

# Przygotowanie i analiza danych

# Wstęp

- dane z konkursu **KDD Mining Cup 1999**.
- KDD Cup jest światowym konkursem **Eksploracji Danych i Odkrywania Wiedzy** organizowanym przez *ACM SIGKDD - Special Interest Group on Knowledge Discovery and Data Mining*.
- **ZADANIE KONKURSU**
  - zbudowanie modelu prognostycznego (tj. klasyfikatora) zdolnego do rozróżniania połączeń “złych”, zwanych atakami lub włamaniami, oraz “dobrych” (normalnych) połączeń.

# Zestaw danych

- Zestaw danych obejmuje różnego rodzaju ataki symulowane na wojskowym środowisku sieciowym

# Symulowane ataki

- Denial of Service Attack (DoS - odmowa usługi)
- Remote to Local (R2L)
- User to Root (U2R)
- Probe

# Liczność zbioru danych

- 24 treningowe rodzaje ataków
- 14 rodzajów ataków w danych testowych
- dane testowe nie pochodzą z tego samego rozkładu probabilistycznego jak dane treningowe
- Zbiór danych zawiera:
  - 8 050 290 rekordów (1 173 MB)
  - 4 940 000 rekordów treningowych (734 MB)
  - 3 110 290 rekordów testowych (430 MB)

# Opis rekordu danych

- łącznie 42 zmienne
  - *duration* – czas połączenia w sekundach,
  - *protocol\_type* – rodzaj protokołu,
  - *service* – rodzaj usługi,
  - *flag* – połączenie normalne lub błędne,
  - *src\_bytes* – liczba bajtów przesłanych od źródła do celu,
  - *dst\_bytes* – liczba bajtów przesłanych od celu do źródła,
  - *land* – 1 jeżeli połączenie jest z tego samego hosta lub portu, 0 w przeciwnym wypadku,
  - *wrong\_fragment* – liczba błędnych fragmentów,
  - *urgens* – liczba pakietów z flagą „urgent”,
  - *hot* – liczba wskaźników „hot”,
  - *num\_failed\_logins* – liczba logowań zakończonych niepowodzeniem,
  - *logged\_in* – 1 jeżeli użytkownik zalogowany, 0 w przeciwnym wypadku,
  - *num\_compromised* – liczba skompromitowanych warunków,
  - *root\_shell* – 1 jeżeli uzyskano dostęp do powłoki Root, 0 w przeciwnym wypadku,
  - *su\_attempted* – 1 jeżeli użyto komendy „su root”, 0 w przeciwnym wypadku,
  - *num\_root* – liczba dostępów do konta „Root”,
  - *num\_file\_creations* – liczba operacji tworzenia pliku,
  - *num\_shells* – liczba monitów powłoki,
  - *num\_access\_files* – liczba operacji kontroli dostępu na plikach,
  - *num\_outbound\_cmds* – liczba komend wychodzących w sesji ftp,
  - *is\_host\_login* – 1 jeżeli zalogowany na konto „hot”, 0 w przeciwnym wypadku,
  - *is\_guest\_login* – 1 jeżeli zalogowano na koncie gościa, 0 w przeciwnym wypadku,
  - *count* – liczba połączeń do tego samego hosta w przedziale ostatnich 2 sekund,
  - *class* – klasa ruchu.

# Opis rekordu danych

- 5 zmiennych dotyczących tego samego hosta
  - *serror\_rate* – procent połączeń z błędem “SYN”,
  - *rerror\_rate* – procent połączeń z błędem “REJ”,
  - *same\_srv\_rate* – procent połączeń do tego samego serwisu,
  - *diff\_srv\_rate* – procent połączeń do innego serwisu,
  - *srv\_count* – liczba połączeń do tej samej usługi w przeciągu ostatnich 2 sekund.

# Opis rekordu danych

- Kolejnych 13 dotyczy tej samej usługi
  - *srv\_error\_rate* - procent połączeń z błędem „SYN”,
  - *srv\_rerror\_rate* - procent połączeń z błędem „REJ”,
  - *srv\_diff\_host\_rate* – procent połączeń do innego hosta,
  - *dst\_host\_rate* – procent połączeń do tego samego hosta docelowego,
  - *dst\_host\_count* – liczba połączeń do tego samego hosta docelowego, (35) *dst\_host\_srv\_count* – liczba połączeń do tego samego hosta docelowego i tej samej usługi,
  - *dst\_host\_same\_srv\_rate* – procent połączeń do tego samego hosta docelowego i do tej samej usługi,
  - *dst\_host\_diff\_srv\_rate* – procent różnych usług połączonych do danego hosta,
  - *dst\_host\_same\_src\_port\_rate* – procent połączeń do hosta z tym samym portem źródłowym,
  - *dst\_host\_srv\_diff\_host\_rate* – procent połączeń do tej samej usługi od różnych hostów,
  - *dst\_host\_error\_rate* – procent błędnych połączeń do hosta z błędem „SYN”,
  - *dst\_host\_srv\_error\_rate* – procent błędnych połączeń do hosta i tej samej usługi z błędem „SYN”,
  - *dst\_host\_rerror\_rate* – procent błędnych połączeń do hosta z błędem „REJ”.



# Przygotowanie danych

- pakiet STATISTICA 10 PL
- 10-procentowe zbiory danych ze zbiorów uczącego i testowego przygotowanego na konkurs KDD CUP 1999
  - *kddcup.data\_10\_percent.gz* - zbiór uczący zawierający **494 023** rekordów
  - *kddcup.testdata.unlabeled\_10\_percent* – zbiór testowy zawierający **311 030** rekordów

# Przygotowanie danych

- odczyt danych
- przypisanie danym właściwych etykiet
- import do projektu w STATISTICA 10

Podstawowe Edycja Wzrost Wstaw Format Statystyka Data Mining Wykresy Narzędzia Dane Pomoc

Wklej Wytnij Kopia Wyszukaj Wstaw Zamień Powtórz Znajdź Idź do Znajdź, zamień Wartościami losowymi Kopia w dół Kopia w prawo Wypełnianie Zmienne Przypadki Standaryzuj DDE OLE Obiekt Łączą

C:\Users\Agnieszka\Desktop\SSIED projekt\dane\zbior testowy\zbior testowy (311 030 rekordow).xlsx : kddcup.testdata.unlabeled_10_pe															
1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16
duration	protocol_type	service	flag	src_bytes	dst_bytes	land	wrong_fragment	urgent	hot	num_failed_logins	logged_in	num_compromised	root_shell	su_attempted	num_root
1	0	udp	private	SF	105	146	0	0	0	0	0	0	0	0	0
2	0	udp	private	SF	105	146	0	0	0	0	0	0	0	0	0
3	0	udp	private	SF	105	146	0	0	0	0	0	0	0	0	0
4	0	udp	private	SF	105	146	0	0	0	0	0	0	0	0	0
5	0	udp	private	SF	105	146	0	0	0	0	0	0	0	0	0
6	0	udp	private	SF	105	146	0	0	0	0	0	0	0	0	0
7	0	udp	domain	SF	29	0	0	0	0	0	0	0	0	0	0
8	0	udp	private	SF	105	146	0	0	0	0	0	0	0	0	0
9	0	udp	private	SF	105	146	0	0	0	0	0	0	0	0	0
10	0	tcp	http	SF	223	185	0	0	0	0	1	0	0	0	0
11	0	udp	private	SF	105	146	0	0	0	0	0	0	0	0	0
12	0	tcp	http	SF	230	260	0	0	0	0	1	0	0	0	0
13	0	udp	private	SF	105	146	0	0	0	0	0	0	0	0	0
14	0	udp	private	SF	105	146	0	0	0	0	0	0	0	0	0
15	1	tcp	smtp	SF	3170	329	0	0	0	0	1	0	0	0	0
16	0	tcp	http	SF	297	13787	0	0	0	0	1	0	0	0	0
17	0	tcp	http	SF	291	3542	0	0	0	0	1	0	0	0	0
18	0	tcp	http	SF	295	753	0	0	0	0	1	0	0	0	0
19	0	udp	private	SF	105	146	0	0	0	0	0	0	0	0	0
20	0	udp	private	SF	105	146	0	0	0	0	0	0	0	0	0
21	0	tcp	http	SF	268	9235	0	0	0	0	1	0	0	0	0
22	0	udp	private	SF	105	146	0	0	0	0	0	0	0	0	0
23	0	tcp	http	SF	223	185	0	0	0	0	1	0	0	0	0
24	0	tcp	http	SF	227	8841	0	0	0	0	1	0	0	0	0
25	0	tcp	http	SF	222	19564	0	0	0	0	1	0	0	0	0

# Przygotowanie danych – usunięcie powtórzeń

- Usunięcie z arkusza powtarzających się rekordów
- Minimalizacja ryzyka wystąpienia problemu, w którym model nauczy się dobrze rozpoznawać jedynie obiekt najliczniej występujący w zbiorze uczącym

# Przygotowanie danych – czyszczenie danych

- Usunięcie z arkusza rekordów niekompletnych

# Przygotowanie danych – losowanie danych

- **losowy podzbiór przypadków** – losuje 20% podzbioru rekordów ze zbioru uczącego oraz 5% podzbioru rekordów ze zbioru testowego

# Proces eksploracji danych zbioru uczącego

