# 一人称画像と位置に基づくライフログセグメンテーション

瀧本 祥章 山本 修平 西村 拓哉 戸田 浩之

† 日本電信電話株式会社 NTT サービスエボリューション研究所 〒 239-0847 神奈川県横須賀市光の丘 1-1 E-mail: † {yoshiaki.takimoto.ar,shuhei.yamamoto.ea,takuya.nishimura.fk}@hco.ntt.co.jp, hirotoda@acm.org

あらまし スマートフォンやスマートグラスなどのデバイスの普及により,ユーザの状態や,位置情報,周辺環境を 記録したライフログと呼ばれる系列データを得られるようになった.ライフログには多岐にわたる情報が含まれてお り,ユーザの行動内容理解など,活用しようとする研究が盛んに行われている.しかし,ライフログをセグメントと 呼ばれるインデックス可能な単位に分割する標準的なアプローチは未だ存在しない.ここでセグメントとは,料理や 買い物をしているなど,それ単独で意味を持ち,利活用のための検索における基本的な単位となるものである.そこ で本稿では,ユーザの位置情報と,ユーザの視点から継続的に撮影された画像に注目し,ライフログをセグメントに 分割する手法を提案する.評価実験では,NTCIR-13 Lifelog-2 タスクで提供されるユーザ2人の延べ90日分のライ フログと,人手による分割結果を利用して提案手法の検証を行い,提案手法により高精度にセグメントへの分割がで きることを示した.

キーワード ライフログ, DBSCAN, 滞留点抽出, Gated CNN, ImageNet, Places, 移動軌跡

1. はじめに

スマートフォンやスマートグラスなどのデバイスの普及に より,ユーザの状態や,位置情報,周辺環境に関するライフロ グと呼ばれる系列データを得られるようになった.例えば,ス マートデバイスの一種である Fitbit [1] からは心拍数や,歩数, 睡眠の質と時間などのユーザの状態を表すデータから構成され るライフログが得られる.また,スマートフォンアプリである Moves [2] からは walking, running, cycling などのユーザの 移動状況や,ユーザの移動軌跡,滞在地から構成されるライフ ログが得られる.

このように得られるライフログは多様な情報を持つため, ユーザの行動内容理解や,ヘルスケアなどへ活用しようとする 研究が盛んに行われている [14].しかし,ライフログのデータ 量はその性質上,時間と共に増大し,大量のデータから人手で 有用な情報を特定することは困難である.そのため,ライフロ グをセグメントと呼ばれる索引付け可能な単位に分割すること (ライフログセグメンテーション)が求められている.様々な手 法 [8–10,19,20,28,29]が提案されているが,セグメンテーショ ンされるセグメントの性質や数,長さに制約がある場合や,セ グメンテーションを行う上で重要な情報である位置情報を考慮 していないなどの課題がある.ここでセグメントとは,料理や 買い物をしているなど,それ単独で意味を持ち,利活用のため の検索における基本的な単位となるものである.

そこで本稿では,図1のように得られるユーザの位置情報 と,ユーザの視点から継続的に撮影された画像(一人称画像) に注目し,ライフログセグメンテーションを行い,図1中のセ グメント1-6のようなセグメントを抽出する手法を提案する. セグメントはユーザが同一行動を継続している時間や,環境を 表す.例えばセグメント1では,ユーザは会話をしており,セ グメント2ではコンピュータを操作している.これらのセグメ



ントを抽出するため、本稿ではまず、画像類似度アプローチを ベースラインとして提案し、技術的な課題を明らかにする.画 像類似度アプローチでは系列データ中の一人称画像を前後で比 較し、類似度が閾値以下である場合に画像間を分割点として抽 出する.この手法はシンプルであるが、ノイズの影響が大きい、 図1のセグメント4やセグメント5など移動を伴うセグメン トを抽出できないという2つの課題がある.この2つの課題を 解決するため、滞留抽出アプローチでは滞留点抽出技術である D-star [22]を、2群検定アプローチではWelchの2群検定を

活用する.また,Gated CNN アプローチでは前後複数枚の画像を入力として用いることにより,これらの課題を回避する.

評価実験では,提案手法を評価するため,NTCIR-13 Lifelog-2 タスク [13] で提供されるユーザ 2 人の延べ 90 日分のライフ ログと,人手による分割結果を利用して各手法の比較検証を 行った.その結果,ベースラインである画像類似度アプローチ と比較して,他の3つの手法が高い精度を示すことを確認した. 本稿の主な貢献は以下の通りである.

ベースラインとなる画像類似度アプローチの提案とライ
 フログセグメンテーションにおける課題の明確化

• NTCIR-13 Lifelog-2 タスクの LEST (Lifelog Event

表 1: 関連研究と提案手法との相違点

文献	利用データ	提案手法との相違点
Ellis ら [9,10]	音声	場所,環境のみを考慮
Wang ら [29]	画像	セグメントの単位時間が 5 分
Lin ら [20]	映像	場所のみを考慮
Doherty ら [8]	画像	長時間の移動を考慮できない
Li ら [19]	画像	対象セグメントが限られる
Talavera ら [28]	画像	位置情報が未考慮
Luら [21]	映像	位置情報が未考慮
Castro ら [5]	画像	多量の教師データが必要
Poleg ら [23]	画像	対象セグメントが限られる

Segmentation Task) サブタスク [13] の優勝手法を含む, ライ フログセグメンテーションの課題を解決する相異なる 3 つの手 法の提案

• 複数のユーザおよび長期間日常的に収集したデータを用 いた評価実験による提案手法の定性的評価および定量的評価, および有効性の確認

本稿の構成は以下の通りである.まず,2章で関連研究とし て,ライフログセグメンテーションに関する研究を紹介する. 3章で本稿における基本的な概念の定義や問題設定ついて述べ る.その後,4章でベースライン手法および提案手法の詳細を 論じ,5章で提案手法を評価するために行った実験について述 べる.最後に,6章で本稿のまとめと今後の課題について述 べる.

# 2. 関連研究

ライフログは近年注目され,関連する研究が多く行われている [14]. その活用のため,セグメントと呼ばれる索引付け可能 な単位に分割する技術であるライフログセグメンテーション が求められており,表1のように様々な手法が提案されてい る [5,8–10,19–21,23,28,29]. 各手法と提案手法との相違点を 順に述べる.

Ellis らは 62 時間分の周囲の音声データをライフログとして 収集し、スペクトル情報を活用して street や restaurant などの 16 個の場所や環境を表すセグメントに分割した [9,10].本稿で は収集する情報が音声データではなく, 位置情報と一人称画像 であるという点や,想定するセグメントとして場所や環境だけ でなく,家事や食事などユーザの行動に基づくセグメントも想 定する点で異なる. Wang らは 6 週間の一人称画像をライフロ グとして収集し, walking outside や meeting などの 6 種類の セグメントに5分単位で分割した [29]. しかし, ライフログに 含まれるセグメントは5分間単位とは限らない.これに対し提 案手法では,このような単位時間を設定しない. Lin らは時間 制約付きのクラスタリングにより、映像から構成されるライフ ログを office など場所を表すセグメントに分割した [20].本稿 では、前述したように家事や食事などユーザの行動を基づくセ グメントも想定する点で異なる. Doherty らは Sense-Cam [16] で収集した1日当たり1785枚の画像からなるデータに対し、 前後間の画像について, MPEG-7のメタデータから得られる色

表 2: 本稿で用いる記号一覧

記号 意味

L ライフログ

- $l_i$  観測データ
- *n* ライフログのデータ数
- $S_i$  観測データ $l_i$ の画像特徴量ベクトル
- $s_{i,j}$  画像特徴量ベクトル  $S_i$ の次元 jの値
- $m{T_i}$  観測データ  $l_i$ の潜在トピック分布

やエッジ情報の類似度を計算し、閾値未満の箇所をセグメント の境界としている [8]. ただし,ユーザが移動中の場合には過剰 に境界を検出してしまうことから、後処理によって互いに近い セグメント境界のうち,最初の境界のみを残す処理を行う.し かし、後処理は予め定めた時間間隔([8]では3分)のみに依存 することから、より長時間の移動を伴うセグメントを正しく検 出できない. これに対し提案手法では, 長時間の移動を伴うセ グメントであっても検出可能である. Li らは Sense-Cam によ り収集した画像系列を時系列データとし、固有値のピークを導 出し、セグメントの境界とする手法を提案した [19]. しかしこ の手法は、全てのセグメントを検出し、分割することを想定し ておらず、ノイズへの対応もしていない. これに対し提案手法 では,全てのセグメントを検出し,分割することを想定してお り、ノイズへの対応もしている. Talavera らは ImageNet [24] で学習した CNN によりグラフカット技術を用いてセグメント の境界を検出する [28]. また, Lu らは映像に映っている物体 に注目し、一人称映像から重要な瞬間を抽出し、一人称映像の 要約を行う [21]. しかしこれらの手法では、画像情報のみを考 慮しており、ライフログセグメンテーションにおいて重要な情 報である位置情報を考慮していない.これに対し,提案手法の 滞留抽出アプローチや Gated CNN アプローチでは位置情報を 考慮する. Castro らは ImageNet で学習した CNN を再学習す ることによりセグメントにおけるユーザの行動内容を含めて推 定を行う [5]. しかし, 想定する行動ごとに教師データが必要 となるため、新たな行動を考慮するために数千枚のデータにラ ベルを付与する必要がある.これに対し,提案手法の滞留抽出 アプローチや2群検定アプローチは教師データを必要とせず, Gated CNN アプローチにおいても行動別の教師データは必要 としない. Poleg らはユーザの頭部の動きに注目し、画像の変 位を用いてユーザの行動を認識する [23]. しかし, 認識できる 行動は sitting や walking などのユーザの移動状態に限られる. これに対し、提案手法は前述のように家事や食事などの行動を 想定する点で異なる.

# 3. 準備

本章では、本稿における基本的な概念の定義や問題設定を行う.本稿で用いる記号は表2の通りである.

**3.1** 基本的な概念

本稿では,位置情報と一人称画像などの系列データから構成 されるライフログをセグメントに分割する.

まず,入力となるライフログを以下のように定義する.



[定義 1] ライフログ  $L = [l_1, ..., l_n]$  はユーザの状態や位置情報,周囲の情報の観測データを時系列順に並べた系列である. ここで,ライフログに含まれる各データ  $l_i \in L$  は時刻情報  $l_i.time$ ,位置情報  $l_i.loc$ ,一人称画像  $l_i.img$  など複数の属性を 持つ.位置情報  $l_i.loc$  は位置を表す二次元座標であり,一人称 画像  $l_i.img$  は自動的に撮影されたユーザ視点の生画像である. [定義 2] セグメントは開始時刻と終了時刻の組であり,ユー ザや周囲の状態について特定の状態が継続した期間を表す. ここで,特定の状態には料理や買い物,散歩などが例として挙 げられる.

## 3.2 問題設定

本稿では,任意期間のライフログ $L = [l_1, ..., l_n]$ を入力し て受け取り,セグメントの分割点を抽出し,セグメント集合を 出力する.ただし,各データ $l_i$ はいずれかのセグメントに属す るものとする.

# 4. 提案手法

本章では、本稿で扱う4つの手法について順に述べる.まず、 ベースラインとなる画像類似度アプローチの内容と課題につ いて述べる.その後、課題について別々のアプローチで解決を 図った3つの提案手法について述べる.

#### 4.1 画像類似度アプローチ

ライフログセグメンテーションの素朴な手法として図 2 に 示すような前後の観測データ間を比較する手法が考えられる. この手法では,観測データ中の画像のみを用いてセグメンテー ションを行う.図 2 ではまず,周囲の状態を把握するため,観 測した一人称画像 $l_i.img$ を GoogLeNet [26] などの,画像のク ラス分類が可能なニューラルネットワークに入力する.そして, 出力された各クラスの分類確率 $S_i = (s_{i,1}, \ldots, s_{i,m})$ を入力画 像の特徴ベクトルとする.その後,前後の観測データ間の画像 特徴量 $S_i, S_{i+1}$ をコサイン類似度を用いて比較し,類似度が 閾値未満であるデータ間を分割点として検出する.これにより, セマンティックギャップ[6]として知られる画像表現と意味内容 の乖離を回避し,観測データの意味内容に基づく比較が可能に なると期待できる.

この手法はシンプルである一方で、ノイズに弱い、位置情報



を考慮できないという2つの課題を持つ.例えば,図3のよう なライフログを得られたときを考える.図3(a)は1つのセ グメントとして抽出すべき部分である.しかし, la.img の撮影 時にカメラを物体が覆ってしまっていることにより, l2.img と *l*<sub>4</sub>.*imq* の間の類似度が大きい一方で, *l*<sub>2</sub>.*imq* と *l*<sub>3</sub>.*imq*, *l*<sub>3</sub>.*imq* と l4.img の間の類似度が小さくなってしまい, 3 つのセグメン トに分割されてしまう. このようにカメラを何らかの物体が覆 う、ユーザが体の向きを変えるなどが原因となるノイズにより、 セグメントが過剰に分割されてしまう.図3(b)に移動を伴う 単一セグメント中のライフログの例を示す. このライフログは ユーザがショッピングをしており、1つのセグメントとして抽 出すべき部分である.しかし、ユーザが移動しているため、撮 影される画像が刻一刻と変化しており、必ずしも類似度が高く なるとは限らない. そのため, 閾値によっては複数個のセグメ ントに分割されてしまう. このように、移動を伴うセグメント では、同一のセグメント中でも画像が変化していくため、正し い分割ができない.

#### 4.2 滞留抽出アプローチ

画像類似度アプローチにおける2つの課題を解決するため, 滞留抽出アプローチでは画像だけでなく,位置情報にも焦点を 当て,分割点を抽出する.本稿では滞留の抽出に D-star [22] を拡張し,用いる.ここで,滞留とはユーザが一定時間以上, 一定範囲内に留まることを表す.また,D-star は移動軌跡から ユーザが滞留を行った地点である滞留点を抽出する技術であり, ノイズの影響を受けにくい,ストリーム処理に対応しているな どの特徴を持つ.

滞留抽出アプローチの概要を図4に示す.まず,滞留抽出ア プローチでは画像類似度アプローチと同様に学習済みネット ワークに一人称画像*l<sub>i</sub>.img*を入力し,出力された各クラスの分 類確率*S<sub>i</sub>*を得る.次に,D-starを画像の類似度を考慮するよ うに拡張した simD-star により,ライフログから移動を伴わな い(滞留している)セグメントをクラスタとして抽出する.な お,simD-starの詳細は後述する.その後,クラスタの前後を 分割点として,移動部分のセグメントおよび滞留部分のセグメ ントを抽出する.

simD-star のアルゴリズムをアルゴリズム 1 に示す.ここで, アルゴリズム中の関数  $d(l_i.loc, l_j.loc)$  は  $l_i.loc > l_j.loc$ の距離



を表し、 $s(N(l_i))$  は近傍データ集合  $N(l_i)$  に含まれるデータの 最早観測時刻と最遅観測時刻の差,即ち観測期間を表す.また、  $sim(S_i, S_j)$  は画像特徴量間の類似度を表し、画像の意味的な 類似度を考慮するため、[27] を参考に以下のように定義する.

$$\operatorname{sim}(\boldsymbol{S}_{i}, \boldsymbol{S}_{j}) = \frac{\sum_{k=1}^{m} \min(s_{i,k}, s_{j,k}) \times \operatorname{idf}_{k}}{\sum_{k=1}^{m} \max(s_{i,k}, s_{j,k}) \times \operatorname{idf}_{k}}$$
(1)

なお,  $idf_k$  は逆文書頻度を表し,

$$\mathtt{idf}_k = \log \frac{n}{\sum_{p=1}^n s_{p,k}} \tag{2}$$

と定義され、各次元の重要度を考慮することを可能とする. simD-star は入力として、ライフログデータ L と、ウィンドウ サイズ q, 距離閾値  $\varepsilon$ , 最短観測期間閾値  $m_{time}$  (DBSCAN [11] の密度閾値 MinPts に相当),最短滞留時間閾値 mstay を受け 取り、ライフログの滞留を伴うセグメントを抽出する.simDstar では時系列順にライフログのデータ *li* の処理を行う. ま ず、*l*<sub>i</sub>の近傍に存在し、かつ、画像特徴量が類似するデータ*l*<sub>i</sub> をウィンドウWから抽出し(4-7行目),各々の近傍データ集 合  $N(l_i), N(l_j)$  に加える (8–9 行目). その後,  $l_{i-q+1}$  につい て、近傍データ集合の観測期間が最短観測期間閾値以上である 場合(11行目)は,近傍データ集合 N(l<sub>i-a+1</sub>) が既存のクラス タと同一のデータを保持する場合には、同一データを保持する クラスタおよび近傍データ集合  $N(l_{i-q+1})$  をすべてマージする (12-13 行目). そうでない,即ち,近傍データ集合 N(li-q+1) が既存のクラスタと同一のデータを保持しない場合には、新た なクラスタを形成する(14-15行目).最後に、観測期間が滞 留時間閾値以上であるクラスタを全て、滞留を伴うセグメント として出力する.

#### 4.3 2 群検定アプローチ

2 群検定アプローチでは,画像類似度アプローチにおける 2 つの課題を解決するため,画像の特徴次元を削減し,対象デー タデータ前後での画像の分布変化に注目する.

2 群検定アプローチの概要を図 5 に示す.まず,画像類似度 アプローチ,滞留抽出アプローチと同様に学習済みネットワー クに一人称画像 *l<sub>i</sub>.img* を入力し,出力された各クラスの分類 確率 *S<sub>i</sub>* を得る.次に,画像の特徴次元を削減し,ノイズの影 響を小さくするため,各画像の画像特徴量 *S<sub>i</sub>* を文書,各特徴 量のうち,確率値 *s<sub>i,j</sub>* が閾値を超える次元をその画像が含む

7	アルゴリズム 1: simD-star
	<b>Input:</b> $L, q, \varepsilon, m_{time}, m_{stay}, \tau$
1	$W \leftarrow \phi$ // スライディングウィンドウ
2	$C \leftarrow \phi$ // クラスタ集合
3	foreach $l_i \in L$ do
4	Push $l_i$ in $W$
5	$N(l_i) \leftarrow \phi$ // 近傍データ集合
6	foreach $l_j \in W$ do
7	$  \text{if } d(l_i.loc,l_j.loc) < \varepsilon \wedge \mathtt{sim}(\boldsymbol{S_i},\boldsymbol{S_j}) > \tau \text{ then} \\$
8	Add $l_i$ to $\boldsymbol{N}(l_j)$
9	Add $l_j$ to $\boldsymbol{N}(l_i)$
10	Shift $l_{i-q+1}$ from W
11	if $s(N(l_{i-q+1})) \ge m_{time}$ then
12	$ ext{ if } \exists oldsymbol{C} \in \mathcal{C}.oldsymbol{N}(l_{i-q+1}) \cap oldsymbol{C}  eq \phi  ext{ then}$
13	同一データを持つクラスタと $m{N}(l_{i-q+1})$ を全て
	L マージ
14	else
15	Add $N(l_{i-q+1})$ to $C$

16 return  $\{C|s(C) \ge m_{stay}\}$ 



単語とみなし、全てのデータに対して LDA(Latent Dirichlet Allocation)[4] を用いて潜在トピック分布  $T_i = (t_{i,1}, ..., t_{i,K})$ を抽出した。その後、データ  $l_i$  に対して、前後 q 個ずつの データを対象とする 2 つのスライディングウィンドウ window1 , window2 を用意し、各ウィンドウ中の潜在トピック分布集合  $\{T_{i-q}, ..., T_i\}, \{T_i, ..., T_{i+q}\}$  についてトピック次元ごとに 平均  $E_i(t_1), ..., E_i(t_K)$  と分散  $V_i(t_1), ..., V_i(t_K)$  を算出する。そ の後、Welch の t 検定により、p 値  $P_i = (p_{i,1}, ..., p_{i,K})$  を算 出する。最後に、次元ごとに求めた p 値の総和を計算し、値が 大きい予め定めた定数個のデータをセグメント終了データとし て抽出する。

## 4.4 Gated CNN アプローチ

Gated CNN アプローチは近年系列データを扱う上で注目さ れているネットワークである Gated CNN [7] に,対象時刻の データとその前後のデータを入力することにより,対象時刻の データがセグメント終了のデータであるか推定する手法である. なお入力するデータは,画像や位置情報に縛られずに,様々な 情報を扱うことができ,入力を変化させてもネットワーク構成



図 6: Gated CNN アプローチで用いたネットワーク構造

の変更を必要としない.図6にニューラルネットワークの構成 を示す.

まず前処理部では、画像や位置情報などの入力データから複 数のデータを生成する.入力データの例として,画像類似度ア プローチや2群検定アプローチと同様の、画像特徴量 $S_i$ や潜 在トピック分布 T<sub>i</sub>や,前後の画像特徴量および潜在トピック 分布を比較した結果である  $\cos(S_{i-1}, S_i)$ ,  $\cos(T_{i-1}, T_i)$ , 緯 度経度情報に DBSCAN を適用したものが挙げられる.次に, 前処理部で生成されたデータは2層の Full Connect 層によっ てベクトル Vi に変換する. その後, ライフログの流れを考慮 するため、前後 q 個のベクトル  $V_{i-q}, \ldots, V_{i+q}$  と共に Gated Unit を複数回通される. Gated Unit は,入力から2つの同次 元のベクトルを出力し、一方に Sigmoid 関数を適用してから 要素積を取ることにより、必要な情報のみを出力する.本研究 では、精度向上のため Sigmoid 関数を適用しないベクトルに 対し, Layer Normalization [3] を適用している. 最後に, Full Connect 層と Softmax 層によってデータ  $l_i$  がセグメント終了 のデータであるか判定する.

# 5. 評価実験

本章では,提案手法の有用性を確認するため,NTCIR-13 Lifelog-2 タスク [13] で提供されたユーザ2人の延べ90日分の ライフログと,人手による分割結果を用いた評価実験について 述べる.まず,用いたデータについて説明する.次に,評価指 標について述べる.最後に行った実験の結果について述べる.

5.1 データセット

本稿で用いるデータセットは NTCIR-13Lifelog-2 タスクで 提供されるユーザ2人(ユーザ1,ユーザ2)のデータである. データの収集期間はユーザ1が2016年8月8日から10月5日 の59日間であり,ユーザ2が2016年9月9日から10月11日 の9月20日と10月9日を除く31日間である.ユーザが起き ている間,1分当たり1つのユーザ視点の画像および位置情報, Activity (Moves アプリで収集された walking や cycling など ユーザの移動状況を表すラベル)を含むライフログデータを収 集しており,1日当たり1,250個から1,500個のデータが存在 する. なお, プライバシー保護の観点から, 顔やデバイスの画 面が画像に映り込んでいる場合にはぼがしが入れられ, かつ, 全ての画像が 1024 × 768 の解像度にリサイズされている. ま た, 位置情報も自宅や職場については, GPS による絶対位置で はなく, それぞれ Home, Work と意味的な位置情報に Moves アプリケーション [2] によって置換されている.

前述のように画像特徴量が複数の手法で必要となるため,深層 学習技術により抽出した.抽出に用いたモデルは ImageNet [24] で学習した GoogLeNet [26] および AlexNet [18] (1000 次元), Places365 [30] で学習した GoogLeNet, AlexNet, VGG [25] および ResNet [15] である (365 次元). これらは *Caffe* [17] 上で動作する学習済みモデルが github<sup>(注1)(注2)</sup>上で公開されて いる.

次に,2 群検定アプローチや Gated CNN アプローチで用い た潜在トピックを LDA によって抽出した.抽出方法は 4.3 節 の通りであり,単語とみなす確率値の閾値は 0.1 と設定した. また予備実験の結果から,最も高い性能を示したことからト ピック数を 10 に設定した.

## 5.2 評価指標

本稿では,NTCIR-13 Lifelog-2 タスクの LEST サブタスク による評価指標に基づいて,Precision, Recall, F1 score を用 いて評価を行った.各値の算出方法は以下の通りである.

$$Precision = \frac{\sum_{i=1}^{n} \sum_{j=1}^{n} And \left( f(l_i, l_j), GT(l_i, l_j) \right)}{\sum_{i=1}^{n} \sum_{j=1}^{n} f(l_i, l_j)} \quad (3)$$

$$\text{Recall} = \frac{\sum_{i=1}^{n} \sum_{j=1}^{n} And \left( f(l_i, l_j), GT(l_i, l_j) \right)}{\sum_{i=1}^{n} \sum_{j=1}^{n} GT(l_i, l_j)} \quad (4)$$

F1 score = 
$$\frac{2 \times \text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$
 (5)

ここで  $f(l_i, l_j)$  および  $GT(l_i, l_j)$  はデータ  $l_i \ge l_j$  が提案手法に よる分割または正解データにおいて,同一セグメントに属して いる場合に1 (*True*),属していない場合に0 (*False*)の二値 をとる関数である.また,And(x, y) は論理積を表し,*True*, 即ち,x = y = 1のときに1となり,*False*のときに0となる.

なお、Precision は各データの組が同一セグメントに含まれ ると判定した結果のうち、実際に同一セグメントに含まれる組 の割合を表し、細かく分割すると、より大きな Precision を得 られる可能性が高い.一方、Recall は同一セグメントに含まれ るデータ組のうち、同一セグメントに含まれると判定できた組 の割合を表し、分割数が小さくなると、より大きな Recall を得 られる可能性が高い.また、F1 score は Precision と Recall の 調和平均である.

#### 5.3 分割結果

4章 で提案した各手法の比較を行うため,各ユーザのライ フログデータに提案手法を適用した.適用した手法と用いた入 力の一覧は表3の通りであり,パラメータなどは予備実験を元 に表4のように設定した.なお,「2群検定 + 滞留抽出」は滞 留抽出による分割と2群検定による分割の両方を用いて分割し

<sup>(</sup>注1):https://github.com/BVLC/caffe/tree/master/models

<sup>(2):</sup> https://github.com/CSAILVision/places365



た結果である.また,Gated CNN アプローチは入力する特徴 量の追加が容易であることから様々な入力を追加したモデルを 用意した.例えば,「トピック間類似度」は直前データとの潜在 トピック分布のコサイン類似度,「DBSCAN」は緯度経度情報 をDBSCANを適用した結果のクラスタ ID,「Activity」は入力 データに付加されていたデータである.ここで,モデルのユー ザ独立性は、2人のユーザに同一のモデル(パラメータや学習 データが共通)を用いたことを表す.具体的には、チェックが ある場合,ユーザ2人に共通の1つのモデルを用意し、Gated CNN アプローチにおける学習も2人分まとめて行った.チェッ クがない場合には、ユーザ2人別々のモデルを準備し、Gated CNN アプローチにおける学習も個々で行った.

各手法の結果は表5および図7の通りである.これらの結果 から、画像類似度アプローチの Precision が最も大きくなった 一方で、Recall の値が最も小さく、F1 score の値も最小となっ たことがわかる.これは、ノイズや移動を伴うセグメントの影 響により過剰な分割を行ってしまったことが原因と考えられる.

その一方で、他のアプローチでは位置情報やデータの流れを 考慮することにより、同一セグメントであるデータの判定がで きるようになったことから, Recall が大きくなっている. 特に, 滞留抽出アプローチの F1 score が最大であること,2 群検定ア プローチよりも滞留抽出アプローチを組み合わせた方が高精度 であることから、 位置情報の考慮がセグメント分割の精度向上 に大きく寄与することが示唆される. 例えば, ユーザ2の9月 13 日について滞留抽出アプローチの F1 score が他の手法を大 きく上回っている.この日,ユーザ2は家事やショッピングな ど移動を伴うセグメントを主に行っており、これらのセグメン ト中に含まれる画像の類似度は低くなっている. そのため, 画 像に基づく手法では過剰に分割を行ってしまい, Recall や F1 score が小さくなる傾向にある.その一方で,滞留抽出アプロー チでは、これらを移動を伴うセグメントを移動として捉えるこ とにより,過剰な分割を回避することにより比較的正確な分割 が行え, Recall や F1 socre が大きくなる. また, ユーザ1の 8月27日についても滞留抽出アプローチのF1 score が他の手 法を大きく上回っている.この日,ユーザ1は多くの公共交通 機関や自動車による移動や長時間のコンピュータの利用を行っ ていた.これらのセグメント中には周囲の状況変化やユーザの 体勢,コンピュータの画面の変化などが原因により,画像の類 似度が低くなってしまう.そのため,ユーザ2の9月13日と 同様に,滞留抽出アプローチ以外の手法では過剰な分割を行っ てしまい,Recall やF1 score が小さくなる傾向にある.これ に対し,滞留抽出アプローチでは移動の考慮や,逆文書頻度を 考慮した類似度によりこれらの影響を緩和できているため,比 較的正確な分割が行え,Recall やF1 socre が大きくなる.

Gated CNN アプローチ内では,Gated CNN5 が最も高精 度になり,滞留抽出アプローチ,「2 群検定+滞留抽出」に次ぐ 精度となった.これは,Gated CNN によって,複数の特徴量 から有用な特徴を捉えることが可能であることを示す一方で, Gated CNN6 で入力に追加した LDA のトピックの次元数が学 習データに対して大きく,過学習してしまうことを示している.

# 6. ま と め

本稿では、位置情報と一人称画像に注目し、ライフログをセ グメントと呼ばれる索引付け可能な単位に分割するライフログ セグメンテーション技術を4つ提案した.まず、ベースライン となる一人称画像の類似度に基づく画像類似度アプローチを提 案した.その後、画像類似度アプローチにより明らかになった 技術的課題を解決するために、異なる3つのアプローチを提案 した.提案手法の有効性を確認するために、実際に収集したラ イフログを用いた評価実験では、提案手法により、高精度にセ グメントへの分割ができることを示した.また、各アプローチ の比較から位置情報の考慮が重要であることが示唆された.

以下に今後の課題について述べる.本研究ではユーザごとに, ー律のパラメータで分割を行った.しかし,各アプローチの日 ごとの精度が大きく異なっていることから,ライフログは同一 ユーザであっても,日によってその性質は大きく異なることが 推測される.この日に依存する性質の変化を考慮した適応的な セグメンテーションについては今後の課題としたい.また,他 の課題として,より多様なユーザに対する評価実験や,分割結 果を活用して,買い物や食事などのユーザの行動内容を推定す る技術である人間の行動認識技術 [12] などと組み合わせること が挙げられる.

#### 献

[1] Fitbit. https://www.fitbit.com/home.

文

- [2] Moves. https://moves-app.com/.
- [3] J. L. Ba, J. R. Kiros, and G. E. Hinton. Layer normalization. arXiv preprint arXiv:1607.06450, 2016.
- [4] D. M. Blei, A. Y. Ng, and M. I. Jordan. Latent dirichlet allocation. J. Mach. Learn. Res., 3:993–1022, Mar. 2003.
- [5] D. Castro, S. Hickson, V. Bettadapura, E. Thomaz, G. Abowd, H. Christensen, and I. Essa. Predicting daily activities from egocentric images using deep learning. In proceedings of the 2015 ACM International symposium on Wearable Computers, pages 75–82. ACM, 2015.
- [6] R. Datta, D. Joshi, J. Li, and J. Z. Wang. Image retrieval: Ideas, influences, and trends of the new age. ACM Com-

エデル	入力					エデルのコーザ独立姓	
())//	画像類似度	トピック	トピック間類似度	緯度経度	DBSCAN	Activity	
画像類似度	<ul> <li>✓</li> </ul>						$\checkmark$
滞留抽出	<ul> <li>✓</li> </ul>			$\checkmark$			
2 群検定		$\checkmark$					
2 群検定 + 滞留抽出	<ul> <li>✓</li> </ul>	$\checkmark$		$\checkmark$			
Gated CNN1	<ul> <li>✓</li> </ul>						$\checkmark$
Gated CNN2	<ul> <li>✓</li> </ul>					$\checkmark$	$\checkmark$
Gated CNN3	<ul> <li>✓</li> </ul>			$\checkmark$			$\checkmark$
Gated CNN4	<ul> <li>✓</li> </ul>		$\checkmark$	$\checkmark$		$\checkmark$	$\checkmark$
Gated CNN5	<ul> <li>✓</li> </ul>		$\checkmark$	$\checkmark$		$\checkmark$	
Gated CNN6	<ul> <li>✓</li> </ul>	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	

表 4: 設定パラメータ

画像類似度	画像特徴量:GoogLeNet(ImageNet), 閾値:0.04
滞留抽出	画像特徴量 : GoogLeNet(Places365), ウィンドウサ
	イズq:5,距離閾値ε:40m(ユーザ2は120m),最
	短観測期間閾値 $m_{time}$ :3 min, 滞留時間閾値 $m_{stay}$ :
	$5 \min$ ,類似度閾値 $\tau: 0.4$ (ユーザ 2 は $0.3$ )
2 群検定	セグメント終了データ数:50,ウィンドウサイズ q:
	10
Gated CNN	適応学習アルゴリズム:Adam ( $\alpha = 0.00001, \beta_1 =$
	$0.9, \beta_2 = 0.999, \epsilon = 10^{-8})$ , ミニバッチサイズ:10,
	エポック数・200 学習データ・人手による分割結果
	8日分(正解データとは異なる),ウィンドウサイズ

表 5: 各手法の分割精度

モデル	Precision	Recall	$F1 \ score$
画像類似度	0.901	0.352	0.494
滞留抽出	0.559	0.698	0.579
2 群検定	0.768	0.453	0.550
2 群検定 + 滞留抽出	0.762	0.485	0.573
Gated CNN1	0.848	0.421	0.547
Gated CNN2	0.837	0.421	0.545
Gated CNN3	0.855	0.406	0.535
Gated CNN4	0.860	0.407	0.539
Gated CNN5	0.846	0.436	0.561
Gated CNN6	0.790	0.455	0.551

puting Surveys (Csur), 40(2):5, 2008.

- [7] Y. N. Dauphin, A. Fan, M. Auli, and D. Grangier. Language modeling with gated convolutional networks. arXiv preprint arXiv:1612.08083, 2016.
- [8] A. R. Doherty and A. F. Smeaton. Automatically segmenting lifelog data into events. In *Image Analysis for Multimedia Interactive Services, 2008. WIAMIS'08. Ninth International Workshop on*, pages 20–23. IEEE, 2008.
- [9] D. P. Ellis and K. Lee. Minimal-impact audio-based personal archives. In Proceedings of the the 1st ACM workshop on Continuous archival and retrieval of personal experiences, pages 39–47. ACM, 2004.
- [10] D. P. Ellis and K. Lee. Accessing minimal-impact personal audio archives. *IEEE MultiMedia*, 13(4):30–38, 2006.
- [11] M. Ester, H.-P. Kriegel, J. Sander, X. Xu, et al. A densitybased algorithm for discovering clusters in large spatial databases with noise. In *KDD*, pages 226–231, 1996.
- [12] Y. Guan and T. Plötz. Ensembles of deep lstm learners for activity recognition using wearables. Proc. ACM Interact. Mob. Wearable Ubiquitous Technol., 1(2):11:1–11:28, June

2017.

- [13] C. Gurrin, H. Joho, F. Hopfgartner, L. Zhou, D. T. D. Nguyen, R. Gupta, and R. Albatal. Overview of the NTCIR-13 lifelog-2 task. In *The NTCIR-13 Conference*, Tokyo, Japan, 2017.
- [14] C. Gurrin, A. F. Smeaton, A. R. Doherty, et al. Lifelogging: Personal big data. Foundations and Trends® in Information Retrieval, 8(1):1–125, 2014.
- [15] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In 2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2016, Las Vegas, NV, USA, June 27-30, 2016, pages 770–778, 2016.
- [16] S. Hodges, L. Williams, E. Berry, S. Izadi, J. Srinivasan, A. Butler, G. Smyth, N. Kapur, and K. Wood. Sensecam: A retrospective memory aid. *UbiComp 2006: Ubiquitous Computing*, pages 177–193, 2006.
- [17] Y. Jia, E. Shelhamer, J. Donahue, S. Karayev, J. Long, R. Girshick, S. Guadarrama, and T. Darrell. Caffe: Convolutional architecture for fast feature embedding. In *Proceed*ings of the 22nd ACM International Conference on Multimedia, MM '14, pages 675–678, New York, NY, USA, 2014. ACM.
- [18] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. In Proceedings of the 25th International Conference on Neural Information Processing Systems, NIPS'12, pages 1097– 1105. Curran Associates Inc., USA, 2012.
- [19] N. Li, M. Crane, and H. J. Ruskin. Automatically detecting" significant events" on sensecam. International Journal of Wavelets, Multiresolution and Information Processing, 11(06):1350050, 2013.
- [20] W.-H. Lin and A. Hauptmann. Structuring continuous video recordings of everyday life using time-constrained clustering. SPIE, 2006.
- [21] Z. Lu and K. Grauman. Story-driven summarization for egocentric video. In Computer Vision and Pattern Recognition (CVPR), 2013 IEEE Conference on, pages 2714–2721. IEEE, 2013.
- [22] K. Nishida, H. Toda, and Y. Koike. Extracting arbitraryshaped stay regions from geospatial trajectories with outliers and missing points. In ACM SIGSPATIAL International Workshop on Computational Transportation Science (IWCTS), pages 1–6, 2015.
- [23] Y. Poleg, C. Arora, and S. Peleg. Temporal segmentation of egocentric videos. In *Computer Vision and Pattern Recognition (CVPR)*, 2014 IEEE Conference on, pages 2537– 2544. IEEE, 2014.
- [24] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, A. C. Berg, and L. Fei-Fei. ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision (IJCV)*, 115(3):211–252, 2015.
- [25] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. CoRR, abs/1409.1556, 2014.

- [26] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich. Going deeper with convolutions. In 2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pages 1–9, June 2015.
- [27] Y. Takimoto, K. Sugiura, and Y. Ishikawa. Extraction of frequent patterns based on users' interests from semantic trajectories with photographs. In *Proceedings of the 21st International Database Engineering & Applications Symposium*, pages 219–227. ACM, 2017.
- [28] E. Talavera, M. Dimiccoli, M. Bolanos, M. Aghaei, and P. Radeva. R-clustering for egocentric video segmentation. In *Iberian Conference on Pattern Recognition and Image Analysis*, pages 327–336. Springer, 2015.
- [29] Z. Wang, M. D. Hoffman, P. R. Cook, and K. Li. Vferret: content-based similarity search tool for continuous archived video. In *Proceedings of the 3rd ACM workshop on Con*tinuous archival and retrival of personal experences, pages 19–26. ACM, 2006.
- [30] B. Zhou, A. Lapedriza, A. Khosla, A. Oliva, and A. Torralba. Places: A 10 million image database for scene recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2017.