

Which Debts Are Worth the Bank's Effort

Songtao Song

1. Regression discontinuity: banking recovery

After a debt has been legally declared "uncollectable" by a bank, the account is considered to be "charged-off." But that doesn't mean the bank simply walks away from the debt. They still want to collect some of the money they are owed. The bank will score the account to assess the expected recovery amount, that is, the expected amount that the bank may be able to receive from the customer in the future (for a fixed time period such as one year). This amount is a function of the probability of the customer paying, the total debt, and other factors that impact the ability and willingness to pay.

The bank has implemented different recovery strategies at different thresholds (\$1000, \$2000, etc.) where the greater the expected recovery amount, the more effort the bank puts into contacting the customer. For low recovery amounts (Level 0), the bank just adds the customer's contact information to their automatic dialer and emailing system. For higher recovery strategies, the bank incurs more costs as they leverage human resources in more efforts to contact the customer and obtain payments. Each additional level of recovery strategy requires an additional \$50 per customer so that customers in the Recovery Strategy Level 1 cost the company \$50 more than those in Level 0. Customers in Level 2 cost \$50 more than those in Level 1, etc.

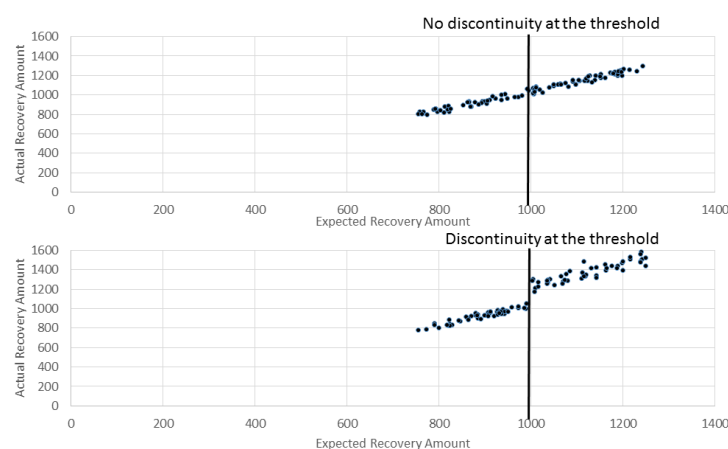


Figure 1. Demonstration of Regression Gap

The big question: does the extra amount that is recovered at the higher strategy level exceed the extra \$50 in costs? In other words, was there a jump (also called a "discontinuity" as shown in Figure 1) of more than \$50 in the amount recovered at the higher strategy level?

First, we'll load the banking dataset and look at the first few rows of data. This puts us in a good position to understand the dataset itself and begin thinking about how to analyze the data.

```
1. # Import modules
2. import pandas as pd
3. import numpy as np
4.
5. # Read in dataset
6. df = pd.read_csv('datasets/bank_data.csv')
7.
8. # Print the first few rows of the DataFrame
9. df.head()
```

2. Graphical exploratory data analysis

The bank has implemented different recovery strategies at different thresholds (\$1000, \$2000, \$3000 and \$5000) where the greater the Expected Recovery Amount, the more effort the bank puts into contacting the customer. Zeroing in on the first transition (between Level 0 and Level 1) means we are focused on the population with Expected Recovery Amounts between \$0 and \$2000 where the transition between Levels occurred at \$1000. We know that the customers in Level 1 (expected recovery amounts between \$1001 and \$2000) received more attention from the bank and, by definition, they had higher Expected Recovery Amounts than the customers in Level 0 (between \$1 and \$1000).

Here's a quick summary of the Levels and thresholds again:

- Level 0: Expected recovery amounts >\$0 and <=\$1000
- Level 1: Expected recovery amounts >\$1000 and <=\$2000
- The threshold of \$1000 separates Level 0 from Level 1

A key question is whether there are other factors besides Expected Recovery Amount that also varied systematically across the \$1000 threshold. For example, does the customer age show a jump (discontinuity) at the \$1000 threshold or does that age vary smoothly? We can examine this by first making a scatter plot of the age as a function of Expected Recovery Amount for a small window of Expected Recovery Amount, \$0 to \$2000. This range covers Levels 0 and 1.

```
1. # Scatter plot of Age vs. Expected Recovery Amount
2. from matplotlib import pyplot as plt
3. %matplotlib inline
4. plt.scatter(x=df['expected_recovery_amount'], y=df['age'], c="g", s=2)
5. plt.xlim(0, 2000)
6. plt.ylim(0, 60)
7. plt.xlabel('Expected Recovery Amount')
8. plt.ylabel('Age')
9. plt.legend(loc=2)
10. plt.show()
```

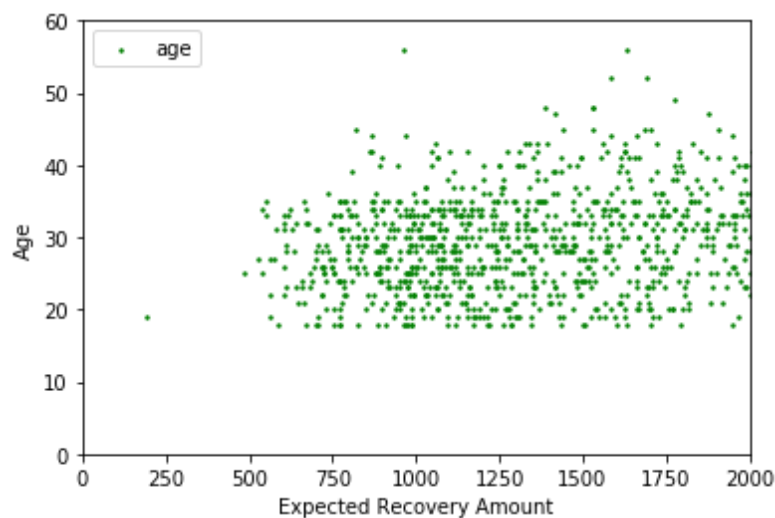


Figure 2 Expected Recovery Amount vs. Age

3. Statistical test: age vs. expected recovery amount

We want to convince ourselves that variables such as age and sex are similar above and below the \$1000 Expected Recovery Amount threshold. This is important because we want to be able to conclude that differences in the actual recovery amount are due to the higher Recovery Strategy and not due to some other difference like age

or sex.

The scatter plot of age versus Expected Recovery Amount as Figure 2 did not show an obvious jump around \$1000. We will be more confident in our conclusions if we do statistical analysis examining the average age of the customers just above and just below the threshold. We can start by exploring the range from \$900 to \$1100.

For determining if there is a difference in the ages just above and just below the threshold, we will use the Kruskal-Wallis test which is a statistical test that makes no distributional assumptions.

```
1. # Import stats module
2. from scipy import stats
3.
4. # Compute average age just below and above the threshold
5. era_900_1100 = df.loc[(df['expected_recovery_amount']<1100) &
6.                        (df['expected_recovery_amount']>=900)]
7. by_recovery_strategy = era_900_1100.groupby(['recovery_strategy'])
8.
9. # Perform Kruskal-Wallis test
10. Level_0_age = era_900_1100.loc[df['recovery_strategy']=="Level 0 Recovery"]['age']
11. Level_1_age = era_900_1100.loc[df['recovery_strategy']=="Level 1 Recovery"]['age']
12. print(stats.kruskal(Level_0_age, Level_1_age))
13.
14. by_recovery_strategy['age'].describe()
```

The test resulted a p-value as 0.06297556896097407, which indicated that we should not reject null hypothesis as no significant difference between Level 0 and 1 on age under significant level 0.05.

4. Statistical test: sex vs. expected recovery amount

We were able to convince ourselves that there is no major jump in the average customer age just above and just below the \$1000 threshold by doing a statistical test as well as exploring it graphically with a scatter plot.

We want to also test that the percentage of customers that are male does not jump as well across the \$1000 threshold. We can start by exploring the range of \$900 to

\$1100 and later adjust this range.

We can examine this question statistically by developing cross-tabs as well as doing chi-square tests of the percentage of customers that are male vs. female.

```
1. # Number of customers in each category
2. crosstab = pd.crosstab(df.loc[(df['expected_recovery_amount'] < 2000) &
3.                               (df['expected_recovery_amount'] >= 0)], ['recovery
   _strategy'],
4.                               df['sex'])
5.
6. # Chi-square test
7. chi2_stat, p_val, dof, ex = stats.chi2_contingency(crosstab)
8. print("Due to the p-value of chi-
   square test on sex is " + str(p_val) + ", we can't reject null hypothesis.")
9.
10. crosstab
```

The result showed that the p-value of chi-square test on gender is 0.39416505436. Thus, we can't reject null hypothesis. In other word, there is no significant difference between Level 0 and Level 1 on genders.

5. Exploratory graphical analysis: recovery amount

We are now reasonably confident that customers just above and just below the \$1000 threshold are, on average, similar in terms of their average age and the percentage that are male. It is now time to focus on the key outcome of interest, the actual recovery amount.

A first step in examining the relationship between the actual recovery amount and the expected recovery amount is to develop a scatter plot where we want to focus our attention at the range just below and just above the threshold. Specifically, we will develop a scatter plot of Expected Recovery Amount (Y) vs. Actual Recovery Amount (X) for Expected Recovery Amounts between \$900 to \$1100. This range covers Levels 0 and 1. A key question is whether or not we see a discontinuity (jump) around the \$1000 threshold.

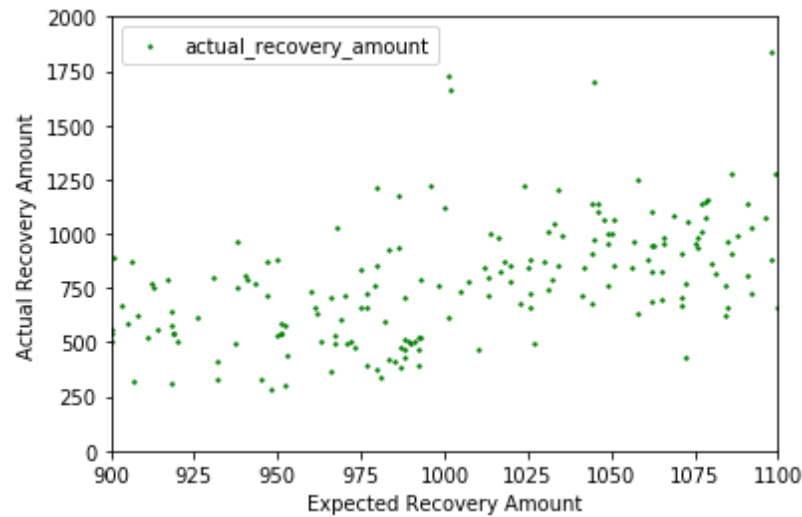


Figure 3. Expected Recovery Amount vs. Actual Recovery Amount

6. Statistical analysis: recovery amount

Figure 3 indicated a possible jump around \$1000 threshold. Just as we did with age, we can perform statistical tests to see if the actual recovery amount has a discontinuity above the \$1000 threshold. We are going to do this for two different windows of the expected recovery amount \$900 to \$1100 and for a narrow range of \$950 to \$1050 to see if our results are consistent.

Again, the statistical test we will use is the Kruskal-Wallis test, a test that makes no assumptions about the distribution of the actual recovery amount.

We will first compute the average actual recovery amount for those customers just below and just above the threshold using a range from \$900 to \$1100. Then we will perform a Kruskal-Wallis test to see if the actual recovery amounts are different just above and just below the threshold. Once we do that, we will repeat these steps for a smaller window of \$950 to \$1050.

```
1. # Compute average actual recovery amount just below and above the threshold
2. print(by_recovery_strategy['actual_recovery_amount'].describe().unstack())
3.
4. # Perform Kruskal-Wallis test
5. Level_0_actual = era_900_1100.loc[df['recovery_strategy']=='Level 0 Recovery
   ']['actual_recovery_amount']
```

```
6. Level_1_actual = era_900_1100.loc[df['recovery_strategy']=='Level 1 Recovery']
   ['actual_recovery_amount']
7. print(stats.kruskal(Level_0_actual, Level_1_actual))
8.
9. # Repeat for a smaller range of $950 to $1050
10. era_950_1050 = df.loc[(df['expected_recovery_amount']<1050) &
11.                        (df['expected_recovery_amount']>=950)]
12. Level_0_actual = era_950_1050.loc[df['recovery_strategy']=='Level 0 Recovery']
   ['actual_recovery_amount']
13. Level_1_actual = era_950_1050.loc[df['recovery_strategy']=='Level 1 Recovery']
   ['actual_recovery_amount']
14. print(stats.kruskal(Level_0_actual, Level_1_actual))
```

The result showed that the p-value for both windows are 6.177308752803109e-16 and 3.80575314300276e-08 respectively, which indicated a significant difference between Level 0 and Level 1 on actual recovery amount.

7. Regression modeling: no threshold

We now want to take a regression-based approach to estimate the impact of the program at the \$1000 threshold using the data that is just above and just below the threshold. In order to do that, we will build two models. The first model does not have a threshold while the second model will include a threshold.

The first model predicts the actual recovery amount (outcome or dependent variable) as a function of the expected recovery amount (input or independent variable). We expect that there will be a strong positive relationship between these two variables.

We will examine the adjusted R-squared to see the percent of variance that is explained by the model. In this model, we are not trying to represent the threshold but simply trying to see how the variable used for assigning the customers (expected recovery amount) relates to the outcome variable (actual recovery amount).

```
1. # Import statsmodels
2. import statsmodels.api as sm
3.
4. # Define X and y
5. X = era_900_1100['expected_recovery_amount']
```



```

6. y = era_900_1100['actual_recovery_amount']
7. X = sm.add_constant(X)
8.
9. # Build linear regression model
10. model = sm.OLS(y, X).fit()
11. predictions = model.predict(X)
12.
13. # Print out the model summary statistics
14. print(model.summary())

```

And the results were:

```

=====
                        OLS Regression Results
=====
Dep. Variable:      actual_recovery_amount    R-squared:                0.261
Model:              OLS                     Adj. R-squared:           0.256
Method:             Least Squares           F-statistic:             63.78
Date:               Wed, 20 Mar 2019         Prob (F-statistic):      1.56e-13
Time:               01:44:08                Log-Likelihood:          -1278.9
No. Observations:   183                     AIC:                    2562.
Df Residuals:       181                     BIC:                    2568.
Df Model:           1
Covariance Type:    nonrobust
=====
                        coef      std err          t      P>|t|      [0.025      0.975]
-----
const                -1978.7597    347.741     -5.690     0.000    -2664.907    -1292.612
expected_recovery_amount  2.7577      0.345      7.986     0.000      2.076      3.439
=====
Omnibus:             64.493    Durbin-Watson:           1.777
Prob(Omnibus):       0.000    Jarque-Bera (JB):         185.818
Skew:                1.463    Prob(JB):                 4.47e-41
Kurtosis:            6.977    Cond. No.:                1.80e+04
=====

```

8. Regression modeling: adding true threshold

From the first model, we see that the regression coefficient is statistically significant for the expected recovery amount and the adjusted R-squared value was about 0.26. As we saw from the graph, on average the actual recovery amount increases as the expected recovery amount increases. We could add polynomial terms of expected recovery amount (such as the squared value of expected recovery amount) to the model but, for the purposes of this practice, let's stick with using just the linear term.

The second model adds an indicator of the true threshold to the model. If there was no impact of the higher recovery strategy on the actual recovery amount, then we would expect that the relationship between the expected recovery amount and the

actual recovery amount would be continuous. In this case, we know the true threshold is at \$1000.

We will create an indicator variable (either a 0 or a 1) that represents whether or not the expected recovery amount was greater than \$1000. When we add the true threshold to the model, the regression coefficient for the true threshold represents the additional amount recovered due to the higher recovery strategy. That is to say, the regression coefficient for the true threshold measures the size of the discontinuity for customers just above and just below the threshold. If the higher recovery strategy did help recovery more money, then the regression coefficient of the true threshold will be greater than zero. If the higher recovery strategy did not help recover more money than the regression coefficient will not be statistically significant.

```
1. # Create indicator (0 or 1) for expected recovery amount >= $1000
2. df['indicator_1000'] = np.where(df['expected_recovery_amount']<1000, 0, 1)
3. era_900_1100 = df.loc[(df['expected_recovery_amount']<1100) &
4.                        (df['expected_recovery_amount']>=900)]
5.
6. # Define X and y
7. X = era_900_1100[['expected_recovery_amount', 'indicator_1000']]
8. y = era_900_1100['actual_recovery_amount']
9. X = sm.add_constant(X)
10.
11. # Build linear regression model
12. model = sm.OLS(y,X).fit()
13.
14. # Print the model summary
15. model.summary()
```

The results were:

```

OLS Regression Results
=====
Dep. Variable:    actual_recovery_amount    R-squared:            0.314
Model:            OLS                      Adj. R-squared:       0.307
Method:           Least Squares           F-statistic:          41.22
Date:             Wed, 20 Mar 2019         Prob (F-statistic):    1.83e-15
Time:             12:47:53                 Log-Likelihood:        -1272.0
No. Observations: 183                     AIC:                   2550.
Df Residuals:     180                     BIC:                   2560.
Df Model:         2
Covariance Type:  nonrobust
=====
               coef      std err          t      P>|t|      [0.025      0.975]
-----
const                3.3440      626.274      0.005      0.996    -1232.440     1239.128
expected_recovery_amount  0.6430       0.655      0.981      0.328      -0.650       1.936
indicator_1000        277.6344      74.043      3.750      0.000      131.530     423.739
=====
Omnibus:            65.977    Durbin-Watson:          1.906
Prob(Omnibus):      0.000    Jarque-Bera (JB):        186.537
Skew:               1.510    Prob(JB):                 3.12e-41
Kurtosis:           6.917    Cond. No.                  3.37e+04
=====

Warnings:
[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.
[2] The condition number is large, 3.37e+04. This might indicate that there are
strong multicollinearity or other numerical problems.

```

9. Regression modeling: adjusting the window

The regression coefficient for the true threshold was statistically significant with an estimated impact of around \$278 and a 95 percent confidence interval of \$132 to \$424. This is much larger than the incremental cost of running the higher recovery strategy which was \$50 per customer. At this point, we are feeling reasonably confident that the higher recovery strategy is worth the additional costs of the program for customers just above and just below the threshold.

Before showing this to our managers, we want to convince ourselves that this result wasn't due just to us choosing a window of \$900 to \$1100 for the expected recovery amount. If the higher recovery strategy really had an impact of an extra few hundred dollars, then we should see a similar regression coefficient if we choose a slightly bigger or a slightly smaller window for the expected recovery amount. Let's repeat this analysis for the window of expected recovery amount from \$950 to \$1050 to see if we get similar results.

```

1. # Redefine era_950_1050 so the indicator variable is included
2. era_950_1050 = df.loc[(df['expected_recovery_amount']<1050) &

```

```

3.             (df['expected_recovery_amount']>=950)]
4.
5. # Define X and y
6. X = era_950_1050[['expected_recovery_amount','indicator_1000']]
7. y = era_950_1050['actual_recovery_amount']
8. X = sm.add_constant(X)
9.
10. # Build linear regression model
11. model = sm.OLS(y,X).fit()
12.
13. # Print the model summary
14. model.summary()

```

The results were:

```

=====
                        OLS Regression Results
=====
Dep. Variable:      actual_recovery_amount      R-squared:                0.283
Model:              OLS                        Adj. R-squared:           0.269
Method:             Least Squares              F-statistic:             18.99
Date:               Wed, 20 Mar 2019            Prob (F-statistic):      1.12e-07
Time:               12:50:18                    Log-Likelihood:          -692.92
No. Observations:   99                        AIC:                     1392.
Df Residuals:       96                        BIC:                     1400.
Df Model:           2
Covariance Type:    nonrobust
=====

```

	coef	std err	t	P> t	[0.025	0.975]
const	-279.5243	1840.707	-0.152	0.880	-3933.298	3374.250
expected_recovery_amount	0.9189	1.886	0.487	0.627	-2.825	4.663
indicator_1000	286.5337	111.352	2.573	0.012	65.502	507.566

```

=====
Omnibus:            39.302      Durbin-Watson:           1.955
Prob(Omnibus):      0.000      Jarque-Bera (JB):         82.258
Skew:               1.564      Prob(JB):                 1.37e-18
Kurtosis:           6.186      Cond. No.                  6.81e+04
=====

```

Warnings:

[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.

[2] The condition number is large, 6.81e+04. This might indicate that there are strong multicollinearity or other numerical problems.

10. Conclusion

Whether we use a wide window (\$900 to \$1100) or a narrower window (\$950 to \$1050), the incremental recovery amount at the higher recovery strategy is much greater than the \$50 per customer it costs for the higher recovery strategy. So, we can say that the higher recovery strategy is worth the extra \$50 per customer that the bank is spending.