

ユーザ行動に基づくモバイルネットワーク故障検知と予測

Mobile Network Failure Detection and Forecasting with Multiple User Behavior

大木 基至*¹ 竹内 孝*² 植松 幸生*¹ 上田 修功*²
 Motoyuki Oki Koh Takeuchi Yukio Uematsu Naonori Ueda

*¹NTT コミュニケーションズ株式会社 *²NTT コミュニケーション科学基礎研究所
 NTT Communications Corporation NTT Communication Science Laboratories

Providing stable and high-quality service is a critical issue for mobile network service providers. However, due to an unexpectedly huge amount of data traffic exceeding network capacity of a provider, a mobile network service experiences severe failures such as network troubles, performance deterioration, and slow throughput. Then, the service users often detect service outages before the service provider detects them. They can immediately publish their impressions on the service through social media and search for failure information on the web. In this paper, we propose a machine learning approach that incorporates multiple user behavior data into detecting and forecasting failure events. The approach is based on novel feature extraction methods and a model ensemble method that combines outputs of supervised and unsupervised learning models from multiple user behavior datasets. We demonstrate the effectiveness of the approach by extensive experiments with real-world failure events.

1. はじめに

モバイルネットワークのサービス提供者にとって、ネットワーク故障の早期検知あるいは予測は、高品質サービスの安定提供を実現するための重要な課題となっている。これはトラフィックの急増 [Cisco 17] によるスループット低下が引き起こすサービス品質の劣化や、通信設備の経年劣化等による深刻なネットワーク故障を予め知ることができれば、事前の対策措置の実行や迅速な復旧が可能となるためである。

サービス提供者では、ネットワーク故障を発見するためトラフィックやシステムログの監視が常時行われている [Gill 11]。ただ、このような監視は膨大なデータの確認が必要なため、故障発見までに遅延が生じる場合がある。また、コールセンターでの受付情報を用いた故障発見も行われているが、通話対応に時間がかかるため、常時監視と同様に遅延する可能性が高い。

一方、モバイルネットワークのサービス利用者（以下、ユーザ）は、サービスの使用不能や品質劣化についての SNS へ投稿しており、Twitter での故障への反応はコールセンターへの連絡数よりも早いと報告されている [Qiu 10]。Tweet データからネットワーク故障を検知するための研究やシステムの提案が行われているが、Tweet データのみを利用しており実用可能な性能は得られていない [Takeshita 15, Maru 16]。

ユーザの故障への反応は Twitter のみならず、Web 上のアクセス数や検索ログ数に反映されていると考えられる。そのため、これらの複数種類データの同時解析が実現できれば、より精度の高い早期の故障検知の実現が期待される。複数のユーザ行動データを用いた解析は、実社会上的イベント検知（暴動など） [Kallus 14, Zhao 16] で行われているが、現在までネットワーク故障検知の分野では行われていない。

そこで本研究では、複数のユーザ行動データから特徴抽出を行い、アンサンブル教師あり機械学習モデルを用いて故障検知と予測を行うための方法を提案する。本研究の全体像を図 1 に示す。ユーザ行動データにはソーシャルデータ（Twitter, ニュース記事, RSS）やテレコムデータ（Web アクセスログ,

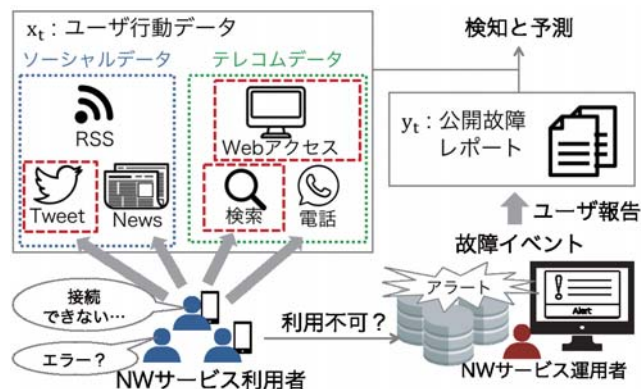


図 1: ユーザ行動に基づく故障検知と予測の全体像。本研究では赤の点線で囲まれたデータを使用する。

検索ログ, コールセンター受付情報) などの様々な候補が挙げられるが、本研究ではサービスに関する Tweet データ, サービスの公式ページへの Web アクセスログ, サービス名を含む検索クエリ数を用いて、現実のネットワーク故障の検知・予測を行う。本研究で使用するデータの例として、あるサービス故障発生前後 4 日間に観測された 1 時間ごとの総頻度を図 2 に示す。赤枠は故障発生期間を表す。故障発生時に頻繁に観測される典型的な単語やユーザにアクセスされる URL を示す。故障発生時にはそれぞれのユーザ行動データで通常時とは異なる振る舞いや変化が確認できる。

提案法では、複数のユーザ行動データを入力データとし、故障発生の予測スコアを毎分ごとに出力する。入力データは毎分単位で集計されるため、非常にスパースとなり、さらに故障の発生数は非常に少数であるため、故障に対応するラベルが非常に少ない不均衡となる。またユーザ行動データによって、スコアの予測に適した機械学習モデルは異なると考えられる。そこで提案法では、データのスパースネスとラベル不均衡の問題に適した特徴量変換を行い、複数のモデルの出力を組み合わせることで予測スコアを出力するモデルアンサンブル法を行う。現実

連絡先: 大木基至, NTT コミュニケーションズ株式会社, 東京都港区芝浦 3-4-1 グランパークタワー 16F, m.ooki@ntt.com

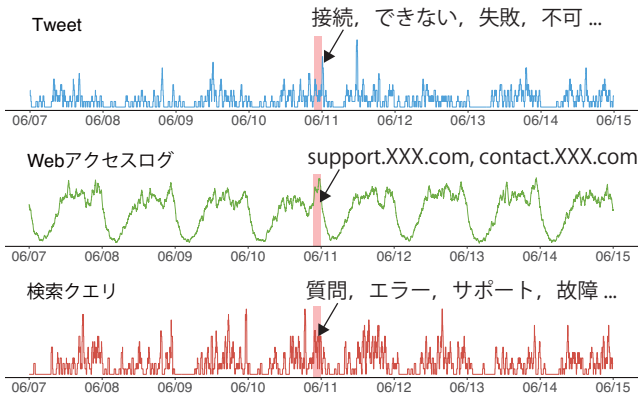


図 2: あるモバイルネットワークサービスに関する Tweet 数, 公式ページへ Web アクセス数, 検索クエリ数. 赤枠は故障の発生期間を表す.

表 1: 九つの公開された故障情報

故障 ID	発生期間 (分)	日付と時刻
1	188	2016/11/13, 08:22 a.m.
2	60	2016/2/04, 06:00 a.m.
3	60	2016/2/13, 10:52 p.m.
4	60	2016/3/24, 02:00 a.m.
5	150	2016/6/10, 09:20 p.m.
6	512	2016/6/15, 09:20 p.m.
7	100	2016/7/06, 08:50 a.m.
8	950	2016/7/27, 10:00 p.m.
9	1554	2016/12/25, 01:46 a.m.

のモバイルネットワークサービスの故障データを用いた故障検知および予測の実験を行い, 既存の単一のデータを用いた場合と比較し, 提案法による大幅な精度向上を示す.

2. データセット

本章では, 本論文で使用する故障情報と三つのユーザ行動データについて説明する.

2.1 故障情報

インターネット上で公開されているモバイルネットワークサービスで発生した 2015 年 11 月から 2016 年 12 月の九つの故障情報を収集した. 表 1 に収集した故障 ID, 故障の発生期間 (分), 日付と時刻を示す. 発生期間が 60 から 1554 分まで変化するの, サービスの故障はそれぞれ異なる原因で発生し, 復旧に要する時間も異なるためである.

2.2 ユーザ行動データ

本研究では, Tweet と Web アクセスログと検索クエリの三つのデータをユーザ行動データとして使用する. 故障データと同期間でそれらを収集した. Tweet データは, Twitter API を通じて, サービス名の略称を含む合計 72,070 ツイートを収集した. 各 Tweet に対して, 形態素解析を行い, 名詞, 動詞, 形容詞の単語集合を作成した. そのうち, 頻度が 2 以下の単語やストップワードを削除し, 合計 5750 単語を使用した. Web アクセスログデータとは, モバイルネットワークサービスに関するトップページ, 購入ページ, サポートページを含む合計 29,520,772 件の Web ページへのアクセスログである. アクセスログは (アクセス時刻 = “2016-03-20 11:54:12”, URL

表 2: 訓練とテストデータの評価データセット

テスト ID	訓練 ID	ラベル比率 [%]	スパース率 [%]		
			\mathcal{T}	\mathcal{W}	\mathcal{Q}
5	1 - 4	0.12	99.55	99.06	99.85
6	2 - 5	0.10	99.56	99.06	99.85
7	2 - 6	0.27	99.57	99.42	99.85
8	2 - 7	0.31	99.58	99.42	99.85
9	4 - 8	0.57	99.58	99.58	99.84

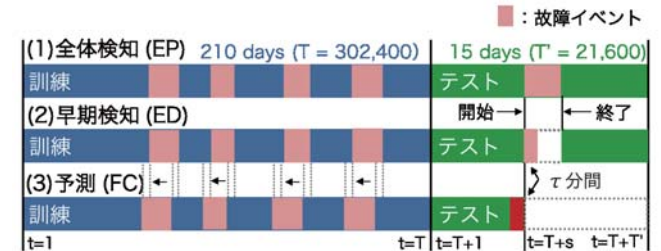


図 3: 全期間 (EP), 早期検知 (ED), 将来予測 (FC) の三つの問題設定

= “http://contact.XXX.com”, ページタイトル = “Contact Form”) の組で表現される. その中から, アクセス数が 2 以下のユニークページを削除し, 合計 6,378 件のページを使用した. 検索クエリデータとは, 検索エンジンの運用ログから, モバイルネットワークサービス名の略称を含む合計 82,136 件の検索データである. 検索数が 2 以下のユニークな検索クエリを削除し, 合計 1209 個のユニークな検索クエリを収集した.

3. 問題設定

ある時刻 t のユーザ行動データの特徴ベクトルを $\mathbf{x}_t \in \mathbb{R}^M$, 各時刻 t の故障の発生有無を表すラベルを $y_t \in \{-1, 1\}$ で表し, 訓練データとテストデータを $\{(\mathbf{x}_t, y_t)\}_{t=1}^T$ と $\{\mathbf{x}_t\}_{t=T+1}^{T+T'}$ とする. テスト期間として, ある一つの故障発生時の前後 1 週間を使用する. 訓練データはテストデータより過去に発生した故障 ID を複数含む 210 日間のデータとする. 表 1 にある故障情報から, 後半の故障 5 件をそれぞれ 1 件ずつテストデータに含む評価データセットを作成した. それらを表 2 に示す. 訓練データに含まれる故障 ID を訓練 ID, 訓練データ内の故障発生 (つまり, $y_t = 1$) 比率をラベル比率, 訓練データ内の特徴ベクトルの 0 の比率をスパース率とする. ラベル比率が非常に小さいのは, 故障発生頻度は正常時 (つまり, $y_t = -1$) と比べ, 非常に少ないためである.

本研究では, 三つの問題設定を定義する. それらを図 3 に示す. それらは, テストデータの全範囲の故障を検知する問題 (全体検知: EP), テストデータの故障開始時刻から τ 分間の故障を検知する問題 (早期検知: ED), テストデータの τ 分後の故障発生を予測する問題 (予測: FC) の三つである. 本研究では, 各問題設定において訓練データの特徴ベクトル \mathbf{x}_t から予測スコアを求める機械学習モデルを訓練し, テストデータを用いて学習モデルを評価する.

4. 故障検知と予測

本章では, 故障検知と予測を行うための機械学習モデルと比較し, 複数のユーザ行動データと機械学習モデルを組み合わ

せるモデルアンサンブル法を提案する。故障検知と予測性能の評価実験を通じて、提案法の有用性を検証する。

4.1 検知性能のモデル比較

はじめに、各ユーザ行動データの時系列性を用いてスパース比率を改善するために、各特徴ごとに単純移動平均による特徴変換を行った [Botezatu 16]。次に、不均衡データの分類性能の向上が確認されている Bi-Normal Separation (BNS) [Forman 03] による特徴変換を行った。特徴変換されたデータを用いて、最適な機械学習モデルを選択するために、四つの教師あり分類モデル: ロジスティック回帰 (LR), 決定木のアダプスト (ADA), ランダムフォレスト (RF), 三層パーセプトロン (NN) と二つの教師なし異常検知モデル: 1 クラス SVM (OCS), オートエンコーダー (AE) の計六つのモデルを比較する。実装は scikit-learn と Keras を使用した。各モデルのパラメータは、5 交差検証法を用いて、訓練データを 5 分割し、検証データでの AUC の平均値が最も高いパラメータとした。

三つのユーザ行動データと六つの機械学習モデルを用いて、五つの故障 ID に対する早期検知 (ED) と全体検知 (EP) の AUC の結果を表 3 に示す。便宜上、Tweet データを T , Web アクセスログデータを W , 検索クエリデータを Q と表記する。“All” は五つの故障 ID の平均値と標準偏差を表す。ED では Tweet データを用いたオートエンコーダー, EP では Web アクセスログデータを用いた LR が最も高い性能を示した。また、Tweet データでは EP より ED のほうが平均的に性能が高く、Web アクセスログデータでは、ED より EP のほうが平均的に性能が高いことを確認した。この結果から、ユーザは最初に Twitter から故障に関する情報を発信し、その後、故障情報を手に入れるため、Web ページへのアクセスを行う一連の行動がある可能性が示唆される。また、故障 ID : 5 では、Web アクセスログデータを用いた RF, 故障 ID : 6 の ED では、検索クエリデータを用いた RF が最も高い性能を示した。この結果から、各故障 ID で故障の種類や発生時間が異なるため、機械学習モデルとユーザ行動データの最適な組み合わせはそれぞれの故障で異なることが考えられる。

4.2 モデルアンサンブル法

複数のユーザ行動データと機械学習モデルを活用するために、2 段階のモデルアンサンブル法を行う。前節の実験から一つ以上の故障で最も性能が高かったモデル LR, RF, AE を用いる。モデルアンサンブル法は、以下で定義される。

$$(\text{Level 1}) \quad \mathbf{z}_t = (f_{M_1}^{D_1}(\mathbf{x}_t), f_{M_2}^{D_2}(\mathbf{x}_t), \dots, f_{M_3}^{D_q}(\mathbf{x}_t)) \quad (1)$$

$$(\text{Level 2}) \quad \hat{p}_t = \text{sig}(\mathbf{w}^T \mathbf{z}_t + b) \quad (2)$$

sig はシグモイド関数 $\text{sig}(a) = \frac{1}{1 + \exp(-a)}$, $\mathbf{w} \in \mathbb{R}^{3 \times q}$ は 2 段階目のモデル係数, b はバイアス項である。 $f_{M_j}^{D_j}(\mathbf{x}_t)$ は、時刻 t のユーザ行動データ $D_j, j = (1, \dots, q)$ の三つのモデル M_1 (=LR) と M_2 (=RF) と M_3 (=AE) の出力結果である。この定義から、最終予測スコア \hat{p}_t は各ユーザ行動データと機械学習モデルの出力結果の線形結合で算出される。準ニュートン法を用いてモデル係数 \mathbf{w} とバイアス項 b を推定するための損失関数を以下で定義する。

$$L(\mathbf{w}, b) = - \sum_{t=1}^T \{(1 - y_t) \log(1 - p_t) + y_t \log p_t\} + \psi(\mathbf{w}, b, C) \quad (3)$$

$\psi(\mathbf{w}, b, C)$ は正則化項を表す。実験では、リッジ正則化 $\psi(\mathbf{w}, b, C) = \frac{C}{2} (\|\mathbf{w}\|^2 + b^2)$ を用いる。

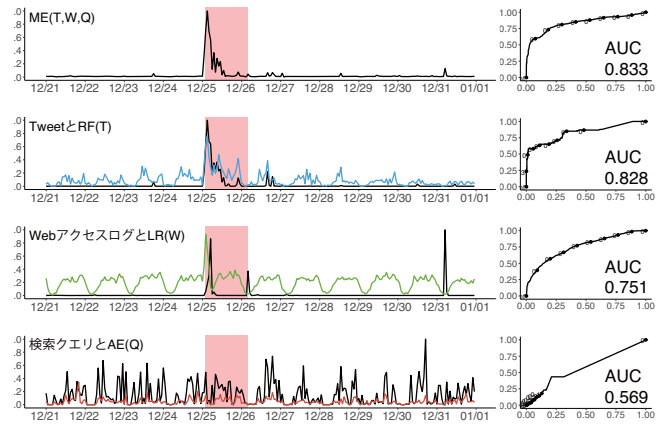


図 4: 故障 ID : 9 のユーザ行動データの推移と各モデルの予測スコア

比較手法として、三つのユーザ行動データを連結させて一つの特徴ベクトルからリッジ正則化付きのロジスティック回帰を学習させるデータアンサンブル法 (DE) を用いる。EP と ED に関して、すべての故障 ID の平均値と標準偏差の結果を表 4 に示す。ED では Tweet と Web アクセスログのモデルアンサンブル法, EP では三つのユーザ行動データのモデルアンサンブル法が最も高い性能を示した。また、単一のユーザ行動データを用いたモデルアンサンブル法より、複数のユーザ行動データを用いたモデルアンサンブル法が高い結果を示した。一方、データアンサンブル法では性能改善が確認されなかった。これは、特徴次元数が増大し、過学習を引き起こした可能性が考えられる。

故障 ID : 9 の予測スコアと AUC の値を図 4 に示す。1 行目はモデルアンサンブル法, 2 から 4 行目は各ユーザ行動データの推移と最も性能が高かった機械学習モデルの結果を表す。赤枠は故障発生期間を表す。モデルアンサンブル法の結果から、故障発生前後の偽陽性 (2016/12/27 の Tweet や 2016/12/31 の Web アクセスログなど) を減らす効果を確認した。また、モデルアンサンブル法は Tweet データと Web アクセスログデータの故障時の予測スコアを活用した効果で、急増でかつロングテールな予測スコアを表現し、AUC が向上したと考えられる。

4.3 予測性能のモデル比較

本節では、三つ目の問題設定である τ 分後の故障予測に関する性能を比較する。故障ラベル y_t を τ 分先にずらした訓練データ $\{(\mathbf{x}_t, y_{t+\tau})\}_{t=1}^{T-\tau}$ を使用して機械学習モデルの訓練を行う。テストデータは図 3 の (3) のとおり $\{\mathbf{x}_t\}_{t=T-\tau+1}^{T+s}$ を使用する。機械学習モデルとして、Tweet データと Web アクセスログデータを用いたモデルアンサンブル法 (ME(T,W), ME(T), ME(W)), データアンサンブル法 DE(T,W), LR(T), LR(W) の六つを比較する。予測のための時間幅は $\tau \in \{10, 30, 50\}$ とする。

各時間幅 τ ごとの AUC の平均値と標準偏差を表 5 に示す。太字は各時間ごとで最も高いスコアを表す。 $\tau = 10$ と 30 では ME(T,W), $\tau = 50$ で ME(T) が最も高い性能を示した。これらの高い性能から、故障発生前からユーザがサービスの品質劣化を体感していた可能性が考えられる。複数, または一つのユーザ行動データを用いたモデルアンサンブル法は予測問題においても効果的であることが確認できた。

表 3: 各ユーザ行動データと機械学習モデルの AUC の結果. ED は $\tau = 60$ の早期検知の結果を表す.

データ	モデル	ID:5		ID:6		ID:7		ID:8		ID:9		All	
		ED	EP	ED	EP	ED	EP	ED	EP	ED	EP	ED	EP
\mathcal{T}	LR	0.92	0.94	0.82	0.60	0.62	0.74	0.99	0.76	1.00	0.82	0.87 \pm 0.16	0.77 \pm 0.12
	ADA	0.78	0.63	0.02	0.35	0.42	0.45	0.49	0.58	1.00	0.51	0.54 \pm 0.37	0.50 \pm 0.11
	RF	0.83	0.79	0.72	0.76	0.95	0.96	0.99	0.83	1.00	0.83	0.90 \pm 0.12	0.83 \pm 0.08
	NN	0.79	0.65	0.18	0.34	0.29	0.27	0.86	0.58	0.99	0.47	0.62 \pm 0.36	0.46 \pm 0.16
	OCS	0.88	0.69	0.72	0.56	0.50	0.50	0.50	0.56	1.00	0.70	0.72 \pm 0.22	0.60 \pm 0.09
	AE	0.94	0.82	0.86	0.73	0.93	0.90	0.90	0.70	1.00	0.73	0.93 \pm 0.05	0.78 \pm 0.08
\mathcal{W}	LR	0.78	0.90	0.76	0.86	0.99	0.99	1.00	0.99	0.98	0.75	0.90 \pm 0.12	0.90 \pm 0.10
	ADA	0.62	0.80	0.86	0.70	0.82	0.89	0.82	0.93	0.94	0.80	0.81 \pm 0.12	0.82 \pm 0.09
	RF	0.94	0.97	0.82	0.75	0.68	0.73	0.93	0.95	1.00	0.71	0.87 \pm 0.13	0.82 \pm 0.13
	NN	0.57	0.82	0.51	0.71	0.88	0.89	0.98	0.97	0.99	0.75	0.79 \pm 0.23	0.83 \pm 0.11
	OCS	0.60	0.50	0.40	0.68	0.79	0.78	0.04	0.51	0.82	0.52	0.53 \pm 0.32	0.60 \pm 0.13
	AE	0.41	0.50	0.38	0.69	0.81	0.81	0.03	0.50	0.82	0.47	0.49 \pm 0.33	0.59 \pm 0.15
\mathcal{Q}	LR	0.53	0.39	0.86	0.58	0.78	0.74	0.45	0.61	0.60	0.55	0.65 \pm 0.17	0.57 \pm 0.13
	ADA	0.26	0.33	0.16	0.56	0.61	0.62	0.45	0.53	0.27	0.50	0.35 \pm 0.18	0.51 \pm 0.11
	RF	0.73	0.55	0.87	0.70	0.62	0.52	0.36	0.64	0.42	0.48	0.60 \pm 0.21	0.58 \pm 0.09
	NN	0.25	0.32	0.20	0.51	0.30	0.47	0.54	0.52	0.50	0.52	0.36 \pm 0.15	0.47 \pm 0.09
	OCS	0.77	0.68	0.70	0.56	0.66	0.60	0.47	0.57	0.88	0.58	0.70 \pm 0.15	0.60 \pm 0.05
	AE	0.78	0.67	0.85	0.62	0.50	0.62	0.37	0.66	0.71	0.57	0.64 \pm 0.20	0.63 \pm 0.04

表 4: モデルアンサンブル (ME) とデータアンサンブル (DE) の AUC の結果. ED は $\tau = 60$ の早期検知の結果を表す.

データ	ME		DE	
	ED	EP	ED	EP
\mathcal{T}	0.94 \pm 0.09	0.87 \pm 0.08	-	-
\mathcal{W}	0.90 \pm 0.13	0.88 \pm 0.11	-	-
\mathcal{Q}	0.62 \pm 0.16	0.60 \pm 0.08	-	-
(\mathcal{T}, \mathcal{W})	0.95 \pm 0.10	0.90 \pm 0.09	0.88 \pm 0.15	0.88 \pm 0.09
(\mathcal{T}, \mathcal{Q})	0.94 \pm 0.08	0.87 \pm 0.08	0.88 \pm 0.13	0.73 \pm 0.13
(\mathcal{W}, \mathcal{Q})	0.90 \pm 0.13	0.88 \pm 0.11	0.90 \pm 0.13	0.89 \pm 0.10
($\mathcal{T}, \mathcal{W}, \mathcal{Q}$)	0.94 \pm 0.09	0.91 \pm 0.09	0.88 \pm 0.14	0.88 \pm 0.10

表 5: τ 分後の予測性能 (FC) の比較

モデル	$\tau = 10$	$\tau = 30$	$\tau = 50$
ME(\mathcal{T}, \mathcal{W})	0.81 \pm 0.27	0.78 \pm 0.21	0.69 \pm 0.18
ME(\mathcal{T})	0.72 \pm 0.35	0.71 \pm 0.34	0.72 \pm 0.20
ME(\mathcal{W})	0.76 \pm 0.25	0.77 \pm 0.17	0.77 \pm 0.09
DE(\mathcal{T}, \mathcal{W})	0.72 \pm 0.35	0.66 \pm 0.26	0.69 \pm 0.08
LR(\mathcal{T})	0.79 \pm 0.27	0.60 \pm 0.29	0.48 \pm 0.23
LR(\mathcal{W})	0.70 \pm 0.25	0.71 \pm 0.25	0.73 \pm 0.20

5. おわりに

三つのユーザ行動データからモバイルネットワークサービスの故障を検知および予測する機械学習アプローチを提案した。提案法により、過去の故障ラベルとユーザ行動データから特徴抽出を行い、モデルアンサンブル法に基づいて故障検知および予測を行った。モバイルネットワークサービスの実世界の複数のデータを用いた複数の実験を通じて、提案法は全体検知、早期検知、予測の三つの問題において高い性能を示した。

今後は、さらなる他のユーザ行動データの活用、高い検知性能を達成した要因の分析、本手法に基づいたリアルタイム検知

システムの開発とその運用が考えられる。

参考文献

- [Botezatu 16] Botezatu, M.; Giurgiu, I.; Bogojeska, J.; and Wiesmann, D.: Predicting disk replacement towards reliable data centers, In *SIGKDD* (2016).
- [Cisco 17] Cisco Systems Inc.: Cisco visual networking Index: Global mobile data traffic forecast update 2016 - 2021 (2017).
- [Forman 03] Forman, G.: An extensive empirical study of feature selection metrics for text classification, *Journal of Machine Learning Research*, 3, 1289–1305 (2003).
- [Gill 11] Gill, P.; Jain, N.; and Nagappan, N.: Understanding network failures in data centers: measurement, analysis, and implications, In *SIGCOM* (2011).
- [Kallus 14] Kallus, N.: Predicting crowd behavior with big public data, In *WWW* (2014).
- [Maru 16] Maru, C.; Enoki, M.; Nakao, A.; Yamamoto, S.; Yamaguchi, S.; and Oguchi, M.: Development of failure detection system for network control using collective intelligence of social networking service in large-scale disasters, In *HT* (2016).
- [Qiu 10] Qiu, T.; Feng, J.; Ge, Z.; Wang, J.; Xu, J.; and Yates, J.: Listen to me if you can : Tracking user experience of mobile network social media, In *IMC* (2010).
- [Takeshita 15] Takeshita, K.; Yokota, M.; and Nishimatsu, K.: Early network failure detection system by analyzing twitter data, In *IM* (2015).
- [Zhao 16] Zhao, L.; Ye, J.; Chen, F.; Lu, C. T.; and Ramakrishnan, N.: Hierarchical incomplete multi-source feature learning for spatiotemporal event forecasting, In *SIGKDD* (2016).