

PORTOFOLIO DATA SCIENTIST

Analisis pada perusahaan Fintech

Laporan UAS Data Mining



Disusun Oleh :

Muhammad Fikri Septiawan

1301160789

Program Studi S1 Informatika

Fakultas Informatika

Telkom University

2019

A. Latar Belakang Masalah

Fintech atau *Financial Technology* adalah istilah yang digunakan untuk menyebut inovasi dalam bidang jasa keuangan atau finansial. Bisa juga diartikan dengan inovasi finansial yang diberi sentuhan teknologi modern. Atau merupakan segmen di dunia *start-up*. Tujuannya adalah membantu untuk memaksimalkan penggunaan teknologi untuk mengubah, mempertajam atau mempercepat berbagai aspek pelayanan keuangan. *Fintech* mempunyai salah satu proses bisnis yaitu meminjamkan modal kepada suatu usaha yang sudah berjalan maupun yang sedang dibangun.

Proses bisnis tersebut, dibutuhkan suatu pengambilan keputusan apakah keputusan untuk meminjamkan modal kepada suatu usaha tepat atau tidak (mendapatkan profit yang cukup besar maupun jangka waktu pengembalian modal). Masalah pengambilan keputusan tersebut mempunyai suatu solusi yaitu menggunakan *data minning* dengan mengacu pada suatu *dataset fintech*.

B. Tujuan

Adapun tujuan melakukan *data minning* ini sebagai berikut:

1. Untuk pengambilan keputusan peminjaman modal kepada suatu usaha,
2. Mendapatkan model yang tepat seperti syarat tertentu untuk dijadikan pengambilan keputusan pada *point* 1,
3. Mendapatkan informasi tentang perkembangan usaha yang telah diberikan modal oleh *fintech* ini, jika usahanya berkembang maka akan semakin mudah pencairan modal saat dibutuhkan oleh usaha tersebut.

C. Deskripsi Data

Dataset dengan nama *data-uas.xlsx* (bisa diunduh pada <https://github.com/MFikriS/DataScientist/blob/master/Portofolio2/data-uas.xlsx>) lalu harus di-*convert* menjadi *data-uas.csv* agar dapat diolah oleh bahasa pemrograman yaitu bahasa pemrograman python. Sesuai soal UAS digunakanlah baris data ke 500-1500 (1000 data), tidak mempunyai label dan memiliki 11 fitur (kolom) yaitu:

branch	cabang di Fintech, umumnya di level kecamatan.
cutoff_date	tanggal batas.
area	area yang mencakup cabang-cabang di bawahnya, umumnya di level kabupaten.
region	region yang mencakup area-area di bawahnya, level provinsi.
first_date_disbursement	tanggal pencairan pinjaman pertama di setiap cabang.
active_borrowers	jumlah mitra yang sedang memiliki pinjaman di setiap cabang.

D. Praproses

Hasil codingan menggunakan bahasa pemrograman dapat diunduh pada <https://github.com/MFikriS/DataScientist/blob/master/Portofolio2/uas.ipynb>

1. Mengubah tipe data untuk fitur *active_borrowers*, *active_agent*, *delinquency_rate*, *outstanding*, *weekly_disbursement*, *weekly_new_borrower_per_bp* dari bertipe *object* menjadi *numerical*. Mengubah tipe data untuk *cutoff_date*, *first_date_disbursement* menjadi *date*. Hasil pengubahan dapat dilihat pada gambar 4.

```
In [201]: data['active_borrowers'] = pd.to_numeric(data['active_borrowers'], errors='coerce')
data['active_agent'] = pd.to_numeric(data['active_agent'], errors='coerce')
data['outstanding'] = pd.to_numeric(data['outstanding'], errors='coerce')
data['weekly_disbursement'] = pd.to_numeric(data['weekly_disbursement'], errors='coerce')
data['weekly_new_borrower_per_bp'] = pd.to_numeric(data['weekly_new_borrower_per_bp'], errors='coerce')

In [202]: data.dtypes
Out[202]: branch                object
cutoff_date                  object
area                        object
region                     object
first_date_disbursement      object
active_borrowers            float64
active_agent                float64
delinquency_rate            object
outstanding                 float64
weekly_disbursement         float64
weekly_new_borrower_per_bp  float64
dtype: object

In [129]: dataset['cutoff_date'] = pd.to_datetime(dataset.cutoff_date)
dataset['first_date_disbursement'] = pd.to_datetime(dataset.first_date_disbursement)
```

Gambar 4. Convert tipe data dan hasil convert

Hal tersebut dilakukan untuk proses perhitungan yang akan digunakan pada algoritma *data minning*.

2. Pemilihan baris data

Pemilihan baris data dilakukan untuk pemilihan pola-pola data yang tepat. Melakukan proses tersebut karena setiap baris data memiliki pola-pola yang berbeda-beda. Pemilihan data dapat dilihat pada gambar 5.

```
In [9]: dataset = data.iloc[500:1500,0:11]
```

Gambar 5. Pemilihan baris data

3. Mengubah tipe data fitur *delinquency_rate*, *region* menjadi tipe *categorical* karena fitur tersebut bertipe data *object* atau *string* (susah untuk diolah). Dilakukan untuk mengetahui tingkat pengaruh atribut dalam pengolahan data. Pra proses ini dapat dilihat pada gambar 6.

```

In [12]: a = dataset.delinquency_rate

In [13]: a.value_counts()
Out[13]:
0%      760
2.5%-5%   77
0%-1%     64
1%-2.5%   57
5%-10%    48
>10%     44
Name: delinquency_rate, dtype: int64

In [14]: def conditions(a):
    if (a['delinquency_rate'] == "0%"):
        return 0
    elif (a['delinquency_rate'] == "0%-1%"):
        return 1
    elif (a['delinquency_rate'] == "1%-2.5%"):
        return 2
    elif (a['delinquency_rate'] == "2.5%-5%"):
        return 3
    elif (a['delinquency_rate'] == "5%-10%"):
        return 4
    else:
        return 5

In [15]: dataset['delinquency_rate'] = dataset.apply(conditions, axis=1)
dataset['delinquency_rate'] = dataset['delinquency_rate']

In [34]: g = dataset.region

In [347]: def conditionsg(g):
    if (g['region'] == "REGION_00"):
        return 0
    elif (g['region'] == "REGION_01"):
        return 1
    elif (g['region'] == "REGION_02"):
        return 2
    elif (g['region'] == "REGION_03"):
        return 3
    else:
        return 4

In [348]: dataset['region'] = dataset.apply(conditionsg, axis=1)
dataset['region'] = dataset['region']

```

Gambar 6. Convert tipe data fitur delinquency_rate

4. Cek Missing Value pada atribut atau fitur dan mengisinya dengan teknik tertentu. Hal ini dilakukan untuk proses perhitungan agar menghasilkan perhitungan yang tepat.

```

CEK MISSING VALUE

In [16]: dataset.delinquency_rate.isnull().values.any()
Out[16]: False

In [17]: dataset.active_borrowers.isnull().values.any()
Out[17]: False

In [18]: dataset.active_agent.isnull().values.any()
Out[18]: False

In [19]: dataset.outstanding.isnull().values.any()
Out[19]: False

In [20]: dataset.weekly_disbursement.isnull().values.any()
Out[20]: False

In [21]: dataset.weekly_new_borrower_per_bj.isnull().values.any()
Out[21]: True

MENGGISI VALUE PADA FITUR YANG TERDAPAT MISSING VALUE

In [22]: dataset['weekly_new_borrower_per_bj'] = dataset['weekly_new_borrower_per_bj'].fillna(dataset['weekly_new_borrower_per_bj'].mean())
dataset['weekly_new_borrower_per_bj'] = round(dataset['weekly_new_borrower_per_bj'])

```

Gambar 7. Cek missing value dan pengisian missing value

5. Membuat kolom atau fitur baru yaitu jangka waktu pinjam karena fitur tersebut berpeluang untuk pengaruh pada pengolahan data. Fitur baru tersebut didapat dari pengurangan fitur *cutoff_date* dikurangi fitur *first_date_disbursement*. Pra proses tersebut dapat dilihat pada gambar 8.

```

In [131]: dataset['jangka_waktu_pinjam'] = dataset['cutoff_date'] - dataset['first_date_disbursement']

In [132]: dataset['jangka_waktu_pinjam'] = dataset['jangka_waktu_pinjam']

In [133]: dataset.corr()

Out[133]:

```

	active_borrowers	active_agent	delinquency_rate	outstanding	weekly_disbursement	weekly_new_borrower_per_bp
active_borrowers	1.000000	0.316628	0.362378	0.978388	0.225795	-0.515201
active_agent	0.316628	1.000000	0.417417	0.902278	0.107484	-0.503721
delinquency_rate	0.362378	0.417417	1.000000	0.526228	-0.172203	-0.318178
outstanding	0.978388	0.902278	0.526228	1.000000	0.251287	-0.499113
weekly_disbursement	0.225795	0.107484	-0.172203	0.251287	1.000000	0.227136
weekly_new_borrower_per_bp	-0.515201	-0.503721	-0.318178	-0.499113	0.227136	1.000000

Gambar 8. Pembuatan fitur baru yaitu jangka waktu pinjam

- Melakukan drop fitur yang kurang berpengaruh pada pengolahan data. Hal ini dilakukan karena untuk meningkatkan nilai akurasi maupun *score*, dapat dilihat pada gambar 9.

UJI COBA 1 DENGAN MEN-DROP FITUR BRANCH, CUTOFF_DATE, AREA, REGION, FIRST_DATE_DISBURSEMENT

```

In [26]: dataset1 = dataset.drop(['branch', 'cutoff_date', 'area', 'region', 'first_date_disbursement'], axis=1)

```

Gambar 9. Drop fitur yang kurang berpengaruh

E. Analisis Pemilihan Algoritma

Berdasarkan kondisi dataset di atas yaitu tidak adanya label, maka metode yang tepat untuk dipilih adalah metode *clustering*. Metode *clustering* digunakan untuk mengelompokkan satu data ke data lain berdasarkan suatu atribut atau fitur. Metode *clustering* contoh algoritmanya yaitu *K-means*, *Agglomeratif Hierarchical Clustering*, dan lainnya. Dari contoh algoritma *clustering* tersebut, dipilihlah algoritma *K-means* karena datanya mayoritas bertipe data *numerical* dan ngacak.

F. Analisis Penentuan Parameter

- Pemilihan atribut atau fitur,
- Pemilihan baris data (*row*),
- Pemilihan metode klasifikasi atau *clustering*,
- Pemilihan algoritmanya, jumlah *cluster*,
- Nilai *silhouette score*.

G. Hasil Percobaan

- Pemilihan atribut atau fitur

Fitur yang sangat berpengaruh pada pengolahan data *minning* disini adalah *active_borrowers*, *active_agent*, *delinquency_rate*, *outstanding*, *jangka_waktu_pinjam*.

```

In [134]: dataset4 = dataset.drop(['branch', 'cutoff_date', 'area', 'region', 'first_date_disbursement', 'weekly_new_borrower_per_bp'], axis=1)

In [135]: dataset4.corr()

Out[135]:

```

	active_borrowers	active_agent	delinquency_rate	outstanding	weekly_disbursement	jangka_waktu_pinjam
active_borrowers	1.000000	0.316628	0.362378	0.978388	0.225795	0.528332
active_agent	0.316628	1.000000	0.417417	0.902278	0.107484	0.546383
delinquency_rate	0.362378	0.417417	1.000000	0.526228	-0.172203	0.577493
outstanding	0.978388	0.902278	0.526228	1.000000	0.251287	0.796218
weekly_disbursement	0.225795	0.107484	-0.172203	0.251287	1.000000	0.029486
jangka_waktu_pinjam	0.528332	0.546383	0.577493	0.796218	0.029486	1.000000

```

In [136]: dataset4 = dataset4.drop(['weekly_disbursement'], axis=1)

```

Gambar 10. Drop fitur yang kurang berpengaruh dan yang tepat

2. Pemilihan baris data (row).

Dipilih row data ke- 500-999 (500 data) karena mempengaruhi nilai *silhouette score*.

```
In [162]: dataset46 = dataset4.iloc[500:1000,0:5]

In [164]: dataset46

Out[164]:
```

	active_borrowers	active_agents	delinquency_rate	notdelinquent	jumlah_waktu_proses
500	1213.0	3.0	0	1800.0	300
501	410.0	3.0	0	1000.0	133
502	400.0	3.0	0	1200.0	112
503	420.0	3.0	0	1000.0	130
504	400.0	2.0	0	1300.0	119
505	370.0	2.0	0	1100.0	77
506	37.0	2.0	0	900.0	21
507	2700.0	12.0	0	3000.0	530
508	1000.0	7.0	0	4000.0	371
509	420.0	3.0	0	1000.0	112
510	545.0	2.0	0	400.0	35
511	104.0	1.0	0	300.0	42
512	1394.0	0.0	0	2000.0	407
513	1131.0	5.0	0	2400.0	371
514	100.0	2.0	0	500.0	50
515	300.0	0.0	0	800.0	47

Gambar 11. Pemilihan baris data yang tepat

3. Pemilihan metode klasifikasi atau *clustering*.

Menggunakan *clustering* karena dataset tidak mempunyai label, jika menggunakan klasifikasi maka harus ada label (*supervised*).

4. Pemilihan algoritma *clustering* dan jumlah *cluster*.

Algoritma disini menggunakan *Agromerativ Hierarchical Clustering* dan *K-Means* dan jumlah cluster yang bagus yaitu sebanyak 2.

```
In [166]: kmeans = KMeans(n_clusters=2, max_iter=10000)
dataset46['cluster'] = kmeans.fit_predict(dataset46)

from sklearn.metrics import silhouette_score
score = silhouette_score(dataset46, dataset46['cluster'], metric='euclidean')
score

Out[166]: 0.648117911295182

In [170]: kmeans = KMeans(n_clusters=2, max_iter=10000)
dataset46['cluster'] = kmeans.fit_predict(dataset46)

from sklearn.metrics import silhouette_score
score = silhouette_score(dataset46, dataset46['cluster'], metric='euclidean')
score

Out[170]: 0.6423098670608772
```

Gambar 12. Algoritma Kmeans

```
In [203]: from sklearn.cluster import AgglomerativeClustering
hc = AgglomerativeClustering(n_clusters=2, affinity='euclidean', linkage='ward')
y_hc = hc.fit_predict(dataset46)

In [204]: from sklearn.metrics import silhouette_score
score = silhouette_score(dataset46, y_hc, metric='euclidean')
score

Out[204]: 0.6363861009174950

In [205]: from sklearn.cluster import AgglomerativeClustering
hc = AgglomerativeClustering(n_clusters=2, affinity='euclidean', linkage='ward')
y_hc = hc.fit_predict(dataset46)

from sklearn.metrics import silhouette_score
score = silhouette_score(dataset46, y_hc, metric='euclidean')
score

Out[205]: 0.638930400910295

In [ ]:
```

Gambar 13. Algoritma Agromerativ Hierarchical Clustering

5. Nilai *silhouette score*.

Sebagai penilaian kualitas dari *cluster*, semakin mendekati nilai 1 maka semakin bagus kualitas *cluster*.

```
In [166]: kmeans = KMeans(n_clusters=2, max_iter=10000)
dataset46['cluster'] = kmeans.fit_predict(dataset46)

from sklearn.metrics import silhouette_score
score = silhouette_score(dataset46, dataset46['cluster'], metric='euclidean')
score

Out[166]: 0.6493175112255162

In [176]: kmeans = KMeans(n_clusters=2, max_iter=10000)
dataset46['cluster'] = kmeans.fit_predict(dataset46)

from sklearn.metrics import silhouette_score
score = silhouette_score(dataset46, dataset46['cluster'], metric='euclidean')
score

Out[176]: 0.642369867968772
```

Gambar 14. Algoritma Kmeans

```
In [203]: from sklearn.cluster import AgglomerativeClustering
hc = AgglomerativeClustering(n_clusters=2, affinity='euclidean', linkage='ward')
y_hc = hc.fit_predict(dataset46)

In [204]: from sklearn.metrics import silhouette_score
score = silhouette_score(dataset46, y_hc, metric='euclidean')
score

Out[204]: 0.6385358408016295

In [205]: from sklearn.cluster import AgglomerativeClustering
hc = AgglomerativeClustering(n_clusters=2, affinity='euclidean', linkage='ward')
y_hc = hc.fit_predict(dataset46)

from sklearn.metrics import silhouette_score
score = silhouette_score(dataset46, y_hc, metric='euclidean')
score

Out[205]: 0.6385358408016295

In [ ]:
```

Gambar 15. Algoritma Agromerativ Hierarchical Clustering

Berdasarkan gambar 14 dan gambar 15, algoritma *K-Means* lebih bagus daripada algoritma *Agromerativ Hierarchical Clustering* karena nilai *silhouette score* *K-Means* (0.6493175112255162) lebih besar dari nilai *silhouette score* *Agromerativ Hierarchical Clustering* (0.6385358408016295).

Dari hasil percobaan dengan analisis-analisis di atas, terdapat 2 kelompok atau *cluster* yaitu 0 dan 1. *Cluster* 0 berisi potensi *profit* besar bagi *fintech* dan *Cluster* 1 berisi potensi *profit* kecil bagi *fintech*. Dari *clustering* tersebut *fintech* dalam pengambilan keputusan dapat mengacu pada model *cluster* tersebut, Jika pemohon pinjaman merupakan *cluster* 0 maka *fintech* berani meminjamkan modal yang cukup besar dengan potensi *profit* cukup besar.

H. Ringkasan Model Yang Diperoleh

Dari model yang telah dibangun pada point G, didapatkan:

Terdapat 2 kelompok atau *cluster* yaitu 0 dan 1. *Cluster* 0 berisi potensi *profit* besar bagi *fintech* dan *Cluster* 1 berisi potensi *profit* kecil bagi *fintech* dan *Cluster* 1 berisi potensi *profit* kecil bagi *fintech*. Dari *clustering* tersebut *fintech* dalam pengambilan keputusan dapat mengacu pada model *cluster* tersebut, Jika pemohon pinjaman merupakan *cluster* 0 maka *fintech* berani meminjamkan modal yang cukup besar dengan potensi *profit* cukup besar juga.

I. Interpretasi Model

1. Import library numpy, pandas, dan KMeans

```
In [157]: import numpy as np
import pandas as pd
from sklearn.cluster import KMeans
```

2. Drop fitur yang kurang berpengaruh dalam proses data *minning*, yaitu *branch*, *cutoff_date*, *area*, *region*, *first_date*, *disbursement*, *weekly_new_borrower_per_bp*, *weekly_disbursement*

```
In [154]: dataset4 = dataset.drop(['branch', 'cutoff_date', 'area', 'region', 'first_date', 'disbursement', 'weekly_new_borrower_per_bp'], axis=1)
```

```
In [156]: dataset4 = dataset4.drop(['weekly_disbursement'], axis=1)
```

3. Pemilihan *row* atau baris data yang tepat yaitu baris ke-500-999

```
In [162]: dataset46 = dataset4.iloc[0:500,0:5]
```

4. Model Algoritma *K-Means Clustering* dengan jumlah 2 melakukan iterasi maksimum 10000 dan penilaian menggunakan *silhouette score*

```
In [166]: kmeans = KMeans(n_clusters=2, max_iter=10000)
dataset46['cluster'] = kmeans.fit_predict(dataset46)

from sklearn.metrics import silhouette_score
score = silhouette_score(dataset46, dataset46['cluster'], metric='euclidean')
score
```

```
Out[166]: 0.6493175112255162
```

```
In [176]: kmeans = KMeans(n_clusters=3, max_iter=10000)
dataset46['cluster'] = kmeans.fit_predict(dataset46)

from sklearn.metrics import silhouette_score
score = silhouette_score(dataset46, dataset46['cluster'], metric='euclidean')
score
```

```
Out[176]: 0.6429690870560772
```

5. Hasil *cluster*

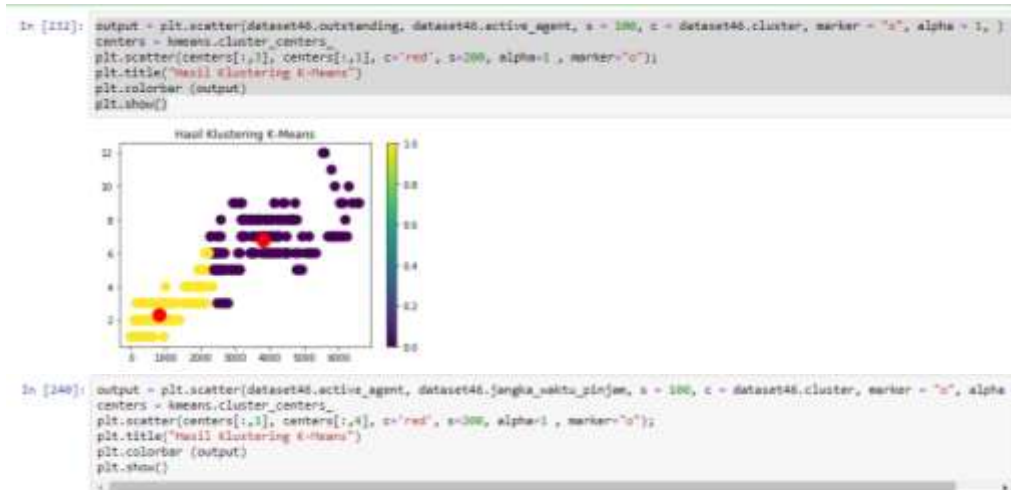
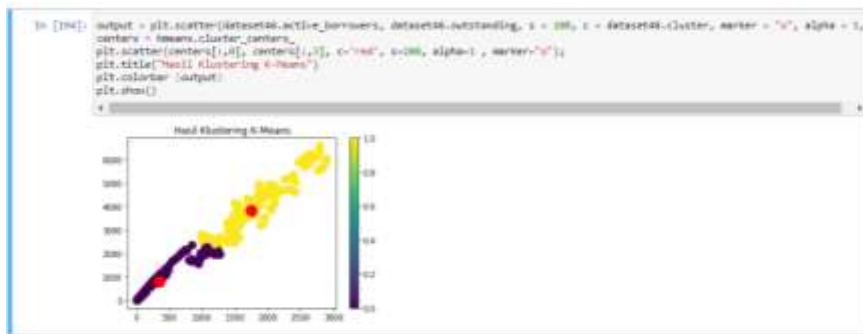
```
In [186]: dataset46['cluster'] = dataset46.cluster
```

```
In [187]: dataset46
```

```
Out[187]:
```

	active_borrowers	active_agents	delinquency_rate	contending	avglo_weeks_ginsen	cluster
800	5213.0	3.0	0	1000.0	390	0
801	416.0	3.0	0	1000.0	113	0
802	489.0	3.0	0	1200.0	112	0
803	423.0	3.0	0	1000.0	113	0
804	489.0	2.0	0	1300.0	118	0
806	575.0	2.0	0	1100.0	77	0
808	37.0	2.0	0	100.0	21	0
807	2703.0	12.0	0	950.0	939	1
808	5889.0	7.0	5	4640.0	371	1
808	423.0	3.0	0	1000.0	112	0
818	149.0	2.0	0	400.0	36	0
811	116.0	1.0	0	300.0	42	0
812	1204.0	9.0	3	2950.0	497	1
813	1531.0	9.0	0	2490.0	371	1
814	180.0	2.0	0	500.0	56	0
818	289.0	2.0	0	800.0	42	0
819	9899.0	9.0	3	4290.0	516	1
817	2656.0	7.0	0	6250.0	415	1
818	1327.0	5.0	0	1900.0	388	0
819	791.0	2.0	0	300.0	36	0

6. Visualisasi *cluster*



Model diatas menunjukan tujuan pada point B telah tercapai, karena:

1. Model tersebut berisi 2 kelompok atau *cluster*:

- *Cluster 0*

Merupakan potensi *profit* besar bagi *fintech*

- *Cluster 1*

Merupakan potensi *profit* kecil bagi *fintech*

Dari 2 *cluster* tersebut dapat digunakan pengambilan keputusan yaitu pemberian pinjaman modal dengan cukup banyak nominalnya kepada *cluster* 0 karena akan mendapat *profit* yang lebih besar juga daripada peminjaman *cluster* 1. Karena pada *cluster* 0 terdapat *jangka_waktu_pinjam* yang lama yaitu > 300 hari maka bunga yang didapat semakin tinggi, hal ini akan mendapatkan profit yang besar. Selain itu jumlah agen yang aktif bekerja ≥ 6 pada *cluster* 0 akan mengawasi pinjaman yang diberikan, hal ini akan meminimalkan resiko (kerugian) yang ada.

2. *Fintech* memberikan syarat pinjaman yaitu jika jumlah mitra yang sedang memiliki pinjaman di setiap cabang berjumlah sedikit (0-700), sisa pinjaman yang belum terbayar oleh mitra berjumlah sedikit (0-2100) juga, jumlah agen yang aktif bekerja di setiap cabang berjumlah banyak (≥ 6), *jangka_waktu_pinjam* > 300 hari (akan menghasilkan bunga yang besar maka akan mendapatkan profit yang besar juga), maka usaha atau mitra dapat meminjam lagi
3. Perkembangan usaha/mitra dapat dilihat pada model yang telah dibangun. Semakin berkembang usaha/mitra maka semakin patuh atau cepat dalam membayar pinjaman (dilihat dari jumlah mitra yang sedang memiliki pinjaman di setiap cabang yang sedikit dan sisa pinjaman yang belum terbayar oleh mitra berjumlah sedikit pada *cluster* 0). Hal ini akan memberikan *profit* untuk *fintech*.