640;48.09291989707780\$S-0.306458;-0.296881;9.538500\$S192.168.1.110\$S53701\$S216.58.209.202\$S4  
13743  
20150907-05:33:17:126\$S17.25484552238640;48.09291989707780\$S-0.306458;-0.296881;9.538500\$S192.168.1.110\$S53902\$S173.194.122.9\$S44  
13743  
20150907-05:33:17:126\$S17.25484552238640;48.09291989707780\$S-0.306458;-0.296881;9.538500\$S192.168.1.110\$S58845\$S216.58.209.202\$S4  
13743  
20150907-05:33:17:126\$S17.25484552238640;48.09291989707780\$S-0.306458;-0.296881;9.538500\$S192.168.1.110\$S33328\$S216.58.209.202\$S4  
13743  
20150907-05:33:17:126\$S17.25484552238640;48.09291989707780\$S-0.306458;-0.296881;9.538500\$S192.168.1.110\$S43047\$S173.194.122.9\$S44  
13743  
20150907-05:33:17:126\$S17.25484552238640;48.09291989707780\$S-0.306458;-0.296881;9.538500\$S192.168.1.110\$S42252\$S195.28.96.197\$S22  
13743  
20150907-05:33:17:126\$S17.25484552238640;48.09291989707780\$S-0.306458;-0.296881;9.538500\$S192.168.1.110\$S35814\$S216.58.209.202\$S4  
13743  
20150907-05:33:17:126\$S17.25484552238640;48.09291989707780\$S-0.306458;-0.296881;9.538500\$S192.168.1.110\$S34596\$S64.233.167.188\$S5  
13743  
20150907-05:33:17:126\$S17.25484552238640;48.09291989707780\$S-0.306458;-0.296881;9.538500\$S192.168.1.110\$S54793\$S173.194.122.9\$S44  
13743  
20150907-05:33:17:126\$S17.25484552238640;48.09291989707780\$S-0.306458;-0.296881;9.538500\$S192.168.1.110\$S36396\$S173.194.122.30\$S4  
13743  
20150907-05:33:17:126\$S17.25484552238640;48.09291989707780\$S-0.306458;-0.296881;9.538500\$S192.168.1.110\$S54401\$S173.194.122.4\$S80  
13743  
20150907-05:33:17:126\$S17.25484552238640;48.09291989707780\$S-0.306458;-0.296881;9.538500\$S192.168.1.110\$S53969\$S216.58.209.164\$S4  
13743  
20150907-05:33:17:126\$S17.25484552238640;48.09291989707780\$S-0.306458;-0.296881;9.538500\$S192.168.1.110\$S43060\$S216.58.209.202\$S4  
13743  
20150907-05:33:17:126\$S17.25484552238640;48.09291989707780\$S-0.306458;-0.296881;9.538500\$S192.168.1.110\$S55917\$S173.194.122.9\$S44  
13743  
20150907-05:33:17:126\$S17.25484552238640;48.09291989707780\$S-0.306458;-0.296881;9.538500\$S192.168.1.110\$S42048\$S216.58.209.202\$S4  
13743  
20150907-05:33:17:126\$S17.25484552238640;48.09291989707780\$S-0.306458;-0.296881;9.538500\$S192.168.1.110\$S48140\$S173.194.122.9\$S80  
13743  
20150907-05:33:17:126\$S17.25484552238640;48.09291989707780\$S-0.306458;-0.296881;9.538500\$S192.168.1.110\$S33301\$S216.58.209.202\$S4  
13743  
20150907-05:33:17:126\$S17.25484552238640;48.09291989707780\$S-0.306458;-0.296881;9.538500\$S192.168.1.110\$S4437\$S216.58.209.202\$S44  
13743  
20150907-05:33:17:126\$S17.25484552238640;48.09291989707780\$S-0.306458;-0.296881;9.538500\$S192.168.1.110\$S53902\$S173.194.122.9\$S44  
13743  
20150907-05:33:17:126\$S17.25484552238640;48.09291989707780\$S-0.306458;-0.296881;9.538500\$S192.168.1.110\$S53902\$S173.194.122.9\$S44  
13743  
20150907-05:33:17:126\$S17.25484552238640;48.09291989707780\$S-0.306458;-0.296881;9.538500\$S192.168.1.110\$S53902\$S173.194.122.9\$S44  
13743

# METHODOLOGY AND THE FIRST INSIGHT INTO DATA

- **CRISP-DM** CROSS-INDUSTRY STANDARD PROCESS FOR DATA MINING
  - **DEVELOPMENT PHASE**
  - DEPLOYMENT PHASE
- **MOBILE LOGS = HUMAN-GENERATED CLASS**
  - HIGH POTENTIAL IN BIG DATA 3V (VOLUME, VELOCITY, VARIETY)
  - EXTREMELY NOISY FOR SPECIFIC RESEARCH PURPOSES
  - LOW OCCURENCES OF MALWARE RELATED ACTIVITIES
  - EVOLVING CHARACTERISTICS E.G. APPLICATIONS ON DEVICES ARE DYNAMICALLY INSTALLED/UNINSTALLED WITHOUT ANY RESTRICTIONS
- **PRIVACY PRESERVING DATA MINING DUE TO CONTRACT CONFIDENTIALITY AND PERSONAL SENSITIVE INFORMATION**





# EXPLORATORY DATA ANALYSIS (EDA)



# DATA PREPARATION

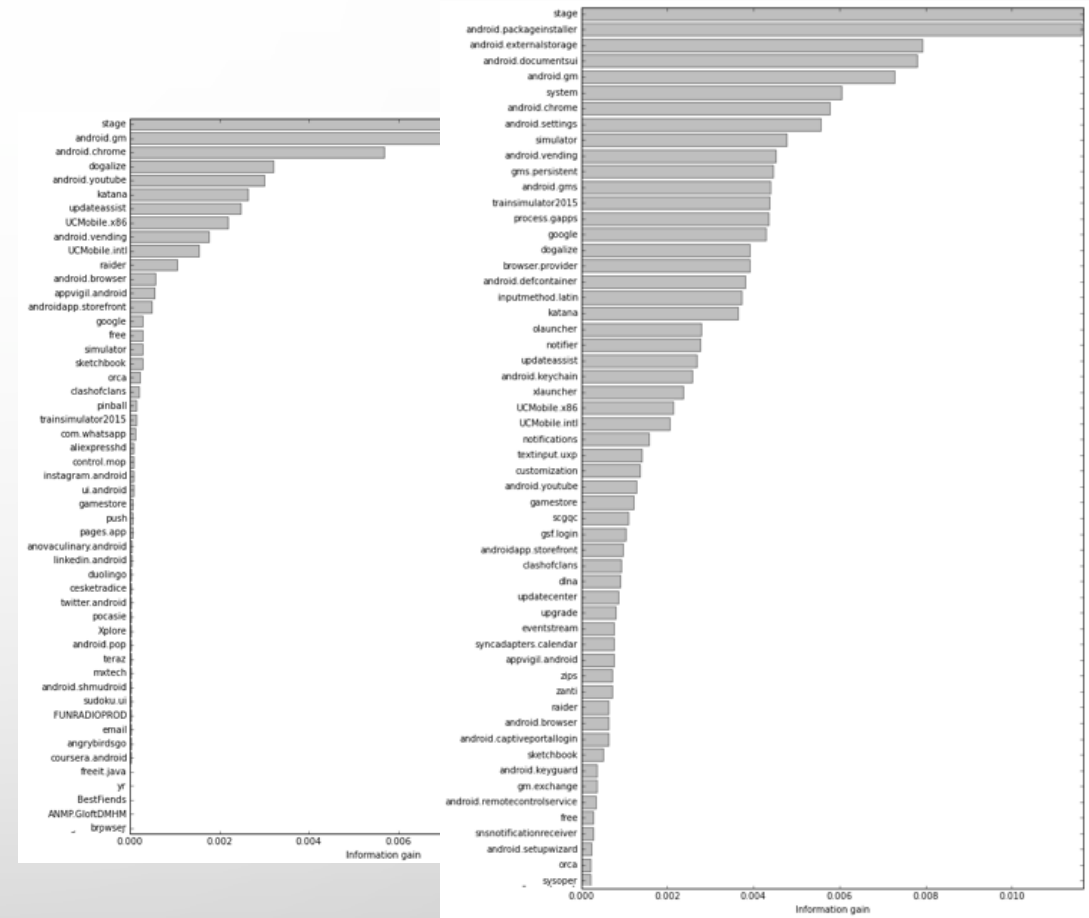
- DATA POOL FOR INCREMENTAL REALIZATION
- TIME-FRAMING SLIDING WINDOWS
- **NOISE FILTERING**, FILTERS CREATION
- **FEATURE ENGINEERING**
  - MOBILE DEVICE BEHAVIOR MODELING
  - FEATURE EXTRACTION FROM RAW MOBILE LOGS
  - FEATURE SELECTION
- **DATA TRANSFORMATION FROM RAW LOGS INTO MACHINE LEARNING DATA**

# FEATURE SELECTION

- ENTROPY

$$H(X) = \sum_{i=1}^n P(x_i) I(x_i) = - \sum_{i=1}^n P(x_i) \log_b P(x_i)$$

- **INFORMATION GAIN (IG)**
- INTRINSIC VALUE (IV)
- GAIN RATIO (IGR)       $IGR = IG/IV$
- THRESHOLDS, EXPLORATORY DATA ANALYSIS





# MODELING: TRAINING OF PREDICTIVE MODELS

- TRADITIONAL IN-MEMORY LEARNING: OFF-LINE LEARNING
- DISTRIBUTED LEARNING: COUPLED WITH INFRASTRUCTURE
- **INCREMENTAL LEARNING: ADAPTIVE ONLINE LEARNING, ONE-PASS OVER DATA**
- SUPERVISED LEARNING, BINARY CLASSIFICATION
  - SUPPORT VECTOR MACHINE (SVM)
  - LOGISTIC REGRESSION
  - FEED-FORWARD NEURON NETWORK
- HIGHLY IMBALANCED ML DATA CLASSES: BOOSTING OF POSITIVE EXAMPLES, REGULARIZATION TO PREVENT OVERFITTING

# RESEARCH RESULTS AFTER SIX MONTHS

- **MODEL PERFORMANCE EVALUATIONS**

- SEQUENTIAL AND CONCURRENT TRAINING DESIGNS
- ACCURACY, PRECISION, RECALL, F1, MCC, RMSE ... HOLD-OUT, K-FOLD, ...
- MATTHEWS CORRELATION COEFFICIENT (MCC)  $\langle 0.89, 0.92 \rangle$

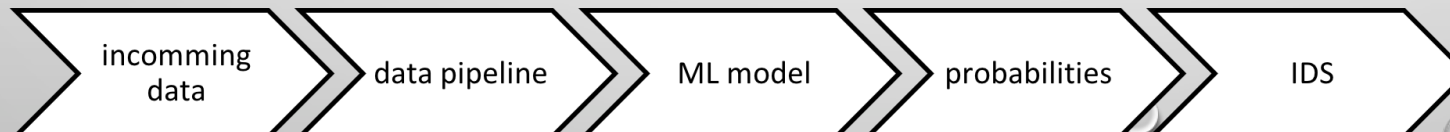
- **MALICIOUS BEHAVIOR DETECTION**

- SATISFIED RESULTS TO DISTINGUISH MALICIOUS BEHAVIOR FROM NORMAL ONE
- DETECTION FOR MALWARE TYPE WITHOUT BINDING TO CONCRETE SPECIFIC ONE

- **MORE TIME AND DATA FOR MORE REAL AND REAL-TIME DETECTIONS**

# BRIEFLY ABOUT OTHER RESEARCH QUESTS

- **EVENT SEQUENCES DETECTION USING NATURAL LANGUAGE PROCESSING (NLP) TECHNIQUES**
  - ATTEMPT OF SEMI-SUPERVISED LEARNING, BIG DATA APPROACH OF DATA TRANSFORMATION (APACHE PIG, HIVE)
  - N-GRAMS, PERPLEXITY MODEL, LOGARITHMIC PROBABILITY, SMOOTHING TECHNIQUES ...
- **PROCESS MINING AND VERIFICATION**
  - EVENTS IN FIXED TIME ORDERS WITH PETRI NET THEORY AND MACHINE LEARNING
  - RESULTS ARE ENCOURAGED, HOWEVER MORE DATA AND TIME IS NECESSARY
- **INTRUSION DETECTION SYSTEM (IDS)**
  - SUCCESSFULLY APPLIED FOR DDOS (DECODERS, RULES)
  - OFFLINE TESTING DUE TO TIME CONSTRAINTS
  - PIONEERING DEPLOYMENT FOR MALICIOUS BEHAVIOR DETECTION IN COLLABORATION WITH IDS



# MACHINE LEARNING AND DATA ANALYTICAL TOOLS

- **VOWPAL WABBIT:** JOHN LANGFORD, MICROSOFT/YAHOO! RESEARCH
  - OUT-OF-CORE FAST ONLINE LEARNING: TERA-FEATURE ( $10^{12}$ ) DATASET ON 1000 NODES IN 1H
  - THE HASHING TRICK 32-BIT MURMURHASH3
  - EXPLOITING MULTI-CORE CPU: PARSING OF INPUT AND LEARNING IN SEPARATE THREADS
  - COMPILED C++ CODE
- **SCIKIT-LEARN:** DAVID COURNAPEAU, INRIA (TELECOM PARISTECH) AND GOOGLE
  - WELL-MAINTAINED AND POPULAR ML TOOL, COMPREHENSIVE ALGORITHM LIBRARY
  - SINCE APRIL 2016, IT IS PROVIDED IN JOINTLY-DEVELOPED ANACONDA FOR CLOUDERA ON HADOOP
- **PYTHON, ANACONDA, NUMPY, SCIPY, PANDAS, MATPLOTLIB ...**

# INFRASTRUCTURES



- **HADOOP CLUSTER - HP BLADE SYSTEM**

- 1X SERVER AND 11X CLIENT 1TB HDD, NODE SPECIFICATION: 2X INTEL® XEON® PROCESSOR E5-2620 (15M CACHE, 2.00 GHZ, 7.20 GT/S INTEL® QPI, 6X CORES, 12X THREADS) + HYPERX THREADING (24 SIMULTANEOUS TASKS PER CLIENT), 32GB RAM, 1TB HDD

- **HPC CLUSTER – IBM SYSTEM**

- 52X IBM DX360 M3, 8X IBM DX360 M4 2X NVIDIA TESLA K20, 2X IBM DX360 M3 NVIDIA TESLA M2070, 2X X3650 M3 MANAGING SERVERS, 4X X3650 M3 DATA-MANAGING SERVERS, X3550 M4 SERVER, INFINIBAND 2X 40 GBPS (IN 52+2+2+4 NODES), 2X DS3512 WITH 72TB DISKS

# RESEARCH DIRECTIONS

- DATA MINING AND MACHINE LEARNING FOR BIG DATA (VOLUME, **VELOCITY, VARIETY, ...**)
  - USER BEHAVIORS, MONITORING SYSTEM, CYBER-SECURITY
  - TEXT, IMAGE AND SPEECH RECOGNITIONS
- RESEARCH DIRECTIONS TOGETHER WITH
  - ACCELERATED PROCESSING: GPU, APACHE SPARK/STORM, HPC AND CLOUD
  - BIG DATA PROCESSING, SEARCHING, KNOWLEDGE AND RELATION DISCOVERY, GRAPH TRAVELLING
  - MOBILE COMPUTING, INTERNET OF THINGS
  - **BASED ON DOMAIN DEMANDS**



# THANK YOU FOR YOUR ATTENTION

DETECTION OF POTENTIALLY DANGEROUS ACTIVITIES  
FROM LOGS OF MOBILE DEVICES  
USING MACHINE LEARNING TECHNIQUES

GIANG.UI@SAVBA.SK  
DEPARTMENT OF PARALLEL AND DISTRIBUTED INFORMATION PROCESSING  
INSTITUTE OF INFORMATICS  
SLOVAK ACADEMY OF SCIENCES



Tempest

IT makes sense



IBM Slovakia