

STATISTICAL DETECTION OF COLLECTIVE DATA FRAUD

Name of author

ABSTRACT

Statistical divergence is widely applied in multimedia processing, basically due to regularity and explainable features displayed in data. However, in a broader range of data realm, these advantages may not out-stand, and therefore a more general approach is required. In data detection, statistical divergence can be used as an similarity measurement based on collective features. In this paper, we present a collective detection technique based on statistical divergence. The technique extracts distribution similarities among data collections, and then uses the statistical divergence to detect collective anomalies. Evaluation shows that it is applicable in the real world.

Index Terms— Statistical divergence, collective, fraud, detection

1. INTRODUCTION

Statistical divergence is widely applied in multimedia processing. Prevalent applications include multimedia event detection [1], content classification [2, 3] and qualification [4, 5]. It has been attracting more attention since the dawn of big data era, basically due to regularity and interpretable features displayed in the data. However, in a broader range of data realm, these advantages may not out-stand (e.g. in online sales data records). It requires a more general approach.

Data manipulation from outside hackers composes another potential threat of data frauds. *Data Manipulation* here, according to a NSA definition, refers to behaviours which “change information contained in those systems, rather than stealing data and holding it for ransom”. In 2013, hackers from Syria put up fake reports via Associated Press’ Twitter account and caused a 150-point drop in the Dow [6].

It is hard to detect a single record that is altered but still remains in correct value scopes, but if sufficient data records are altered to change a final decision, we can still detect malicious data manipulation behaviours. According to our observation, typical manipulations on numerical data will lead to a drift or distortion of its original distribution. To address problems caused by data manipulation, we proposed a novel technique which sorts out manipulated data collections from normal ones by adopting statistical divergence. In this paper, we focus on a concrete data manipulation problem: click farming in online shops, and try to apply our technique to pick out

dishonest behaviours. Our technique maps data collections to points in distribution spaces and reduce the problem to classical point anomaly detection. Optimizations estimate ground truth, mapping each data collection into a single real number within a definite interval. Then a Gaussian classifier can be applied to detect outliers derived from manipulated data. To automatically calculate adaptive threshold for the classifier, we keep two evidence sets for both normal points and anomalies, taking advantage of the property provided by statistical divergence. In the dynamic environments, these evidence sets act intuitively as slide windows and keep up to the evolving features in dynamic scenarios. Our contribution includes: 1) A brief review on data fraud detection and a study on the problem of click farming; 2) Detailed description of both basic and optimized framework of our technique, resolving several technical difficulties such as automated adaptive threshold; 3) Comprehensive experiments that test efficiency of our technique and a comparison with previous work on similar topic.

The rest of the paper is organised as follows: Section 2 states related work on data anomaly detection and briefly introduces click farming. Details of proposed technique are introduced in section 3. Then section 4 presents evaluation results and further findings of the algorithm. Finally, the paper is concluded in section 5.

2. RELATED WORK

2.1. Data Anomaly Detection

Statistical divergence was applied mainly as classifiers on multimedia content [3], especially as kernels in SVMs [2]. As a similarity measurement, it can also be used in qualitative and quantitative analysis in image evaluation [4, 5]. [1] adopted divergence to detect events in multimedia streams.

To detect collective anomalies, [7] adopts the *ART (Adaptive Resonance Theory)* neural networks to detect time-series anomalies. *Box Modeling* is proposed in [8]. *Longest Common Subsequence* was leveraged in [9] as similarity metric for symbolic sequence. Markovian modeling techniques are also popular in this domain [10, 11, 12]. [13] depicts groups in social media as combinations of different “roles” and compare groups according to the proportion of each role within each group.

Wang et al. proposed a technique, *Multinomial Goodness-of-Fit* (MGoF), to analyze likelihood ratio of distributions via

Kullback-Leibler divergence, and is fundamentally a hypothesis test on distributions [14]. MGoF is the best competitor out of the similar techniques, and we use it as our baseline against our approach.

2.2. Real World Problem: Click Farming Detection

Click farming is the behaviour that online sellers use a large number of customer accounts to create fake transaction records and give high remarks on products. There are two types of click farming behaviours: Centralized and Equalized. Equalized click farming refers to scenarios where behaviours are well organised while centralized does not. Current detection techniques for click farming mainly focus on user behaviours. Those techniques require platforms to keep records on user features. However, the detection can be easily bypassed by trained workers and some well programmed applications.

Although it is hard to classify users as honest or malicious, we can still find clues from the sellers' aspect. No matter how much alike between honest users and malicious workers, the fake transaction records will always cause a bias or distortion of the original transaction distribution. Thus, if we can measure the similarity between different transaction distributions, there is still a chance for us to detect dishonest sellers.

3. STATISTICAL DETECTION

3.1. Statistical Divergence Detection with Reference (SDD-R)

Statistical divergence only provides a distance between two or more distributions. In a set of data collections, we can only draw a complete graph where nodes denote data collections and edges refer to the symmetric divergence between two connected nodes. From the graph we can find some points that have apparently larger distances with most of other points and return them as anomalies. This may work if anomalous nodes do not compose a large proportion. However the procedure will be too complicated to work out with large amounts of data. If it is assured that data collections form only one cluster, some optimizations can be applied to reduce complexity.

Alternatively we can provide a frame of reference that generates absolute coordinates rather than the relative ones. This optimization is feasible if data collections form one single cluster in distribution space. This is true in most reality scenarios given that distribution is adopted to depict a macro property which comes out as one universal conclusion. In other words, if multiple distributions are used to describe subgroups of entire sample space, then a conclusive one can be obtained by averaging all these sub-distributions. Therefore, we can use an estimate cluster center as reference and test distances between the reference and each other data collections (Algorithm 1), yielding absolute distances.

Algorithm 1 SDD-R

Input: Data Collections $\mathbb{D} = \{D_1, \dots, D_n\}$

Input: Estimated anomalous probability α

Output: Anomalous Data Collections

```

1: for  $i \leftarrow 1$  to  $n$  do
2:    $P_i \leftarrow$  the distribution of  $D_i$ 
3: end for
4:  $P_R \leftarrow \frac{1}{n} \sum_{i=1}^n P_i$ 
5: for  $i \leftarrow 1$  to  $n$  do
6:    $d_i \leftarrow D(P_i || P_R)$ 
7: end for
8:  $\mathcal{N}(\mu, \sigma) \leftarrow$  Gaussian distribution estimated by  $d_i$ 
9: return  $\{D_i | \frac{d_i - \mu}{\sigma} > 3\}$ 

```

Distribution of all divergences against the reference can be approximated as a Gaussian distribution though the true one may differ a little more from the standard Gaussian than the expected estimation error. That is due to the unknown randomness within real world data. Few assumptions can be applied in real world data sets, no mention that data volume is sometimes relatively low. This topic is out of the domain discussed in this paper and we here only introduce the technique instead of the specific distribution model. Certainly, if stronger assumptions can be included to provide a more precise model, this component in the framework can be replaced to give better results. For the simplicity of our proposal, we deem the distributions of divergences to be Gaussian.

By this approach, time complexity can be reduced from quadratic to linear. Fig. 2 in Section 4.2 demonstrates the result of the above process. Red distribution refer to the distances calculated from normal data collections, blue and green ones are from click-farmed data collections. Clearly, distances of normal data collections assembles together around a small value while anomalous ones lay around a larger distance value.

3.2. Optimization: Statistical Divergence Detection with Evidence(SDD-E)

It is possible to further optimize SDD-R if we can provide this algorithm with evidence (Algorithm 2).

Evidences enables the algorithm to not only refine estimation of real distribution but also build knowledge of anomalous collections, which is similar to the parameter estimation within a certain sample set.

According to the property of statistical divergence, we can infer that the true distribution of divergences calculated from normal data collections are close to but not exactly a Gaussian distribution $\mathcal{N}(\mu, \sigma)$ since for each point, there are both definite upper and lower bounds instead of infinities. Therefore, μ should be slightly larger than zero ($\mu = 0 \iff P_i = P_j, \forall P_i, P_j \in \mathbb{E}_N$, for real world data sets, this is highly unlikely). Time complexity for this algorithm is still linear but

Algorithm 2 SDD-E

Input: Evidence set with normal data collections $\mathbb{E}_N = \{N_1, \dots, N_n\}$

Input: Evidence set with anomalous data collections $\mathbb{E}_A = \{A_1, \dots, A_m\}$

Input: Estimated anomalous probability α

Input: New data collection $\mathbb{D} = \{D_1, \dots, D_l\}$

Output: Anomalous data collections in \mathbb{D}

```
1: for  $i \leftarrow 1$  to  $n$  do
2:    $P_{N_i} \leftarrow$  distribution of  $D_{N_i}$ 
3: end for
4: for  $i \leftarrow 1$  to  $m$  do
5:    $P_{A_i} \leftarrow$  distribution of  $D_{A_i}$ 
6: end for
7:  $P_R \leftarrow \frac{1}{n} \sum_{i=1}^n P_{N_i}$ 
8: for  $i \leftarrow 1$  to  $n$  do
9:    $d_{N_i} \leftarrow D(P_{N_i} || P_R)$ 
10: end for
11: for  $i \leftarrow 1$  to  $m$  do
12:    $d_{A_i} \leftarrow D(P_{A_i} || P_R)$ 
13: end for
14:  $\mathcal{N}_N(\mu_N, \sigma_N) \leftarrow$  normal distribution estimated from  $\{d_{N_1}, \dots, d_{N_n}\}$ 
15:  $\mathcal{N}_A(\mu_A, \sigma_A) \leftarrow$  normal distribution estimated from  $\{d_{A_1}, \dots, d_{A_m}\}$ 
16:  $T \leftarrow$  proper threshold derived from  $\mathcal{N}_N, \mathcal{N}_A$  and  $\alpha$ 
17: for  $i \leftarrow 1$  to  $l$  do
18:    $P_i \leftarrow$  distribution of  $D_i$ 
19:    $d_i \leftarrow D(P_i || P_R)$ 
20: end for
21: return  $\{D_i | d_i > T\}$ 
```

with a larger coefficient.

For certain divergence, it is possible to compare similarity from one distribution against multiple others, such as Jensen-Shannon Divergence. Although it reduces time complexity, it sacrifices unaffordable accuracy because divergence among multiple distribution dilutes differences. Take JSD as an example, suppose $P(1) = P(2) = P(3) = \frac{1}{3}$ and $Q(1) = \frac{1}{6}, Q(2) = \frac{1}{3}, Q(3) = \frac{1}{2}$, then $JSD(P||Q) \approx 0.033$ and $JSD(P, P, P, Q) \approx 0.024$.

This algorithm can be slightly modified to deal with concept drift(for example, trading trend changes over time for online shops as they are often in the process of expanding or dwindling) by turning the two evidence sets as sliding windows and adopting certain update strategies such as *Least Recently Used*(LRU). Time complexity for this optimization is $O(n \cdot (|\mathbb{E}_N| + |\mathbb{E}_A|) \cdot T_D)$, where T_D denotes time complexity of divergence calculation.

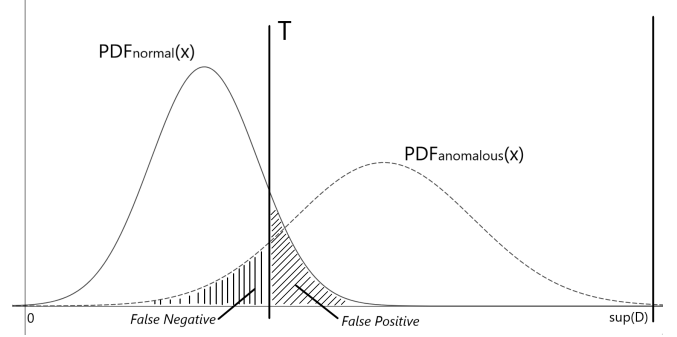


Fig. 1: Optimal threshold should minimize the size of shadow under curve.

3.3. Threshold

One important factor in algorithm SDD-E is the value of threshold. A naive but prevalent approach is to set a fixed value as the threshold(As is shown in Algorithm 1). However, a fixed threshold requires specific analysis in the certain scenario, manual observations and tuning of parameters, which involves lots of human labour.

Fortunately, applying divergence as the distance measurement among data collections provides a fine property. With a reference distribution, divergences of normal data collections form a quasi-Gaussian distribution. The same applies to those anomalous ones.

Moreover, as is verified in experiments, the meta-distribution of anomalous data collections lies in the right-hand-side to the normal one on the real line. As shown in Fig 1, the black curve displays the probability density function (PDF) fitting those divergences calculated from normal data collections; the blue curve displays the PDF derived from anomalous data collections. Threshold is chosen to minimize total errors(both false negative and false positive).

Then the optimal threshold T is calculated by minimizing total errors. However, this is not accurate enough, it implicates an assumption that chances are the same for a new data collection to be either anomalous or not. If we can determine the probability for a new data collection to be anomalous in any segment of data sequence, the equation should be modified as minimizing expected errors, where we use α to denote anomaly probability.

Moreover, with an estimated anomaly probability, SDD-R can be also optimized by ranking all data collections according to their divergence value and select first $n \cdot \alpha$ ones with highest values as anomalies.

4. EVALUATION

Our algorithm was implemented and interpreted in Python 3.5.2. All experiments were tested on Ubuntu 16.04. In the following experiments, we figured out properties of real world

data and performance of our technique against anomalous data collections. We also made a comparison among variations of SDD algorithms and MGoF.¹

4.1. Methodology

We adopted Koubei sellers' transaction records² in experiments. It was provided by Alibaba Tian Chi big data competition where all records were collected from real world business scenarios. Two types of click-farmed data was generated according to patterns described in section 2.2. In our experiments, we use ν to denote the magnitude coefficient of click farming. Hence $|D_{anomalous}| = (1 + \nu)|D_{normal}|$. In the following experiments without extra illustration, we adopted $\nu = 1$.

One defect of this data set is that the detailed time stamp is aligned at each hour of the day due to desensitization. We constructed an enhanced data set by assigning every time stamp a random value for minutes and seconds. Therefore, the enhanced data set should be closer to the reality.

Divergence metric adopted in each SDD algorithms was Jensen-Shannon divergence if no specific notation is made. However, MGoF used only Kullback-Leibler divergence due to its special mechanism. We use a "+" to denote algorithms optimized by a given α .

4.2. Experiments on Koubei Data Set

We first tested our algorithms to see whether and why the algorithm works. Anomalies were random selected days replaced by corresponding click farmed version. To play the role of purchasing platform, we investigated two levels of transaction distribution. The first level is to simply draw a histogram aligned to time spans. The second level is to draw a histogram on the sub-volumes in each time span (i.e. a histogram on frequencies in the first level histogram). To test SDD-E, we randomly selected 30 correct days and 10 click farmed days as normal and anomalous evidence respectively. Here $\alpha = 0.2$. The results are shown in Table 1 and 2.

When classifying toward 1st level histograms, centralized click farming behaviours can be easily discovered. It is because normal collections share a similar distribution while centralized click-farmed ones abruptly violated the original shape. When it came to 2nd level histograms, equalized click farming were also effectively discovered. It can be clearly seen in Fig. 2 that distribution of divergence of both click farming types shows an obvious deviation from the normal one.

The result showed that our technique outperformed MGoF in every real world case. SDD-E provided best performance, yet it consumed the most computing power. Comparison

¹ All resources and more detailed experiment results can be retrieved online: <https://github.com/TramsWang/StatisticalAnomalyDetection>

² <https://tianchi.aliyun.com/competition/information.htm?raceId=231591>

Table 1: Performance on Raw Data

	Centralized						Equalized							
	1st Level			2nd Level			1st Level			2nd Level				
	Pre(%)	Rec(%)	F1(%)	T(ms)	Pre(%)	Rec(%)	Pre(%)	Rec(%)	F1(%)	T(ms)	Pre(%)	Rec(%)	F1(%)	T(ms)
SDD-R	89.51	48.75	63.12	266.77	21.97	99.38	6.67	0.63	1.14	249.15	21.22	72.50	32.83	7.81
SDD-R+	91.25	91.25	91.25	265.96	61.88	61.88	9.38	9.38	9.38	247.31	44.38	44.38	44.38	6.92
SDD-E Static	92.46	68.75	78.86	292.50	36.55	98.13	53.26	6.67	0.63	1.14	271.45	36.07	86.25	50.86
SDD-E Static+	85.02	32.50	47.02	293.77	46.24	91.88	61.52	10.00	0.63	1.18	272.71	43.60	76.25	55.48
SDD-E Dynamic	49.11	99.38	65.73	699.97	23.01	99.38	37.37	245.65	10.36	18.13	681.09	22.09	93.13	35.71
SDD-E Dynamic+	73.21	98.75	84.09	701.06	48.02	96.25	64.07	255.43	8.15	6.88	681.89	40.79	78.13	53.59
MGoF	14.08	21.88	17.13	292.14	13.01	4.38	12.50	3.13	5.00	250.42	12.50	3.13	5.00	3.71

Table 2: Performance on Enhanced Data

	Centralized						Equalized									
	1st Level			2nd Level			1st Level			2nd Level						
	Pre(%)	Rec(%)	F1(%)	T(ms)	Pre(%)	Rec(%)	F1(%)	T(ms)	Pre(%)	Rec(%)	F1(%)	T(ms)				
SDD-R	81.54	41.88	55.33	238.61	17.49	92.50	29.42	13.86	6.67	0.63	1.14	239.44	21.15	71.88	32.69	12.08
SDD-R+	91.25	91.25	91.25	236.94	36.88	36.88	36.88	12.98	8.13	8.13	8.13	236.47	38.75	38.75	38.75	11.03
SDD-E Static	88.31	67.50	76.52	257.60	31.99	92.50	47.54	9.74	6.67	0.63	1.14	257.65	34.49	91.25	50.06	9.37
SDD-E Static+	69.67	13.13	22.09	259.30	42.12	73.75	53.62	9.59	10.00	0.63	1.18	258.45	44.17	82.50	57.54	9.32
SDD-E Dynamic	42.93	99.38	59.96	1106.25	20.12	95.63	33.25	277.70	10.57	17.50	13.18	1108.93	20.30	94.38	33.42	271.74
SDD-E Dynamic+	69.18	98.13	81.15	1110.81	33.25	87.50	48.19	292.60	7.06	3.75	4.90	1118.00	37.21	80.63	50.92	283.50
MGoF	13.97	17.50	15.54	294.08	14.66	8.13	10.45	8.01	12.50	3.13	5.00	250.10	18.42	8.75	11.86	7.55

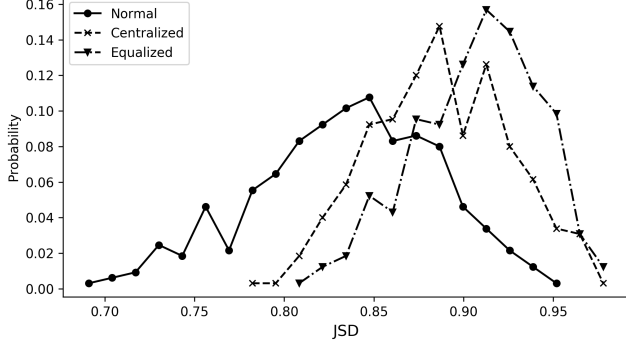


Fig. 2: This figure shows distribution of JSD values (on 2nd level histograms) of normal and two types of click farming data. Divergences were calculated according to a reference averaged among all correct distributions.

among SDD-R revealed improvement of reference as well as the importance of threshold under this technique. It is also clear that dynamic SDD-E is capable of tracing the gradual shift of environment. MGoF turned out to be the worst since it always mark several false positive when c_{th} had not been met and much more false negatives when similar errors occurred too many.

Parameter α improved total accuracy of dynamic SDD-E algorithm by 10-20% as was supposed. It also increased its F1 by more than 20%. α made a great difference in SDD-R as well, which illustrated that divergence sorted almost all collections in correct order according to the averaged reference. However, static SDD-E did not show the same improvement. Since environment drift took greater influence in the result. In comparison with α , adaptive threshold given by evidence sets did not bring the most improvement. But this threshold can be applied together with other optimizations such as slide windows.

4.3. Test against Anomaly Proportion and Magnitude

In this experiment, we tested algorithm performance under various anomaly proportion and magnitude. α ranged from 0.1 to 0.9 when $\nu = 1$ and $\nu \in [0.1, 0.9]$ when $\alpha = 0.1$, other settings remains the same.

Fig. 3 shows that our technique outperformed MGoF and was relatively stable when dealing with all proportions of 1st level centralized anomalies. SDD-E performed even better since it maintains knowledge of both normal and anomalous distributions and calculates the threshold according to the best expectation. However, it relies on the accuracy of distribution estimation. When it came to 2nd level distributions, histograms became much coarser since data available was highly limited and thus its performance suffered dramatically. For the classifiers of MGoF, they compromised to a high error rate. Because more anomalies gathered together and the algo-

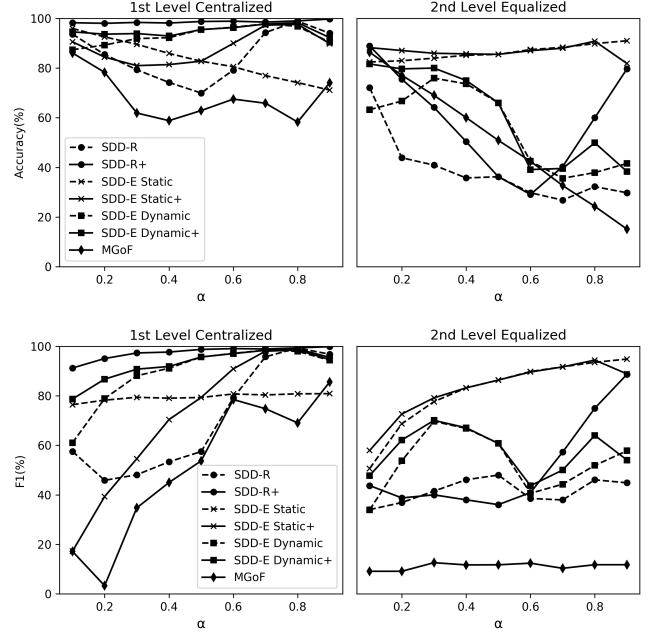


Fig. 3: Accuracy and F1 on Different Anomaly Probabilities

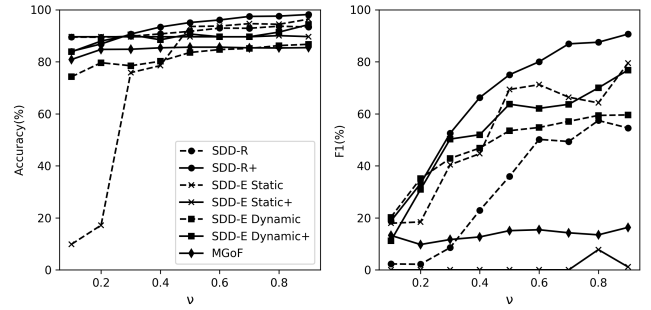


Fig. 4: Accuracy and F1 on Different Anomaly Magnitudes

rithm recognized them as clusters of normal data.

From Fig. 4 we can conclude that our algorithms are still the best, given that they are most sensitive toward tiny anomalous variations. However, static SDD-E did not rise until $\nu > 1$, this is because it suffered from fluctuation on the trade environment at the mean time. MGoF is not sensitive toward minor anomalies either. For a relatively small magnitude of click farming, the classifiers of MGoF quickly degrade to be trivial. The rigid threshold could not automatically rise up and was thus far from optimal.

4.4. Discussion

MGoF's learning procedure of anomalous probability hypothesis is inefficient. To maintain a comprehensive knowledge of anomalies, MGoF has to reserve a single hypothesis entry for every type of them. But in reality, it is always the case that

we face the heterogeneity of outliers. In the Koubei data set, there can be tens of anomalous distributions caused solely by centralized click farming. It takes a long time to discover every possible type of anomaly. Besides, if there happens to be more than c_{th} anomalous distributions of the same type, later discovered collections will no longer declared to be anomalous any more.

However, in SDD-R and SDD-E, that is not a problem since it can map and gather all anomalies together and draw a universal boundary between them and all normal collections. These techniques are suitable to all typical divergence metrics and consume little computation power (except dynamic SDD-E). The only drawback is that they require comprehensive estimation of target distributions. Although other parameters need estimation as well, they are naturally addressable under big data circumstances.

5. CONCLUSION

This paper proposes a series of collective anomaly detection techniques, which helps detect data manipulations in modern data pipelines and data centres. Different from existing algorithms designed for collective anomalies, our approach employs statistical distance as the similarity measurement. We explored several technical points involved in the design of the algorithm and performed a thorough experiment to test its efficiency. The comparison experiment also illustrated the advantages of our technique. It can be concluded that the our technique can efficiently discover anomalies within the data collections and the classifier is sensitive enough toward real world data manipulations.

6. REFERENCES

- [1] Ehsan Amid, Annamaria Mesaros, Kalle J Palomäki, Jorma Laaksonen, and Mikko Kurimo, "Unsupervised feature extraction for multimedia event detection and ranking using audio content," in 2014 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP). IEEE, 2014, pp. 5939–5943.
- [2] Pedro J Moreno, Purdy P Ho, and Nuno Vasconcelos, "A kullback-leibler divergence based kernel for svm classification in multimedia applications," in Advances in neural information processing systems, 2004, pp. 1385–1392.
- [3] Dong-Chul Park, Duc-Hoai Nguyen, Seung-Hwa Beack, and Sancho Park, "Classification of audio signals using gradient-based fuzzy c-means algorithm with divergence measure," in Pacific-Rim Conference on Multimedia. Springer, 2005, pp. 698–708.
- [4] Hang See Pheng, Siti Mariyam Shamsuddin, Wong Yee Leng, and Razana Alwee, "Kullback leibler divergence for image quantitative evaluation," in AIP Conference Proceedings. AIP Publishing, 2016, vol. 1750, p. 020003.
- [5] Jacob Goldberger, Shiri Gordon, and Hayit Greenspan, "An efficient image similarity measure based on approximations of kl-divergence between two gaussian mixtures," in IEEE, 2003, p. 487.
- [6] Maggie Overfelt, "The next big threat in hacking data sabotage,".
- [7] T Caudell and D Newman, "An adaptive resonance architecture to define normality and detect novelties in time series and databases," in IEEE World Congress on Neural Networks, Portland, Oregon, 1993, pp. 166–176.
- [8] Philip K Chan and Matthew V Mahoney, "Modeling multiple time series for anomaly detection," in Data Mining, Fifth IEEE International Conference on. IEEE, 2005, pp. 8–pp.
- [9] Suratna Budalakoti, Ashok N Srivastava, Ram Akella, and Eugene Turkov, "Anomaly detection in large sets of high-dimensional symbol sequences," 2006.
- [10] Nong Ye et al., "A markov chain model of temporal behavior for anomaly detection," in Proceedings of the 2000 IEEE Systems, Man, and Cybernetics Information Assurance and Security Workshop. West Point, NY, 2000, vol. 166, p. 169.
- [11] Christina Warrender, Stephanie Forrest, and Barak Pearlmutter, "Detecting intrusions using system calls: Alternative data models," in Security and Privacy, 1999. Proceedings of the 1999 IEEE Symposium on. IEEE, 1999, pp. 133–145.
- [12] Dmitry Pavlov, "Sequence modeling with mixtures of conditional maximum entropy distributions," in Data Mining, 2003. ICDM 2003. Third IEEE International Conference on. IEEE, 2003, pp. 251–258.
- [13] Rose Yu, Xinran He, and Yan Liu, "Glad: group anomaly detection in social media analysis," ACM Transactions on Knowledge Discovery from Data (TKDD), vol. 10, no. 2, pp. 18, 2015.
- [14] Chengwei Wang, Krishnamurthy Viswanathan, Lakshminarayan Choudur, Vanish Talwar, Wade Satterfield, and Karsten Schwan, "Statistical techniques for online anomaly detection in data centers," in Integrated Network Management (IM), 2011 IFIP/IEEE International Symposium on. IEEE, 2011, pp. 385–392.