# Project Home-a-loan

• • •

Wei Cheung
July 19, 2017

# The Challenge

Part 1: Lead Conversion

Part 2: Loan Processing

Lead →————————————————→ Lock ————————————————→ Fund

**Mission:**
- Predict whether a lead will convert to "lock"
- Supervised learning (Classification)

**Business Values:**
- Know the potential of each customer - know the ones to focus on
- Best Customers profiling - targeted marketing for new customer acquisition
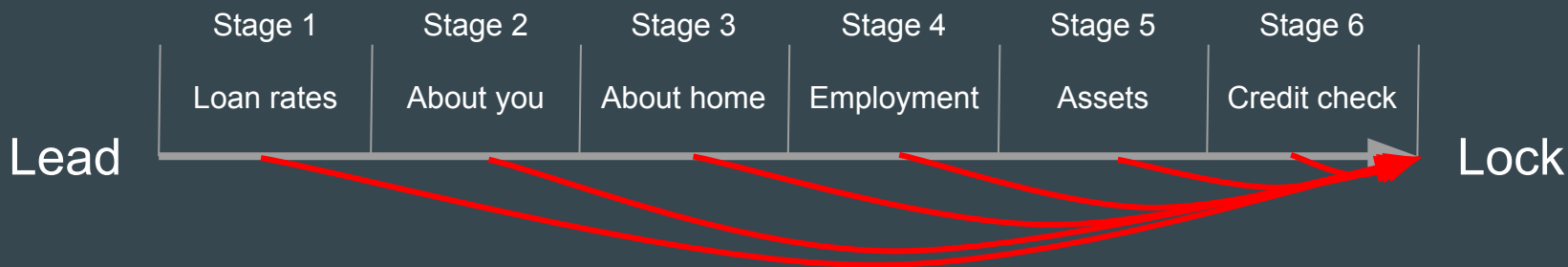
**Mission:**
- Predict "locked-to-funded" time (efficiency)
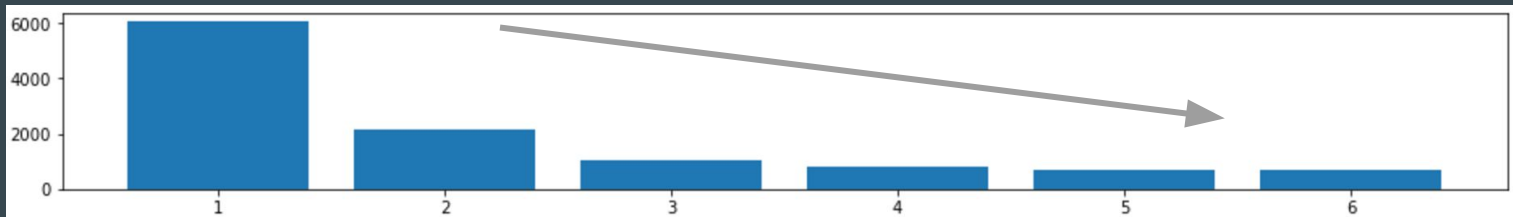- Supervised learning (Regression)

**Business Values:**
- Improve customer experience by providing an expected waiting time
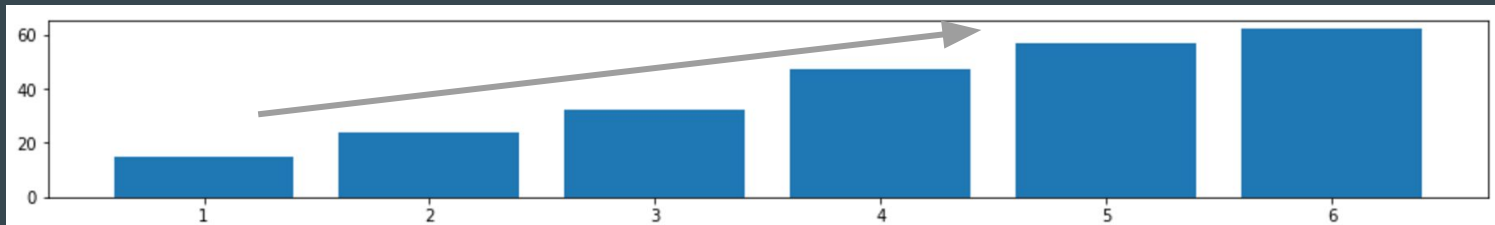- Know areas of improvement for efficiency

# The data

- 9752 cases
- 417 features
- Cleaning / pre-processing:

| Data Type | Examples | Processing |
|---|---|---|
| Numerical | Borrower income, loan amount | None |
| Categorical | Type of home, Education level, City of property, Gender | Dumification |
| Text | Goal of refinancing, Unqualified reason note | Tfidf Vectorization (limiting stop words) |
| Datetime | Created time, Last modified | Categorize (year, quarter, month, dow), Calculate Period (difference of dates/times), Calculate Cohort (quarter/month since initiation) |

# Part 1 - Leads Conversion

Part 1: Lead Conversion

Part 2: Loan Processing

Lead ⟶ Lock ⟶ Fund

**Mission:**
- Predict whether a lead will convert to "lock"
- Supervised learning (Classification)

**Business Values:**
- Know the potential of each customer - know the ones to focus on
- Best Customers profiling - targeted marketing for new customer acquisition

**Mission:**
- Predict "locked-to-funded" time (efficiency)
- Supervised learning (Regression)

**Business Values:**
- Improve customer experience by providing an expected waiting time
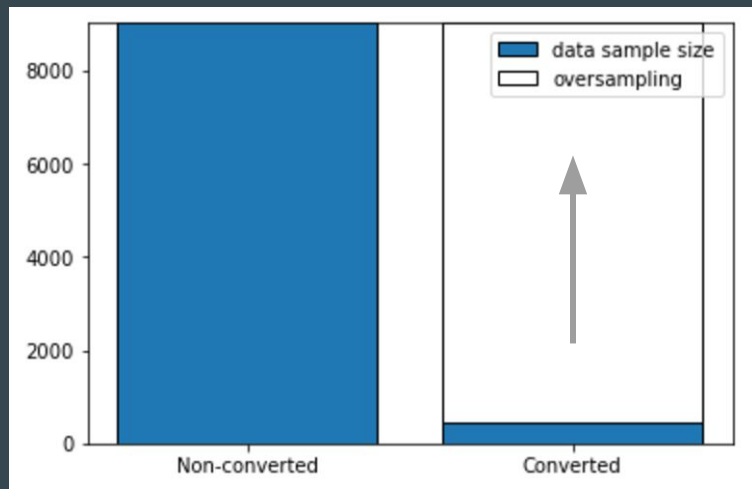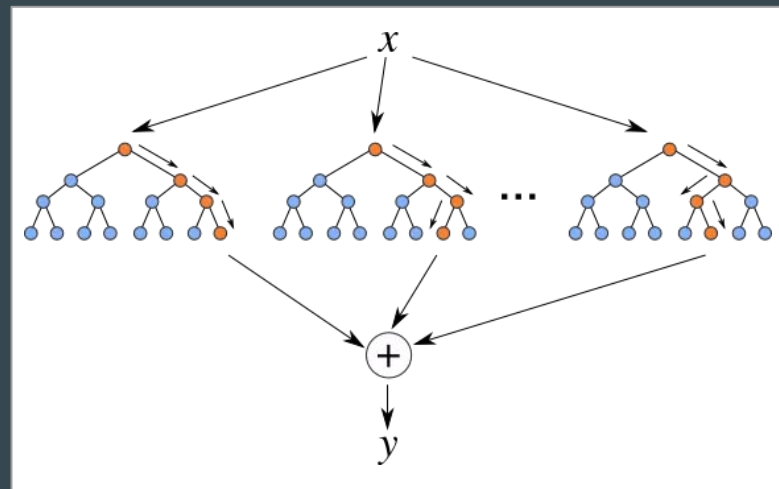- Know areas of improvement for efficiency

# Part 1 - Leads Conversion

# Part 1 - Leads Conversion - Methodology

Random Oversampling for Imbalance Classes

Random Forest Classifier

# Part 1 - Leads Conversion - Results
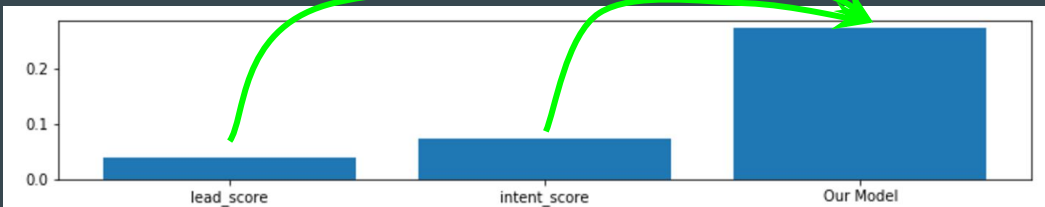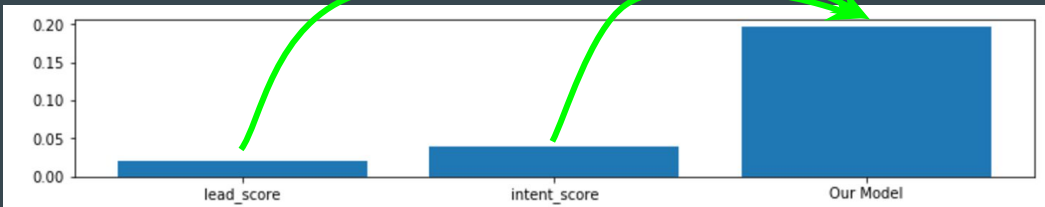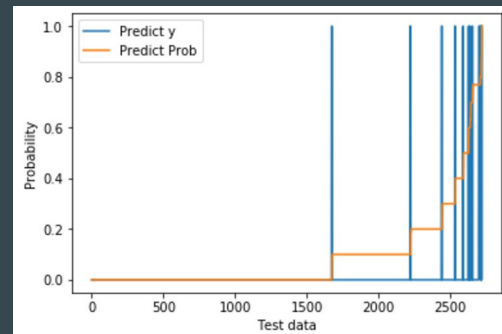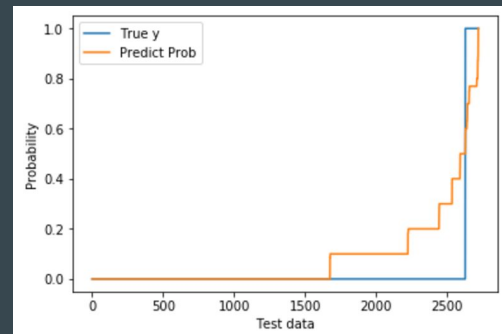
Predicted Probability v.s. True Value

# Part 1 - Leads Conversion - Interpretation



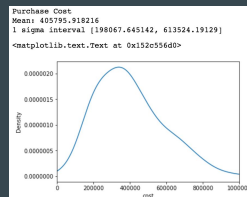| Feature | Borrower Age | Borrower's Annual Income | Borrower current employment | Home purchase cost | Years in home | Property zip code |
|---|---|---|---|---|---|---|
| Mean | 44 | 127,000 | 7.5 | 406,000 | 8 | 98103 |
| Common range | 33 - 54 | 6,300 - 247,000 | 0 - 15 | 190,000 - 614,000 | 0.5 - 15 | 98103, 92691, 98125, 93003 |

# Part 2 - Loan Processing

Part 1: Lead Conversion

Part 2: Loan Processing

Lead ———————————————▶ Lock ———————————————▶ Fund

Mission:
- Predict whether a lead will convert to "lock"
- Supervised learning (Classification)

Mission:
- Predict "locked-to-funded" time (efficiency)
- Supervised learning (Regression)

Business Values:
- Know the potential of each customer - know the ones to focus on
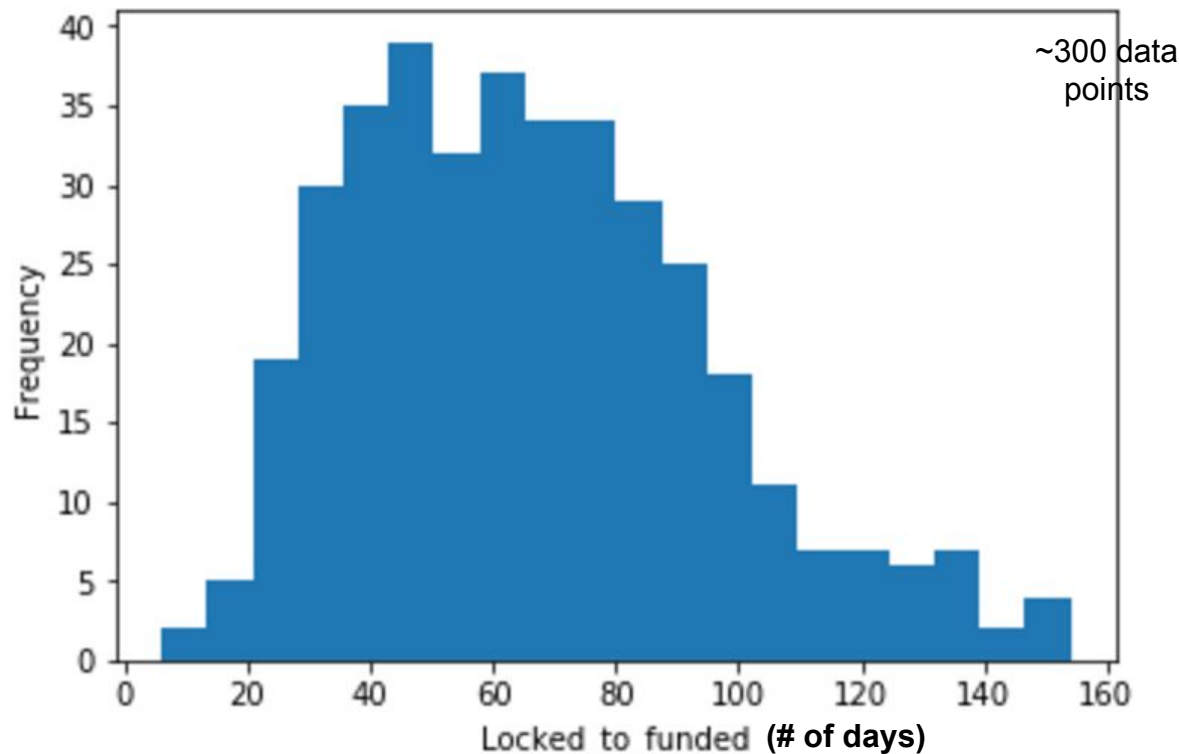- Best Customers profiling - targeted marketing for new customer acquisition

Business Values:
- Improve customer experience by providing an expected waiting time
- Know areas of improvement for efficiency

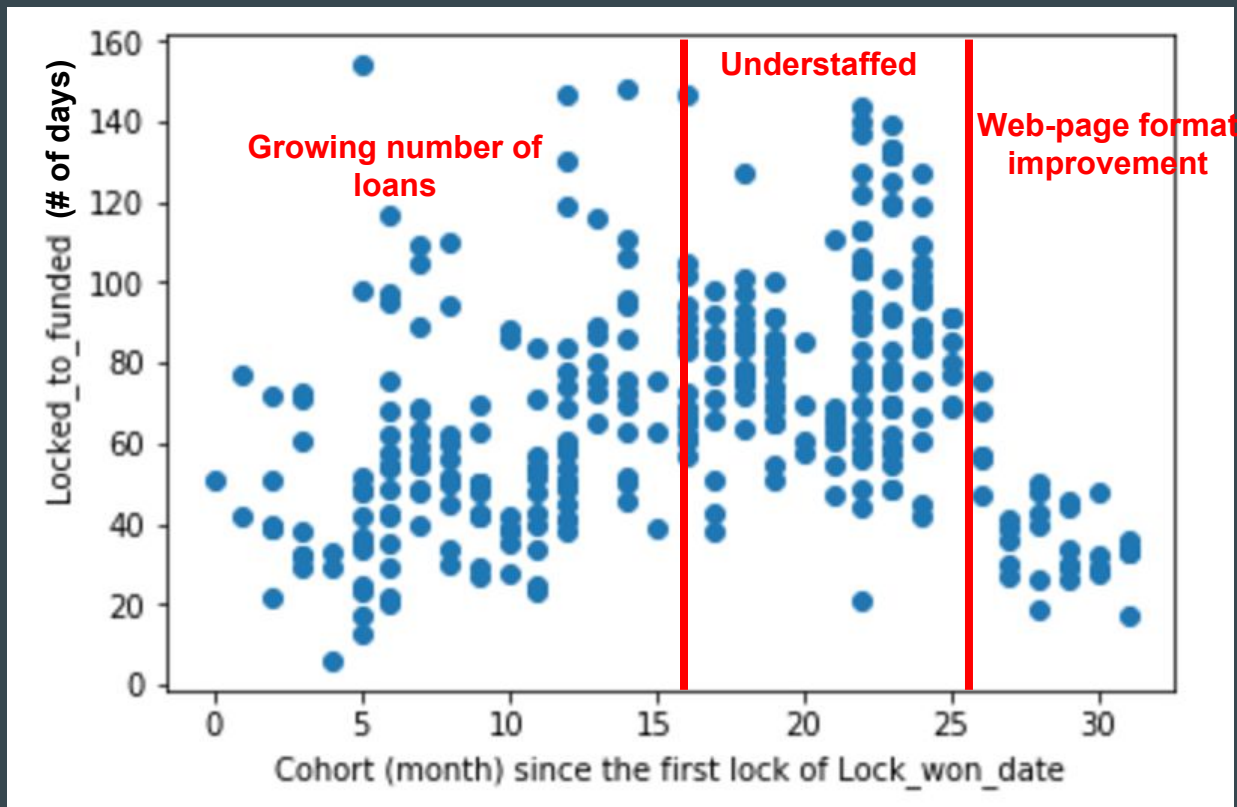# Part 2 - Loan Processing



~300 data points

**Outcome:**
- **Experimented with different features, different algorithms (Linear, L1/L2, KNN, Decision tree, Bagging, RF, Grad/AdaBoosting)**
- **No strong signal (near zero R2)**

**Potential reason:**
- **Data set too small**
- **Signals lie in external factors**

# Part 2 - Loan Processing



(the higher, the worse)

(or, months since the company launched)

# Conclusions

Part 1
- Multistage model construction, Oversampling, Random Forest Classification
- Outperform existing models in Accuracy, F1-score, Precision
- Import feature extracted, Best Customer Profiling

Part 2
- Too few data points (~300) to detect signal for prediction
- Time cohort analysis indicated trends in outcome

Future work
- Implement prediction pipeline into existing workflow
- Include new features (from internal and external sources) into prediction models
- Try different models (neural network) that can take care of wide data sets

# Thank you

Wei Cheung

cheungwz@gmail.com

github.com/weizhic