

Practical Time-Series Clustering for Messy Data in R

Jonathan Page
University of Hawaii Economic Research Organization
uhero.hawaii.edu

16 February, 2018

Problem Definition

GoFundMe-like company in Kenya needs to understand the typical activity patterns in campaign contributions.

Creating a list of campaign archetypes will allow them to reason about the effects of changes to their platform and campaign-specific interventions.

Outline

1. Messy data -> Matrix of trajectories
2. Dynamic Time Warping (DTW) clusters
3. k-Shape (or Shape-based) clusters

Messy Data

Data Cleaning Plan

1. Select tables and columns necessary for analysis
2. Remove test campaigns and users
3. Produce long table of features
4. Create trajectory matrix for each feature

Raw Data

Transactions

campaign_id	contributor_id	amount	payment_time
1	1	64	2016-01-01 00:00:01
1	2	128	2016-01-01 06:00:02
1	3	256	2016-01-02 12:34:56
1	3	512	2016-01-03 06:54:32
1	2	1024	2016-01-05 07:53:10
2	2	2048	2016-01-07 23:59:59

Campaigns

campaign_id	start_time
1	2016-01-01 00:00:00
2	2016-01-05 12:00:00

Intermediate Data

Combined and aggregated (in-time)

campaign_id	balance	amount	contributors	day_of_campaign
1	192	192	2	1
1	448	256	1	2
1	960	512	1	3
1	1984	1024	1	5
2	2048	2048	1	3

Matrix of trajectories

Balance trajectories

b_{it}

192	448	960	960	1984
0	0	2048	2048	2048

Amount trajectories

a_{it}

192	256	512	0	1024
0	0	2048	0	0

Data Cleaning

```
library(tidyverse)
transactions <- read_csv("transactions.csv")
end_of_time <- max(transactions$payment_date)
campaigns <- read_csv(
  "campaigns.csv",
  na = c("", "0000-00-00 00:00:00")
) %>%
mutate(
  days_old = as.numeric(difftime(
    end_of_time,
    date_created,
    units = "days"
  ))
)
```

Data Cleaning

```
transactions <- transactions %>%  
  inner_join(  
    campaigns %>%  
      select(campaign_id, date_created)  
  ) %>%  
  mutate(  
    day_of_campaign = as.numeric(difftime(  
      payment_date,  
      date_created,  
      units = "days"  
    ))  
  ) %>%  
  select(-date_created)
```

Data Cleaning

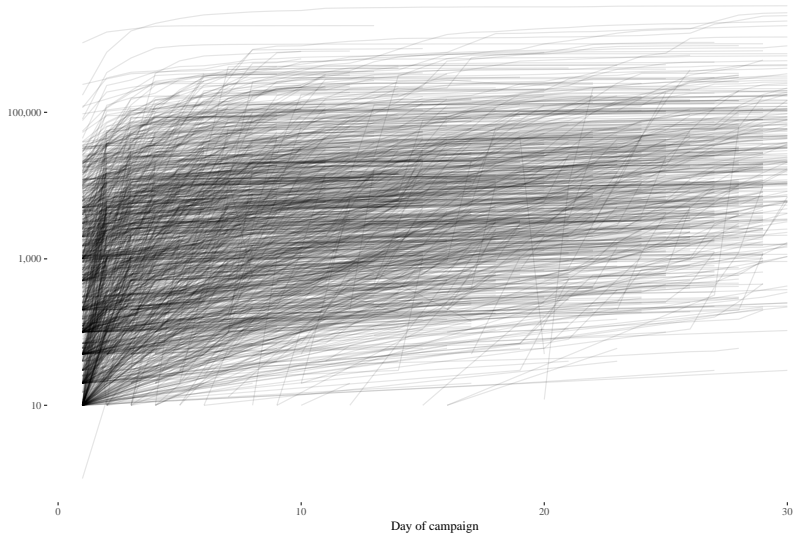
```
clean_transactions <- campaigns %>%  
  filter(days_old >= 30) %>%  
  select(campaign_id) %>%  
  left_join(transactions) %>%  
  group_by(campaign_id) %>%  
  summarize(  
    transaction_count = n(),  
    contributors = length(unique(contributor))  
  ) %>%  
  filter(contributors > 1) %>%  
  select(campaign_id) %>%  
  left_join(transactions) %>%  
  group_by(campaign_id) %>%  
  arrange(payment_date) %>%  
  mutate(campaign_balance = cumsum(amount))
```

Daily Series

```
daily_series <- clean_transactions %>%  
  mutate(day_of_campaign = floor(day_of_campaign) + 1) %>%  
  filter(day_of_campaign <= 30) %>%  
  group_by(campaign_id, day_of_campaign) %>%  
  arrange(id) %>%  
  summarise(  
    balance = last(campaign_balance),  
    amount = sum(amount),  
    transactions = n(),  
  ) %>%  
  filter(day_of_campaign > 0)
```

Long-form Balance Trajectories

Campaign Balance (log-scale)



From sparse series to full matrix

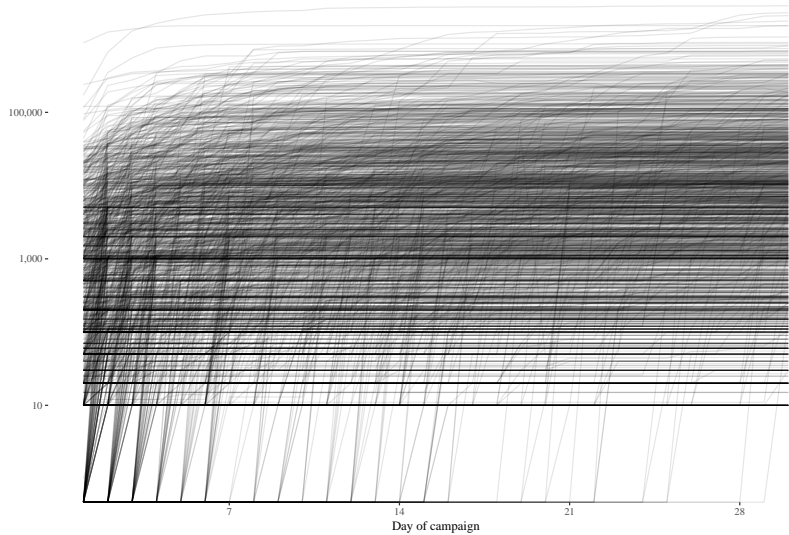
```
balance_traj <- daily_series %>%  
  filter(day_of_campaign %in% 1:30) %>%  
  select(campaign_id, day_of_campaign, balance) %>%  
  spread(day_of_campaign, balance) %>%  
  mutate(`1` = coalesce(`1`, 0)) %>%  
  remove_rownames() %>%  
  column_to_rownames("campaign_id") %>%  
  apply(1, FUN=zoo::na.locf) %>%  
  t()
```

From sparse series to full matrix

```
amount_traj <- daily_series %>%  
  filter(day_of_campaign %in% 1:30) %>%  
  select(campaign_id, day_of_campaign, amount) %>%  
  spread(day_of_campaign, amount, fill = 0) %>%  
  remove_rownames() %>%  
  column_to_rownames("campaign_id")
```

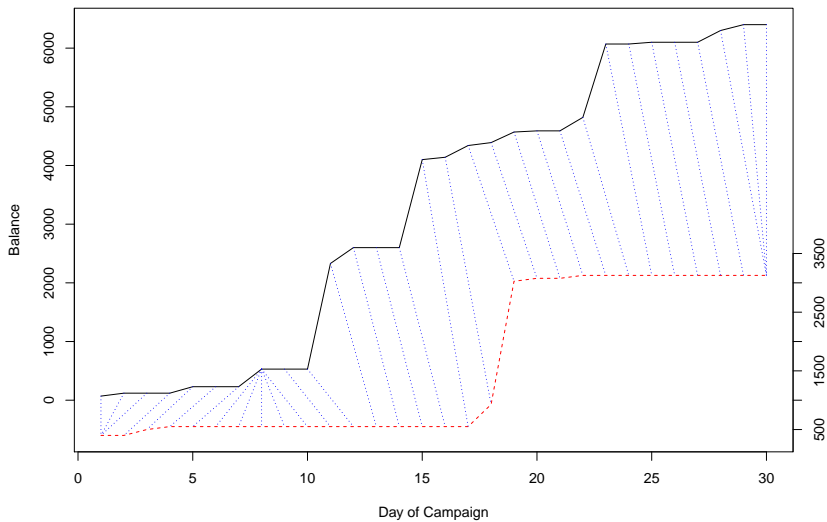
Balance Trajectories

Campaign Balance (log-scale)



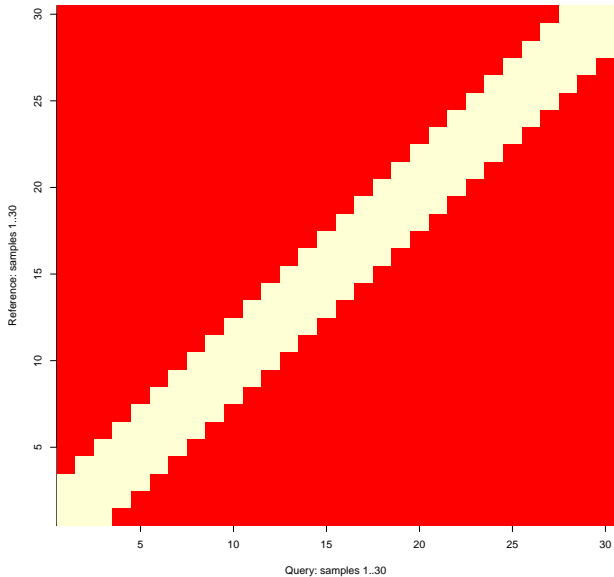
Dynamic Time Warping

Dynamic Time Warping



Sakoe-Chiba Window

Local Cost Matrix



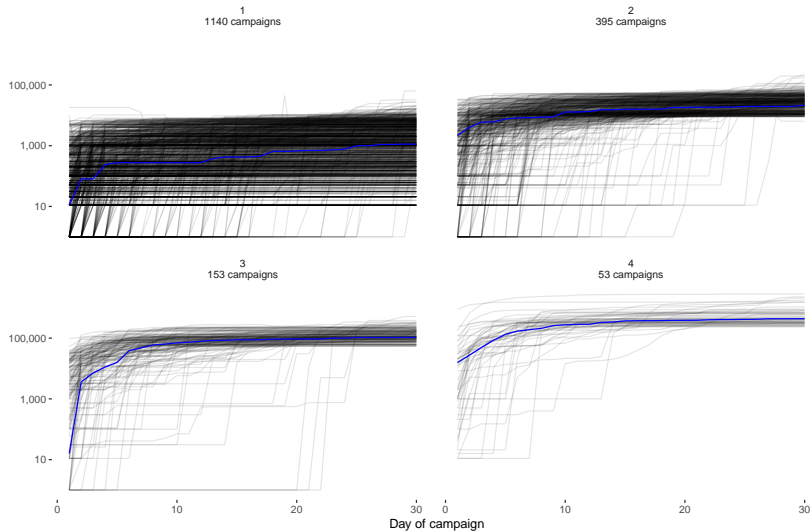
DTW Code

```
pc_dtw <- tsclust(balance_traj,  
  k = 4L,  
  distance = "dtw_basic",  
  seed = 1234,  
  norm = "L2",  
  window.size = 2L)
```

Results, $k=4$

Campaign Balance

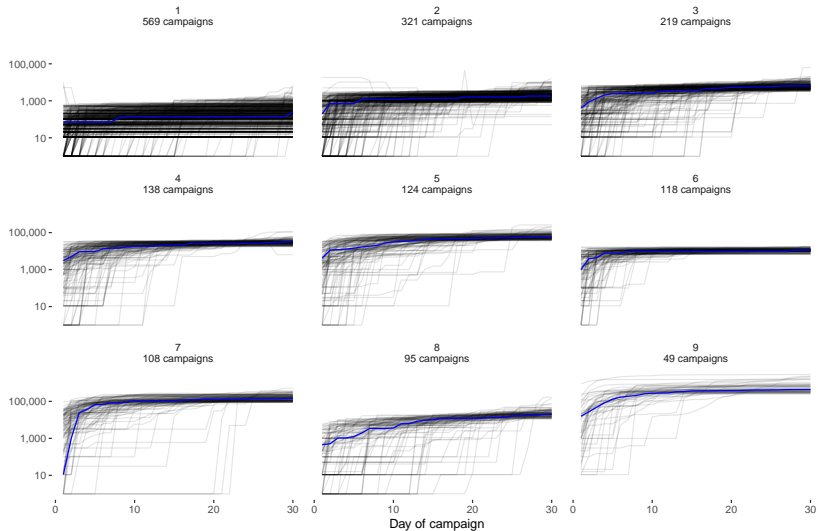
partitional clusters, dtw_basic distance, pam centroids, $k = 4$



Results, $k=9$

Campaign Balance

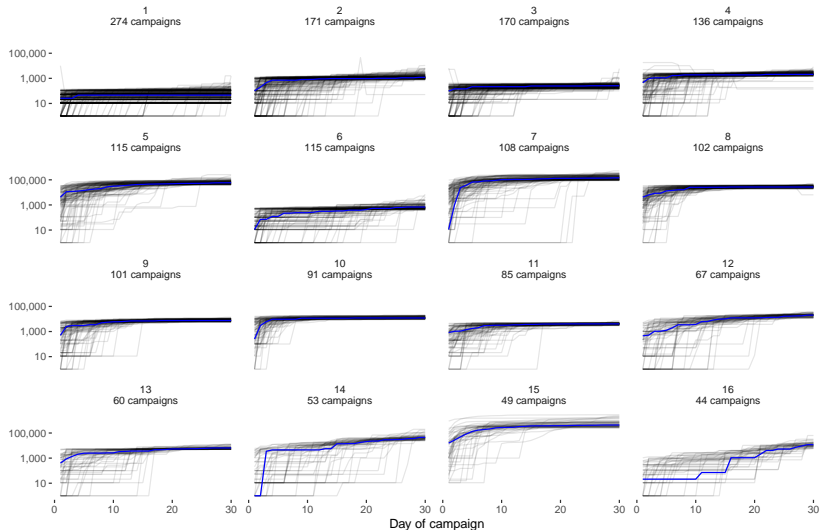
partitional clusters, dtw_basic distance, pam centroids, $k = 9$



Results, k=16

Campaign Balance

partitional clusters, dtw_basic distance, pam centroids, k = 16



k-Shape Clustering

Shape-based distance

Cross-correlation with shift

$$SBD(\vec{x}, \vec{y}) = 1 - \max_s \left(\frac{\vec{x}_{(s)} \cdot \vec{y}}{\sqrt{\|\vec{x}\|^2 \|\vec{y}\|^2}} \right)$$

$$\vec{x}_{(s)} = \begin{cases} \overbrace{(0, \dots, 0, x_1, x_2, \dots, x_{m-s})}^{|s|}, & s \geq 0 \\ (x_{1-s}, \dots, x_{m-1}, x_m, \underbrace{0, \dots, 0}_{|s|}), & s < 0 \end{cases}$$

Paparrizos J, Gravano L (2015). “k-Shape: Efficient and Accurate Clustering of Time Series.” In Proceedings of the 2015 ACM SIGMOD International Conference on Management of Data, {SIGMOD '15}, pp. 1855-1870. ACM, New York, NY, USA. ISBN 978-1-4503-2758-9. doi:10.1145/2723372.2737793.

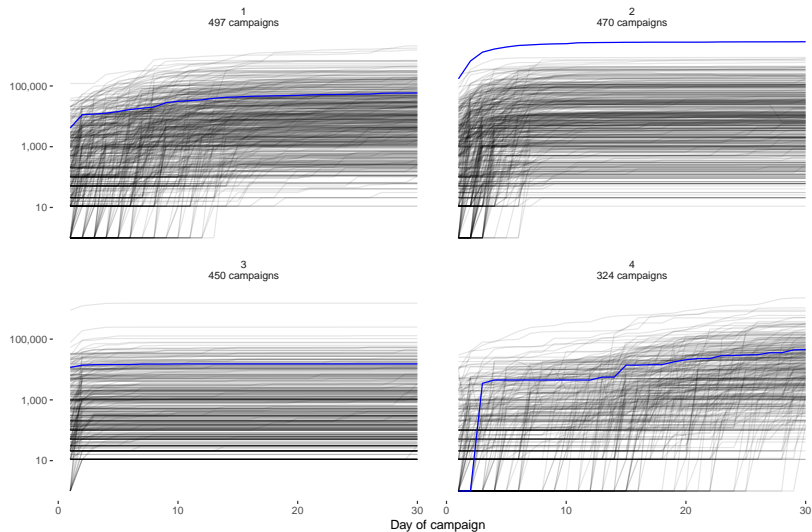
SBD Code

```
pc_sbd4 <- tsclust(  
  balance_traj,  
  type = "p",  
  k = 4L,  
  seed = 1234,  
  distance = "sbd"  
)
```

Results for $k=4$

Campaign Balance

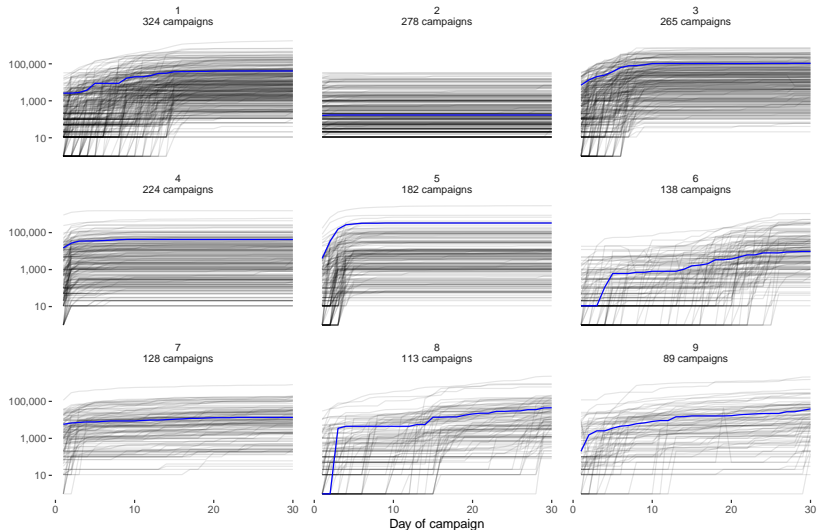
partitional clusters, sbd distance, pam centroids, $k = 4$



Results for $k=9$

Campaign Balance

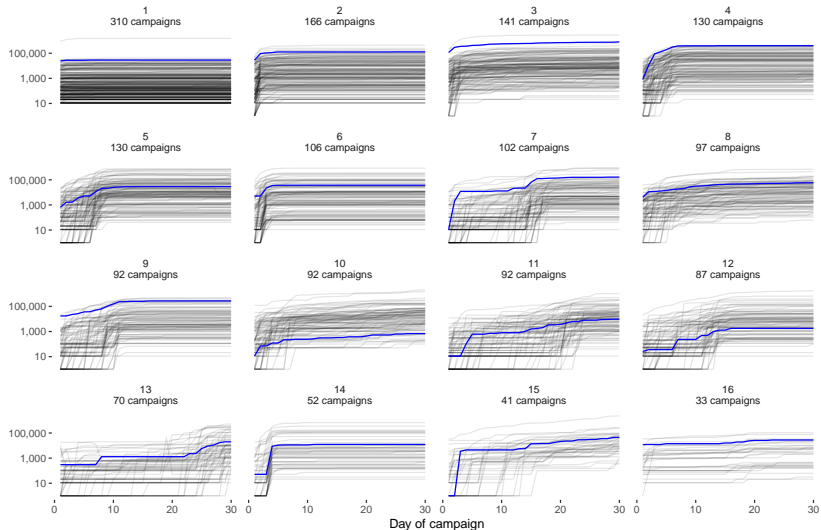
partitional clusters, sbd distance, pam centroids, $k = 9$



Results for $k=16$

Campaign Balance

partitional clusters, sbd distance, pam centroids, $k = 16$



Resources

R packages:

- ▶ `dtwclust`
- ▶ `tidyverse`

Handout Answers

