

Makine Öğrenmesi için Büyük Veri Kümelerinin Kesitinin Çıkartılmasında Yeni Teknikler:

Android Mobil Kötücül Veri Kümelerinin Özlü Bir
Gözden Geçirilmesi*

Gürol Canbek

Orta Doğu Teknik Üniversitesi

<http://gurol.canbek.com>

Şeref Sağıroğlu

Gazi Üniversitesi

Tuğba Taşkaya Temizel

Orta Doğu Teknik Üniversitesi

*Gürol Canbek, Seref Sagiroglu, Tugba Taskaya Temizel, “New Techniques in Profiling Big Datasets for Machine Learning with A Concise Review of Android Mobile Malware Datasets”, International Congress on Big Data, Deep Learning & Fighting Cyber Terrorism (IBIGDELFT 2018), 3–4 December 2018: IEEE



INTERNATIONAL
CONGRESS
ON **BIG DATA**
DEEP LEARNING
and **FIGHTING CYBER TERRORISM**
DECEMBER 3- 4, 2018

Özet

SORUN:

- Büyük veri boyutları:
 - Çokluk, Çabukluk, Çeşitlilik
 - Dürüstlük, Değer, Değişkenlik
- Makine öğrenmesinde büyük verinin kullanımı?

KATKI:

- Bu çalışma veri kümelerinin kesitlerinin çıkartılması için
 - 14 kıstas ile 4 farklı teknik önermektedir.
 - İlk defa: Büyük veri boyutları ile eşleme

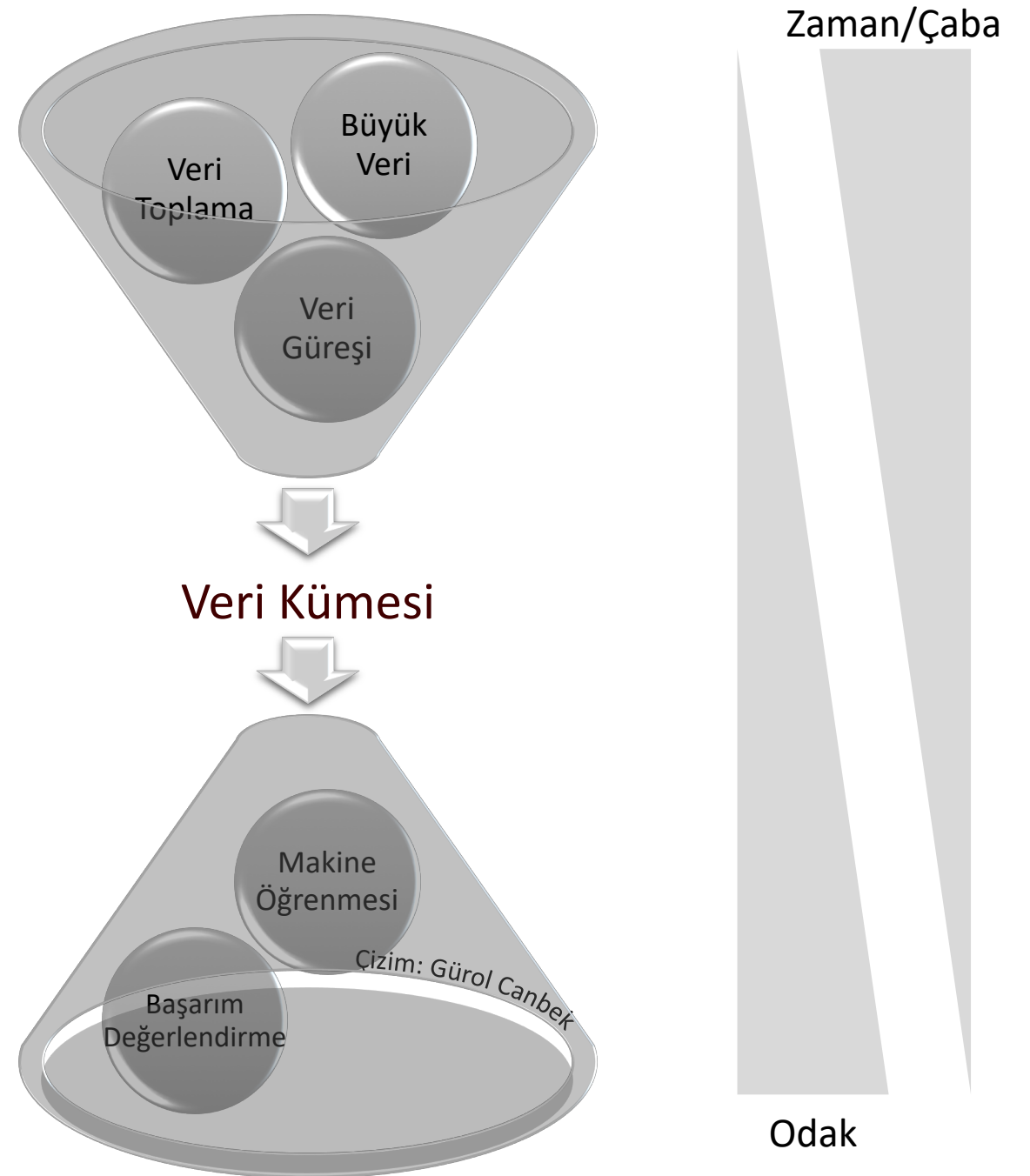
DENEY:

- Android Malware Genome Project, Drebin, Android Malware Dataset, Anroid Botnet ve Virus Total 2018

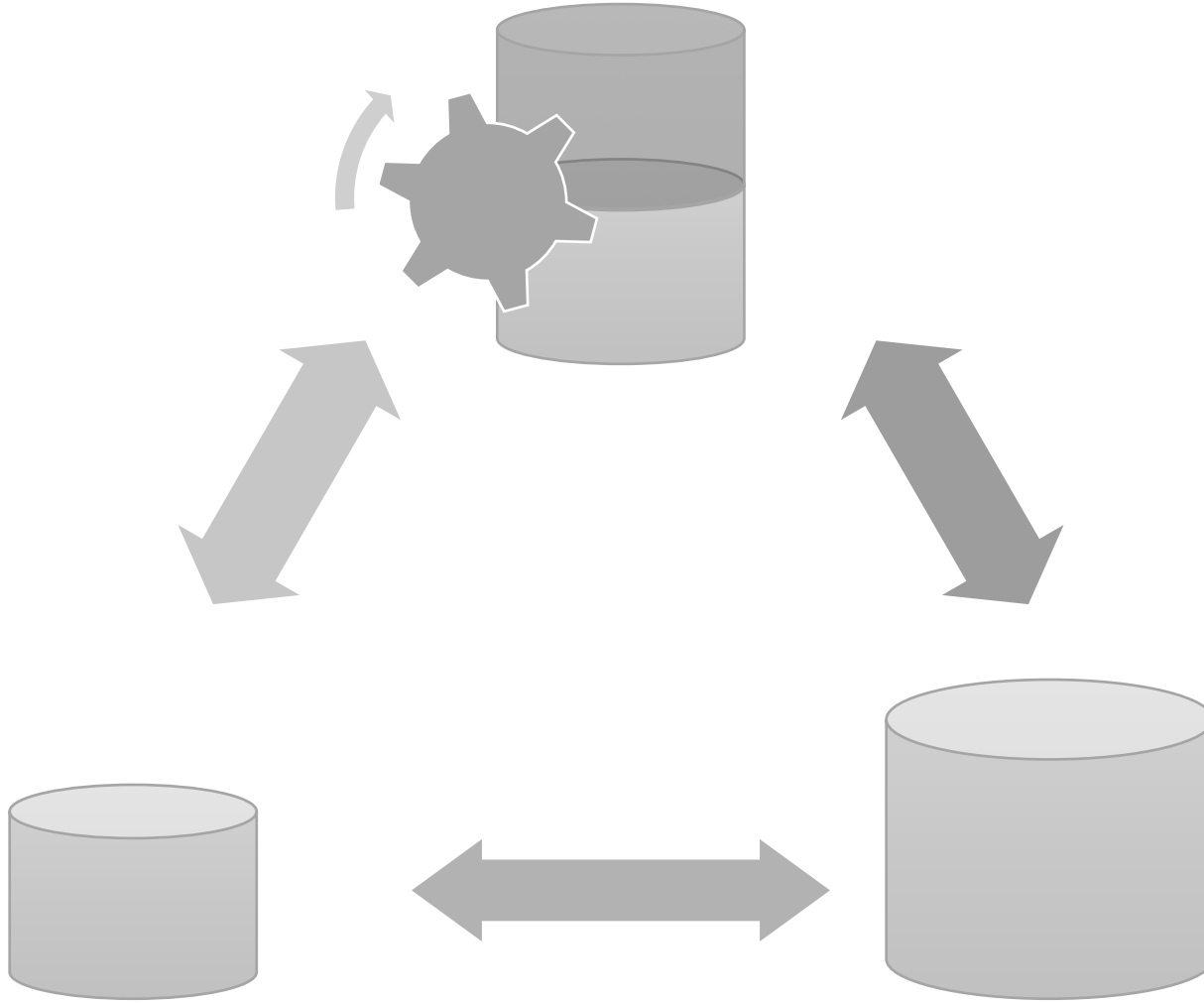
SONUÇ:

- Yöntem, veri kümeleri hakkında karşılaştırılabilir önemli içgörü sunmaktadır.
- Büyük verinin daha görünür, kaliteli ve özümsemiş olmasına uygulanabilir bir katkı sağlamaktadır.

Büyük Veri ve Makine Öğrenme İş Akışı



Veri Kümelerinin Kalitesi



- Eldeki veri kümesi:
 - Kalite artırma
- Farklı veri kümeleri
 - Kalite karşılaştırma

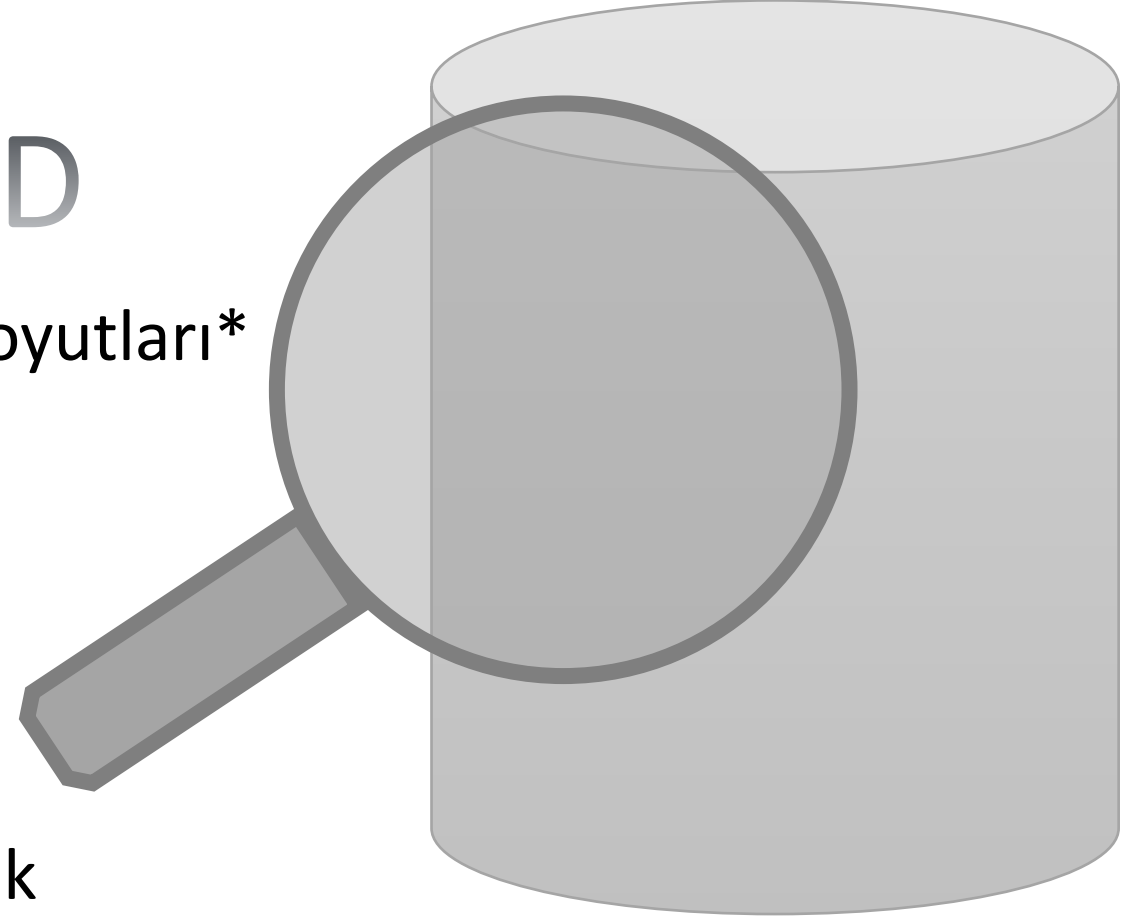
Veri Kümesi Kesit Çıkarma

- İlk planda
- Sistematik bir şekilde
- Önemli içgörü

3Ç, 4D

Büyük Veri Boyutları*

- Çokluk
- Çabukluk
- Çeşitlilik
- Dürüstlük
- Değer
- Değişkenlik
- Dayanak



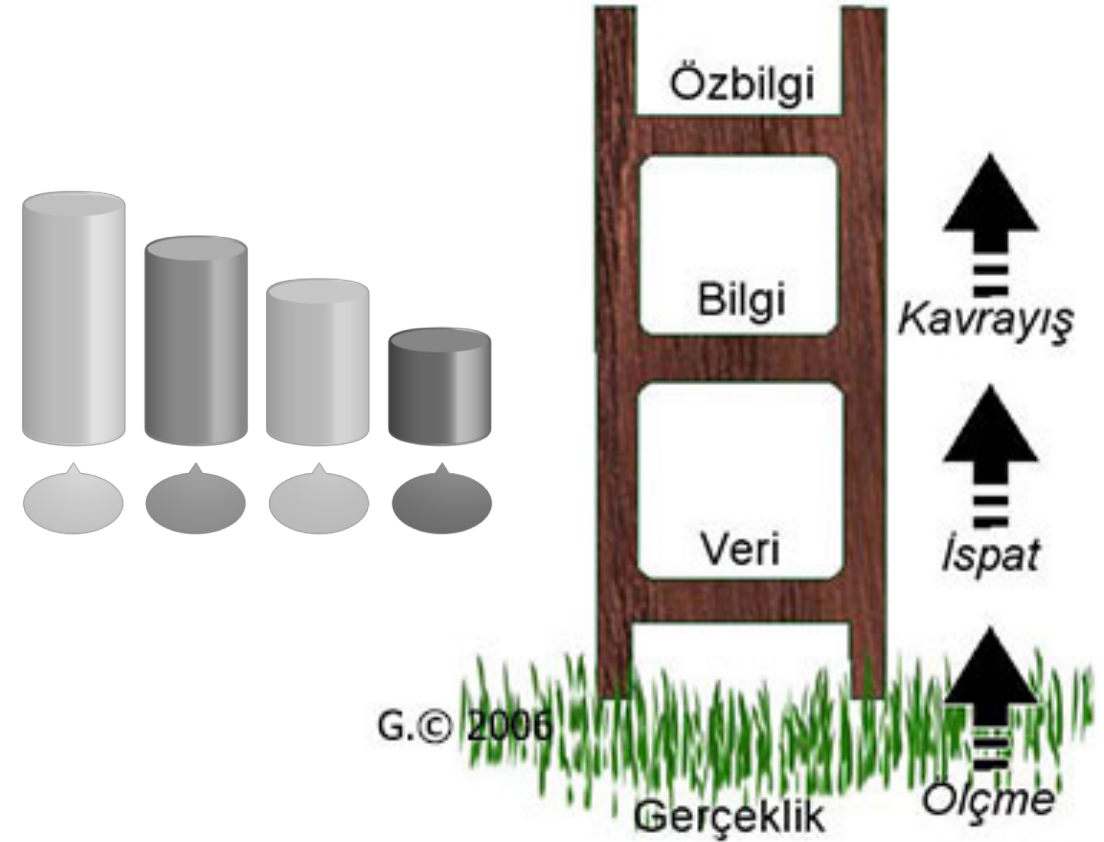
* Volume, Velocity, Variety, Veracity, Value, Variability, Venue için Türkçe karşılık önerileri: Gürol Canbek, Kasım 2018

Gürol Canbek, Seref Sagiroglu, Tugba Taskaya Temizel, “New Techniques in Profiling Big Datasets for Machine Learning with A Concise Review of Android Mobile Malware Datasets”, International Congress on Big Data, Deep Learning & Fighting Cyber Terrorism (IBIGDELFT 2018), 3–4 December 2018: IEEE



Literatür ve Veri/Bilginin Kalitesi

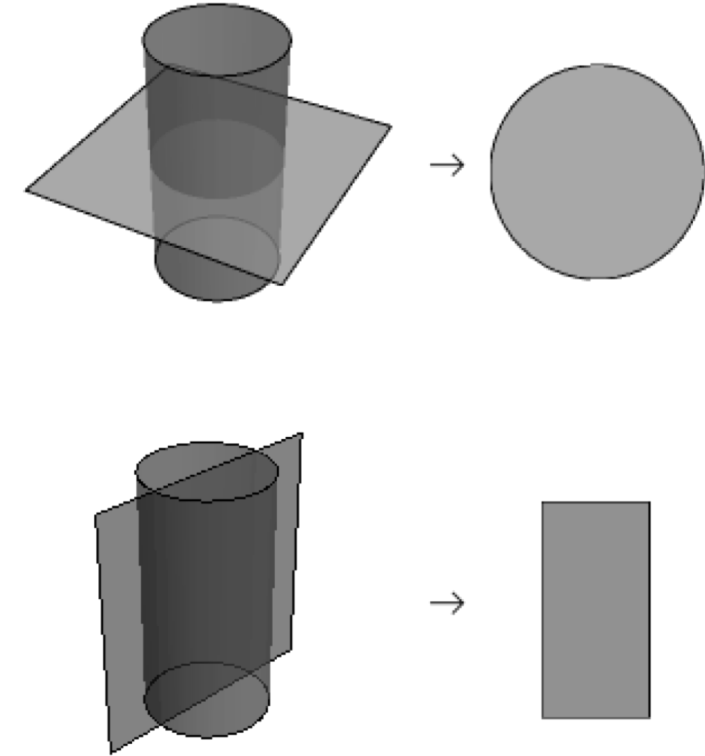
- Literatürde daha önce ele alınmış.
 - İçsel: Doğruluk, tutarlılık, ...
 - Bağlamsal: İlgili, tamlık, ...
 - Gösterimsel: Anlaşılabilirlik, biçim, ...
 - Erişimsel: Kullanabilirlik, güvenlik, ...
- İstatistiksel yöntemler
 - Tanımlayıcı istatistikler
- Sahaya özel kalite bakış açısı
 - Uzamsal veri kümeleri
 - Sağlık verileri
 - Doğa olayları verileri gibi



Canbek, G., Sağiroğlu, Ş. (2006). **Bilgi, Bilgi Güvenliği ve Süreçleri Üzerine Bir İnceleme**. *Politeknik Dergisi*, 9(3), 165–174.

Kesit Çıkarma

- Verinin veya veri kümesinin belirli bir bakış açısından tarif edilmesi.
- Makalede:
 - İlişkisel Veritabanı Sistemleri (RDBMS)
 - Yapısal Sorgulama Dilleri (SQL)
 - İnternet üzerindeki veriler
- Veri kümelerinin makine öğrenmesi açısından kalitesi?



Android Mobil Kötücül Yazılım Veri Kümeleri

Veri Kümesi (Kısaltma)	Yıl	Etüt*
Android Malware Genome Project (AMGP)	2013	%65
Drebin	2014	%22
Android Botnet (ABot)	2015	
Android Malware Dataset (AMD)	2017	
VirusTotal Academic Malware Samples (VT2018)	2018	

* 2009-2018 atmış çalışma. İlave malzeme: github.com/gurol/dsprofiling

Kesit Çıkarma Zemini

- MalWareHouse
 - Android kötücül yazılım çözümleme zemini
 - Farklı veri kümesi havuzlarının veri ambarı yaklaşımı ile birleştirilmesi
 - Kesit çıkarma bilgilerinin elde edilmesi ve görselleştirilmesi
 - Python, MongoDB, R

MalWAREHOUSE



Önerilen Kesit Çıkarma Zümreleri

1. Temel
2. Zaman çizgisi
3. Mükerrer örneklem
4. Yoğunluk/seyreklilik

Temel Kesit Çıkarım

- Üst seviyeden ilk içgörü

Kıstas	Büyük Veri Boyutu	AMGP	Drebin	AMD	ABot	VT2018
Örneklem uzayı büyüklüğü (m)		1,260	5,555	23,743	1,929	4,725
Öznitelik uzayı büyüklüğü (n)	Çokluk	65	94	105	78	111
Fizikî büyüklük (GB)		1.5	6.8	58.1	2.6	18.4
Kötücül ailesi*	Çeşitlilik	49	> 20	71	14	Yok
Kötücül başka biçimleri*	Dürüstlük Dayanak	Yok	Yok	135	Yok	Yok

* Sahaya bağlı kıstaslar

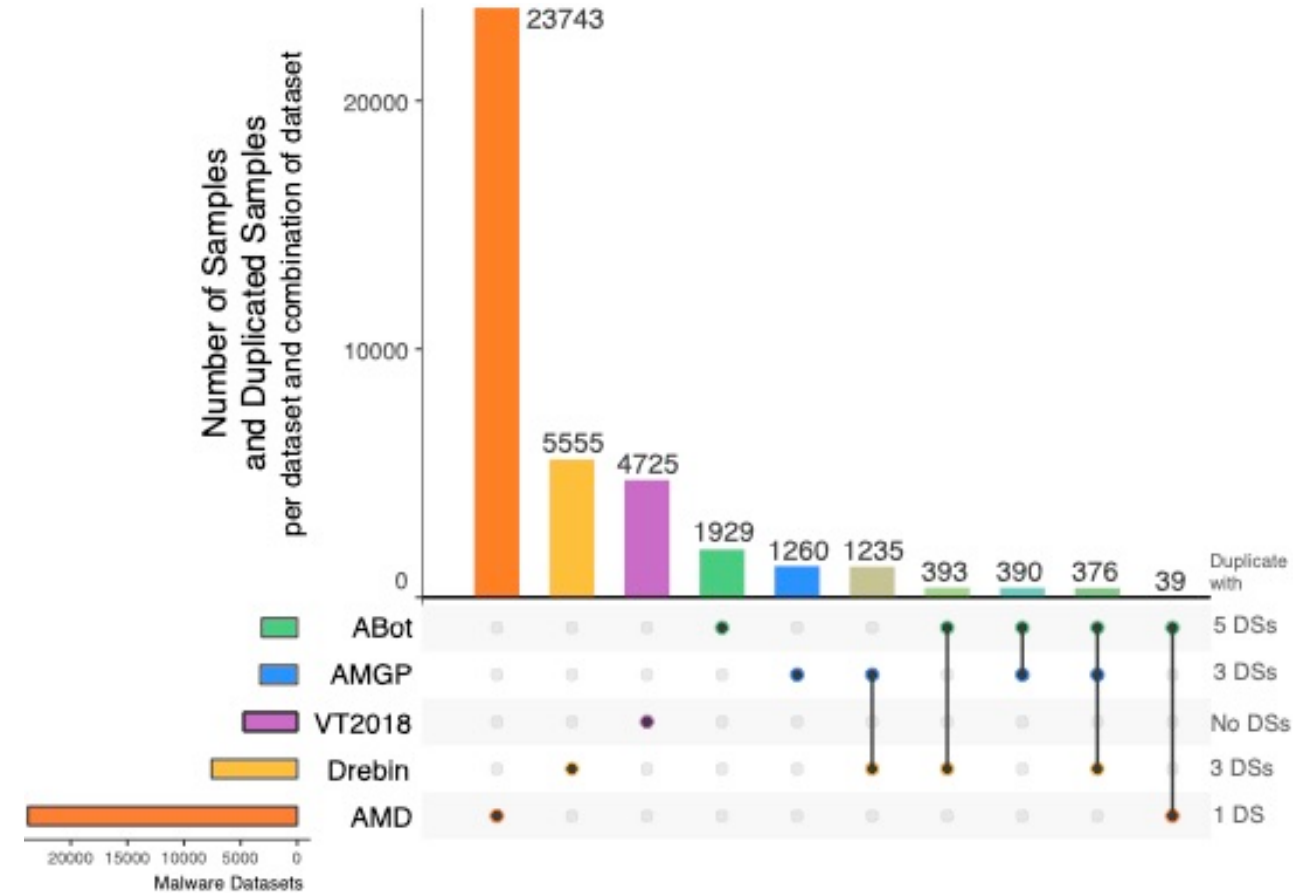
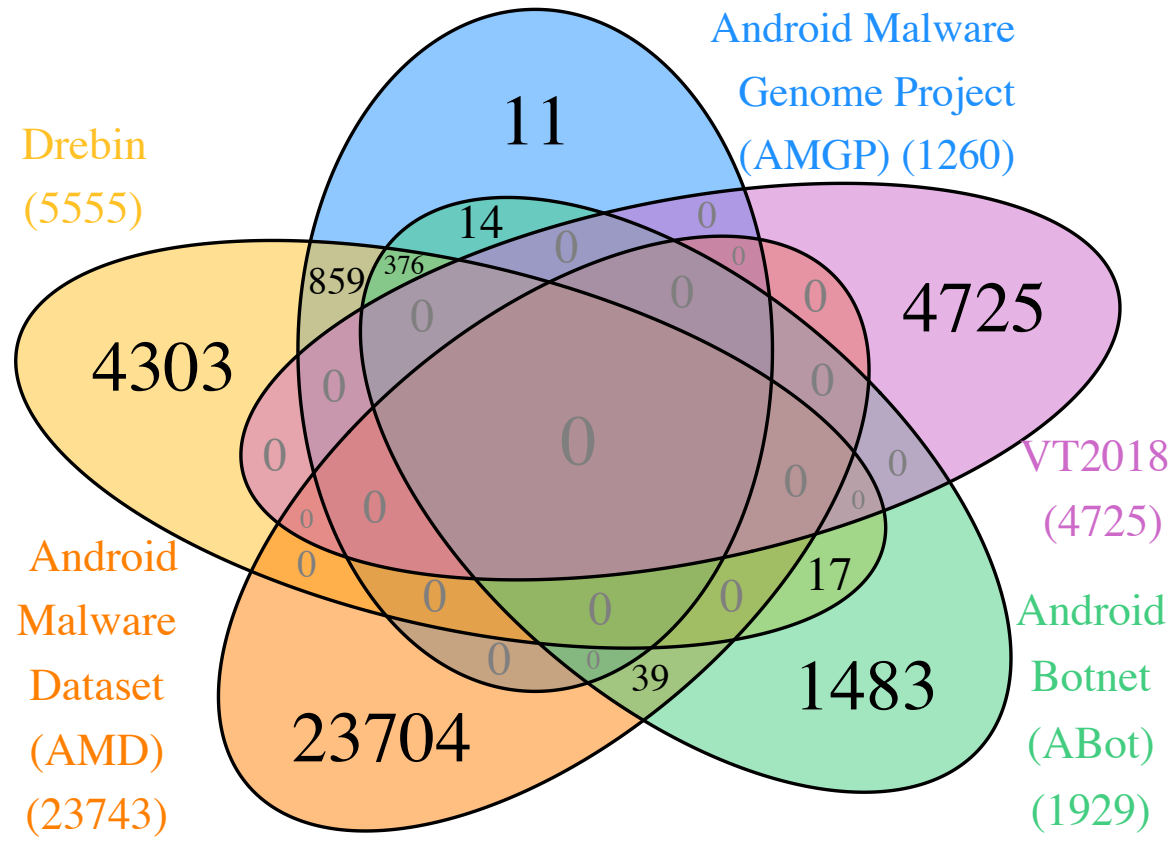
Zaman Çizgisi Kesit Çıkarımı

- **Yaş:** En genç ile en yaşlı örneklem arasındaki fark
- **Tazelik:** En genç örneklemden itibaren geçen zaman aşımı
- API seviyesi menzili (sahaya bağlı)
- Çabukluk, Dürüstlük, Değişkenlik
- Önerilen gösterim yöntemi:

Veri Kümesi	En yaşlı	En genç	API	2008	2009	2010	2011	2012	2013	2014	2015	2016	2017	2018	Tazelik
AMGP	2008	2011	12	2.9 yaşında											-7 yıl
Drebin	2008	2012	16	4											-6
AMD	2008	2016	25	8											-2
ABot	2009	2015	21		7										-3
VT2018	2009	2018	27		9										0

Mükerrer Örneklem Kesit Çıkarımı

- Değer, Dayanak



Yoğunluk/Seyreklik Kesit Çıkarımı

$$yoğunluk = \frac{\text{gözlemlenen öznitelik sayısı}}{m \times n}$$

$$seyreklik = \frac{\text{gözlemlenmeyen öznitelik sayısı}}{m \times n} = 1 - yoğunluk$$

$$yoğunluk_n = \frac{n}{n_{max}}$$

$$seyreklik_n = 1 - density_n$$

Kıstas	AMGP	Drebin	AMD	ABot	VT2018
Yoğunluk	18%	10%	11%	16%	15%
Seyreklik	82%	90%	89%	84%	85%
Öznitelik uzayı yoğunluğu	43%	62%	70%	52%	74%
Öznitelik uzayı seyrekliği	57%	38%	30%	48%	26%
Örneklem uzayı büyüklüğü (m)	1,260	5,555	23,743	1,929	4,725
Öznitelik uzayı büyüklüğü (n)	65	94	105	78	111

Çabukluk, Dürüstlük, Değer

Ayrıca, öznitelik uzayı yoğunluk ve seyrekliği (evrensel veya en büyük öznitelik uzayında göre)

Alınan Sonuçlar

- 14 kıstasın toplu değerlendirilmesi
- Yüksek kesit veri kümeleri:
 - Android Malware Dataset (AMD)
 - VirusTotal Academic Malware Samples (VT2018)
- Düşük kesit veri kümesi
 - Android Malware Genome Project (AMGP)

Kesit	Büyük Veri Boyutu ⁽¹⁾	Kıstas	AMGP	Drebin	AMD	ABot	VT2018
Temel	Çokluk	Örneklem uzayı büyüklüğü (m)	<u>1260</u>	5555	23743	1929	4725
		Öznitelik uzayı büyüklüğü (n) ⁽²⁾	65	94	105	78	111
		Fizikî büyüklük (GB)	1.5	6.8	58.1	2.6	18.4
	Çeşitlilik, Dürüstlük, Dayanak	Kötücül ailesi ⁽³⁾	49	> 20	71	14	N/A
		Kötücül başka biçimleri ⁽³⁾	N/A	N/A	135	N/A	N/A
Zaman Çizgisi	Dürüstlük, Değişkenlik, Çabukluk	Yaş (yıl)	<u>2.9</u>	4.1	7.6	6.6	9
		Tazelik (yıl)	<u>-7</u>	-6	-2	-3	0
		API-seviyesi menzili ⁽³⁾	<u>12</u>	16	25	21	27
Mükerrer Örneklem	Değer, Dayanak	Özgün örneklem sayısı	<u>11</u>	4303	23704	1483	4725
		Mükerrer örneklem sayısı	1249	1252	39	446	0
Yoğunluk / Seyreklik	Çokluk, Dürüstlük, Değer	Yoğunluk ⁽⁴⁾	18%	10%	11%	16%	15%
		Seyreklik ⁽⁴⁾	82%	90%	89%	84%	85%
		Öznitelik uzayı yoğunluğu	<u>43%</u>	62%	70%	52%	74%
		Öznitelik uzayı seyrekliği	57%	38%	30%	48%	26%
Genel sezgisel veri kümesi kesiti			Düşük (-5)	Normal (1)	Yüksek (4)	Normal (1)	Yüksek (4)

(1) Yedi boyut kapsanmıştır: İngilizce: value, variability, variety, velocity, venue, veracity, volume

(2) Öznitelik uzayı yoğunluğu/seyrekliği kıstasında değerlendirilmeye alınmıştır.

(3) Sahaya özel kıstas ancak diğer sahalara uygun bir şekilde uyarlanabilir

(4) Bilgi verici, herhangi bir üstünlük arz etmemektedir

Değerlendirme ve Sonuçlar

Bu çalışmada;

- Dört zümrede ve 14 kıstastan oluşan veri kümesi kesit çıkartımı yöntemi önerilmiştir.
- İlk defa büyük veri boyutları veri kümesi kalitesi anlamında kıstaslarla eşleştirilmiştir.
- Kesitlerin anlaşılabilirliğini artırmak adına farklı görselleştirme yaklaşımları tanıtılmıştır.
- Yöntem, beş Android kötücü yazılım veri kümesinde denenmiş ve veri küme kesitlerinin karşılaştırılması yapılmıştır.

SONUÇ

- Veri kümelerinin ilk planda yüksek/düşük kesit şeklinde etiketlenmesi ilk planda elverişli bir içgörü sağlayabilir.
- Çalışmalarda kullanılan veri kümelerinin kesiti, başarımlar raporumada ifade edilebilir.
- Araştırmacıları kaliteli veri kümesi elde etmek için teşvik sağlar.
- Yeni çıkan veri kümesi o saha içindeki veri kümelerinin yeniden kesitinin çıkartılmasını gerektirir.
- Bir veri kümesinin kesitinin çıkartılması eksik yönlerinin giderilmesini sağlar.

Teşekkürler

Gürol Canbek

<http://gurol.canbek.com/Publications>

Bildiri ile ilgili ilave bilgi, betik ve malzeme

<https://github.com/gurol/dsprofiling>