



UNIVERSITY OF EDINBURGH  
Business School



## Credit scoring with alternative data

Jonathan Crook, Viani Djeundje, Raffaella Calabrese, Mona Hmid  
Credit Research Centre, University of Edinburgh

Credit Scoring & Credit Control XVI Conference August 2019



# Incidence of Unscorable People in US Federal Protection Bureau's CC Panel 2010

Brevoort, Grimm & Kambara (2016)

1-in-48 random sample from a NCRA

188.7m have credit records that can be scored by a commonly used credit scoring model.

Col	Population	No (m)	% of US Pop	Share with Performance	Delinquency Rate
1	Can be Scored	188.7	77.9	89.5	12.2
2	• Thick file	180.7		90.2	12.4
3	• Thin file	8.1	3.3	72.9	6.3
4	• Stale	9.6	4.0	12.3	26.0
5	• Insufficient	9.9	4.1	21.8	22.4
6=4+5	Cannot be scored	19.4	8.1		
7	<b>No credit records</b>	<b>26</b>	<b>10.7</b>		



# Incidence of Unscorable People in US: FICO Data 2015

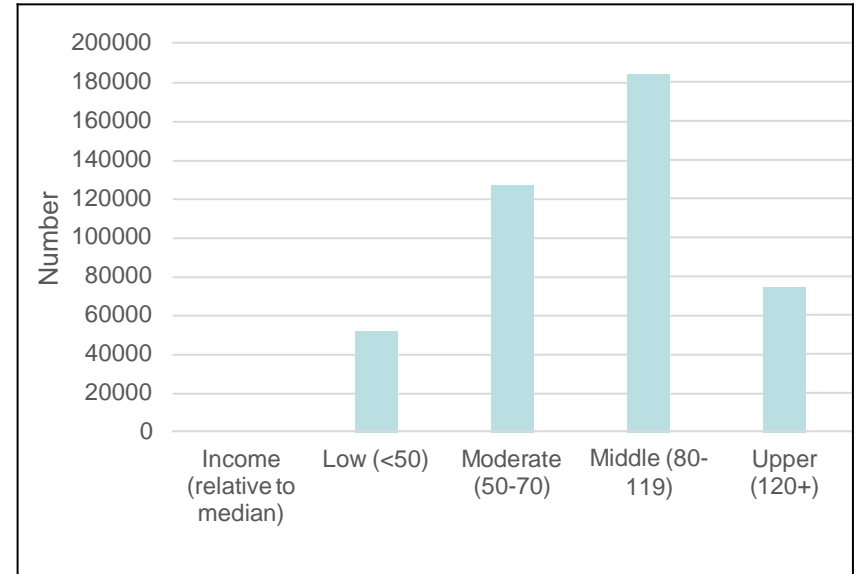
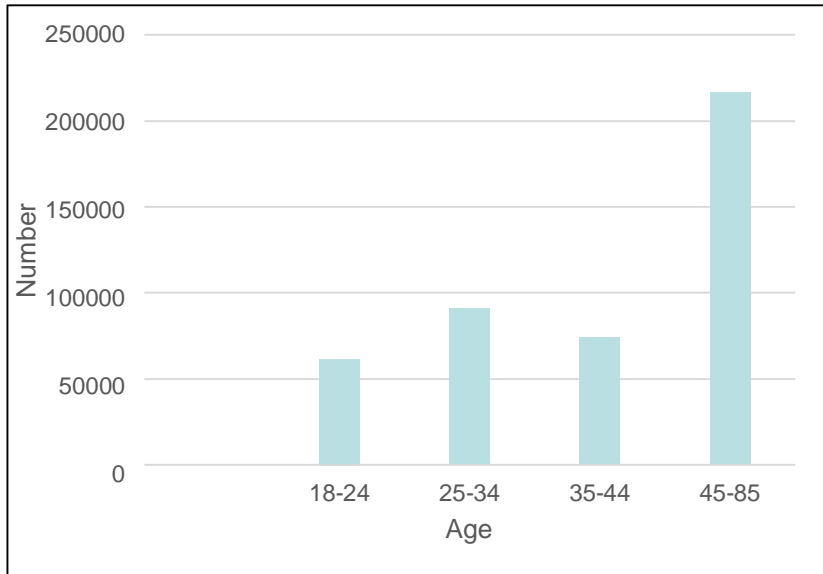
Segment		Millions
Have credit records but insufficient for scoring	No tradeline performance reported in last 6 months	3
	Files with delinquent tradelines, collections or adverse public records	18
	Stale files	7
No credit records		25
<b>Total</b>		<b>53</b>

Source: A. Jennings (2015) *Expanding the credit eligible population in the USA: A Case Study*. Credit Scoring and Credit Control XIV conference, Edinburgh.



## Federal Reserve Bank of New York Consumer Credit Panel/Equifax

Sample: 15,509 individuals who had no credit score due to thin files or no credit file  
[ lack of trade lines to construct score or no credit history] Q1 2007



Source: Smith, M. & Henderson, C. (2017) Beyond thin files. *Social Science Quarterly*, 1-19.



1.7billion estimated number of adults without an account at a financial institution nor facility through a mobile phone provider (31% of adult population) Demirguc-Kunt et al (2017)

Concentrated in lower income countries: China 204m, India 357m Indonesia 102m

A major difficulty is lack of past credit history.

Industry response: use alternative data.

## Examples

- Rental data in UK (Experian)
- Data on utility payments, evictions property values etc (FICO)
- Mobile phone data (CredoLab, Lenddo)
- Psychometrics (Lenddo, )



# Contributions

## Two contributions

1. To show that by using data on email usage and on psychometrics one can gain good separation between good and bad payers.
2. To show the relative contributions to predictive accuracy of these characteristics compared with demographic variables.

## Textual predictors

Author	Predictors	Enhancement in predictive accuracy	Data
Dorfleitener et al (2016)	Spelling errors, length of text, occurrence of keywords in loan purpose statements		
Goa (2018)	Readability, tone, occurrence of deception cues in loan purpose statements		Prosper P2P
Netzer et al (2018)	Two word combinations, types of word groupings	Increase AUC by 2.6%	Prosper P2P
Iver et al (2016)	Soft information (whether posted a pic, number of words used in listing)	Increased AUC from 0.710 to 0.714	Prosper P2P
Berg et al (2018)	Digital foot print variables: use of lower case, email address errors , os of device etc	Increased AUC for scoreables: 0.68 to 0.73 Unscoreables gained 0.68.	E-commerce furniture seller



# Mobile Phone predictors

Bjorkegren & Grissen (2018)	Periodicity of use, fraction of spoken time during the day, variation in usage, autocorrelation between calls and SMS messages	Phone indicators gave higher AUC than bureau alone.	EFL

# Psychometric variables

Klinger et al (2013)	66 undefined variables	Used a lone: AUC 0.7	Micro & SMEs, Peru
Arraiz et al (2017)	undefined	Those accepted under traditional predictors and rejected on Psyc model had higher default rate than those accepted on a traditional model.	Banked entrepreneurs,
Dlogosch et al (2017)	Measures of conscientiousness, emotional stability, openness to experience and integrity.	N/a	Micro-entrepreneurs, Kenya
Liberati & Camillo (2018)	Six factors from a PCA applied to responses to Semiometrie	When added to models containing cash flow and solvency score, AUC increased 0.554 to 0.850.	



# Data

## Data set A

	Demographic	Psychometric	Email
# variables	12	350	53
# cases	1826	1826	33,091

## Data set B

	Demographic	Psychometric	Email
# variables	n/a	n/a	237
# cases			16,358

# Two stage estimation

- Logistic regression
- Stepwise using p-value or AIC
- One for each of:  
demographics, psychometrics, email

$$\hat{\beta}_j \quad j=1, 2, 3$$



Ensemble of separate models

$$\hat{z} = \hat{\beta}_j$$

Logistic regression



**Table 2**  
**Estimated coefficients for the submodel based on demographic data alone**

<i>Variable</i>	<i>Coefficient</i>	<i>p-value</i>
<i>Intercept</i>	-1.6778	0.002
<i>How long has had phone</i>	0.0198	0.418
<i>Number of dependents = 2 (coded 1, 0 otherwise)</i>	-0.4387	0.010
<i>Number of dependents = 6 (coded 1, 0 otherwise)</i>	-0.978	0.606
<i>Hours worked per week</i>	0.0137	0.030
<i>Work experience</i>	-0.0117	0.689
<i>Age in years</i>	0.0079	0.555
<i>Gender (male=1)</i>	-1.9801	0.001
<i>Income</i>	-0.0001	0.287
<i>Age * gender</i>	0.0473	0.005
<i>How long has had phone * work experience</i>	-0.0036	0.099





**Table 4**

**Estimated coefficients for the sub-model based on alternative data alone**

<i>Variable</i>	<i>Coefficient</i>	<i>p-value</i>
<i>Intercept</i>	0.1687	0.705
<i>Time in years to send last 2000 emails</i>	0.6201	0.010
<i>Number of contacts the applicant sent the last 2000 emails to</i>	-0.0054	0.027
<i>Average number of words the applicant used in the subject line of the last 2000 emails</i>	-0.1434	0.108
<i>Fraction of emails sent between 0000hrs and 0600hrs</i>	1.7151	0.004
<i>Fraction of emails sent between 1800hrs and 2400 hrs</i>	1.4781	0.164
<i>Fraction of emails that were sent on Tuesdays</i>	-1.6544	0.048
<i>Fraction of emails that were sent on Thursdays</i>	-2.9411	0.006
<i>Fraction of emails that were sent on Saturdays</i>	-2.6813	0.015
<i>Fraction of emails that were sent on Sundays</i>	-3.6693	0.039
<i>Fraction of emails that were sent to or received from non-top financial product providers</i>	0.7980	0.087
<i>Log of number of emails received from uber.com</i>	23.8613	0.139
<i>Log of number of emails received from uber</i>	24.0732	0.135



# Predictive Performance

## Stage 1

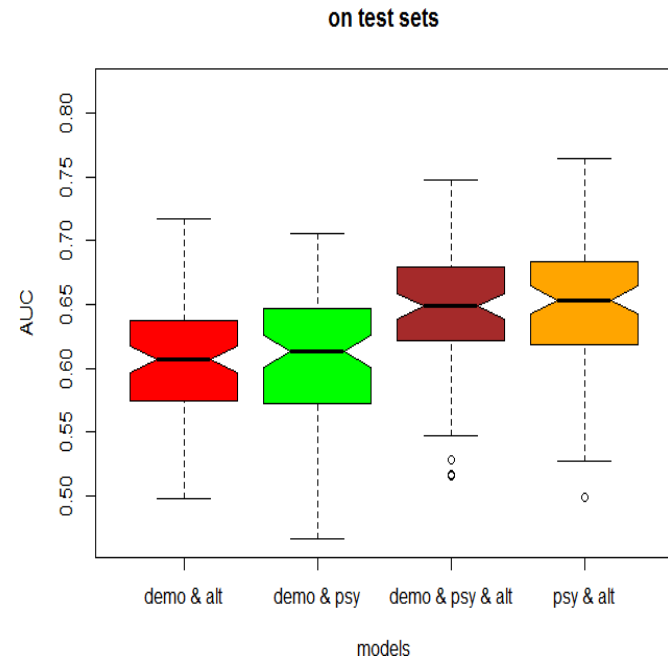
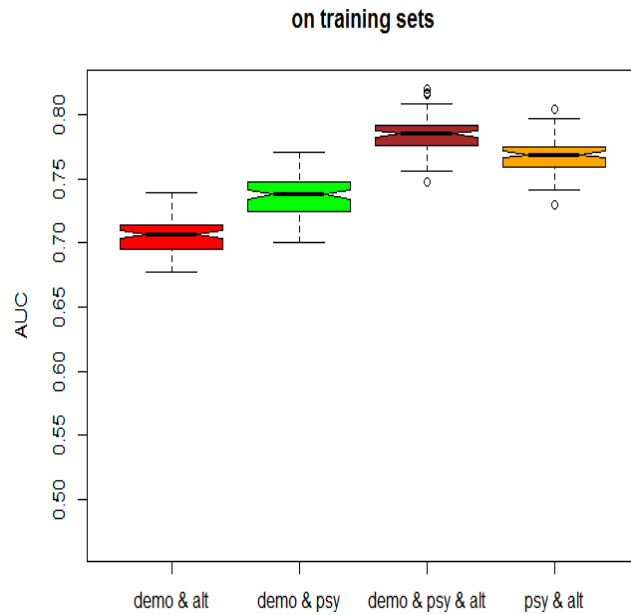
	Demographic model	Psychometric model	Email data model
AUC Training	0.63	0.67	0.67
<b>AUC Test</b>	<b>0.62</b>	<b>0.60</b>	<b>0.58</b>
# training cases	1370	1370	332
# parameters	11	23	13

## Stage 2

	Demographic + Psychometric	Demographic + Email	Psychometric + Email	Demographic + Psychometric + Email
AUC Training	0.66	0.71	0.71	0.73
<b>AUC Test</b>	<b>0.75</b>	<b>0.69</b>	<b>0.67</b>	<b>0.72</b>

# Sensitivity with respect to train/test split

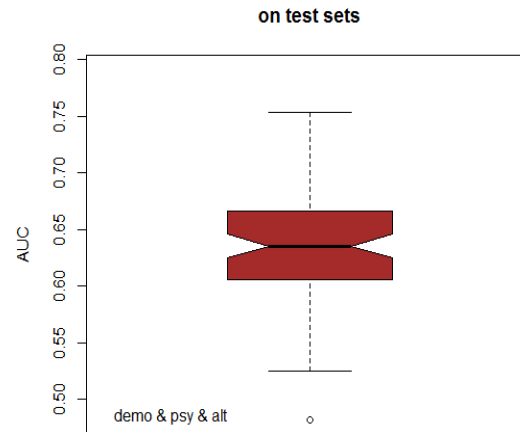
Results from 100 different random samples of 75% training 25% test.



# Imputation

To increase sample size we pooled all 3 datasets and imputed missing values.

- Used MICE (multiple imputations by chained equations) Royston & White 2002. 20 imputates.
- Estimate for each imputate and average using formula from Little & Rubin (2002).
- Estimate separate models for each of demographic variables, psychometric and email datasets.
- Ensemble.



Results similar to the previous analysis.



# Analysis of Data set B

	AUC train	AUC test
Logistic regression	56%	54%
LASSO	64%	57%
Ridge regression	67%	57%
<b>Extreme Gradient Boosting</b>	<b>68%</b>	<b>62%</b>
Oversampling1 + Extreme Gradient Boosting	82%	61%
<b>Oversampling2 + Extreme Gradient Boosting</b>	<b>99%</b>	<b>62%</b>
Neural Networks	97%	60%
PCA + Logistic regression	60%	54%
PCA + LASSO	60%	52%
PCA + Ridge Regression	60%	57%
<b>PCA + Extreme Gradient Boosting</b>	<b>69%</b>	<b>62%</b>
PCA + Neural Networks	73%	57%



# Discussion

- Predictive accuracy from ensemble model greater than that from individual models.
- AUC gained from our models compares well to others in the literature that use alternative data:

Our models (demographic, psychometric, email)	0.75
Berg (digital footprint)	0.73
Dlogisch (psychometric)	0.67
Liberati & Camillo (psychometric)	0.85

## Implications

- Use of alternative data is increasing in lower income countries
- Can gain commercially useful predictive accuracy using psychometric and/or email variables
- Can be useful substitute for lack of credit performance data
- But applicants must be willing to complete a psychometric test before applying and possible gaming