

# Using Machine Learning to Determine Who Should Get a Stimulus Check

---

By Scott L. A. Johnson

# Summary

Explored modeling UCI Adult Data Set based on the 1994 Census using  
Classification models: KNN, Decision Tree, Random Forest, Logistic Regression

Compared performance of models: KNN performed the best

Model citizen who will benefit most from stimulus check: single, younger than 36,  
didn't finish high school

# Outline

- Business Problem
- Data
- Methods
- Results
- Conclusions

# Business Problem

**Assumption\*:** citizens within/under defined income threshold will most benefit from stimulus check

- Want a way to, in the absence a person's immediate income, identify those citizens above\* **[develop PERSON PROFILE based on weighted characteristics]**
- Machine learning algorithm can help us weight the most important characteristics associated with a person's income for use under the following example circumstances:
  - Don't have access to recent tax return
  - Citizen's immediate job/financial circumstances have changed
  - Specific industries/communities have been affected by pandemic/natural disaster

# Business Problem: Questions

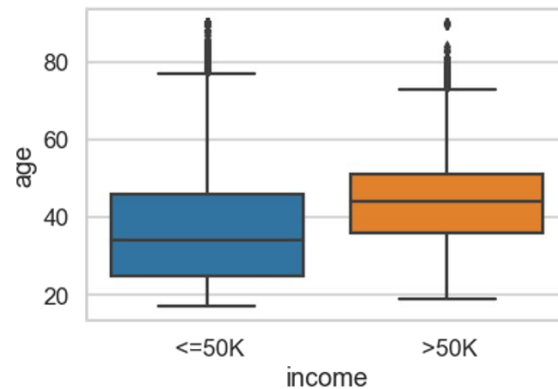
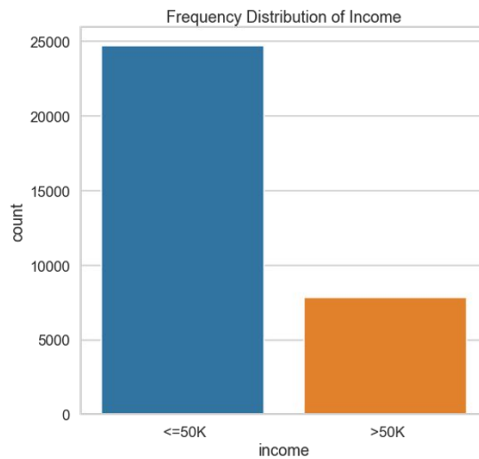
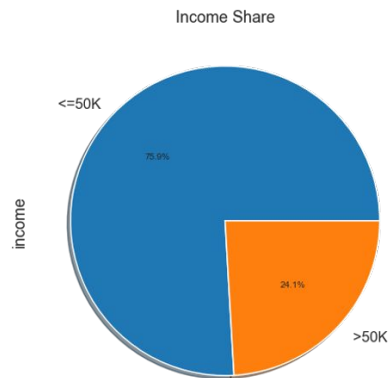
- What factors contribute most to telling us who makes above 50k, least at risk without a stimulus check? What types of people need the most support based on?
- What is the highest precision we can get to minimize the false positives, where we assume a family can make it without a check but can't?
- How do we minimize the overall cost / impact of the pandemic?
- Is there an industry focus we should have? What industries should we focus on?

# Business Problem: Consequences

- False Negative: wrong about who makes under 50k
  - Send check [less resources to allocate to others and support future programs]
  - Citizens may use that check frivolously, but contribute back to the economy
- False Positive:** wrong about who makes over over 50k
  - No check puts citizens at risk of homeless (cost more), lost job
  - We have more resources to reallocate, but the time lost doesn't help you pay your bills
  - Enact policies to alleviate people who we've missed

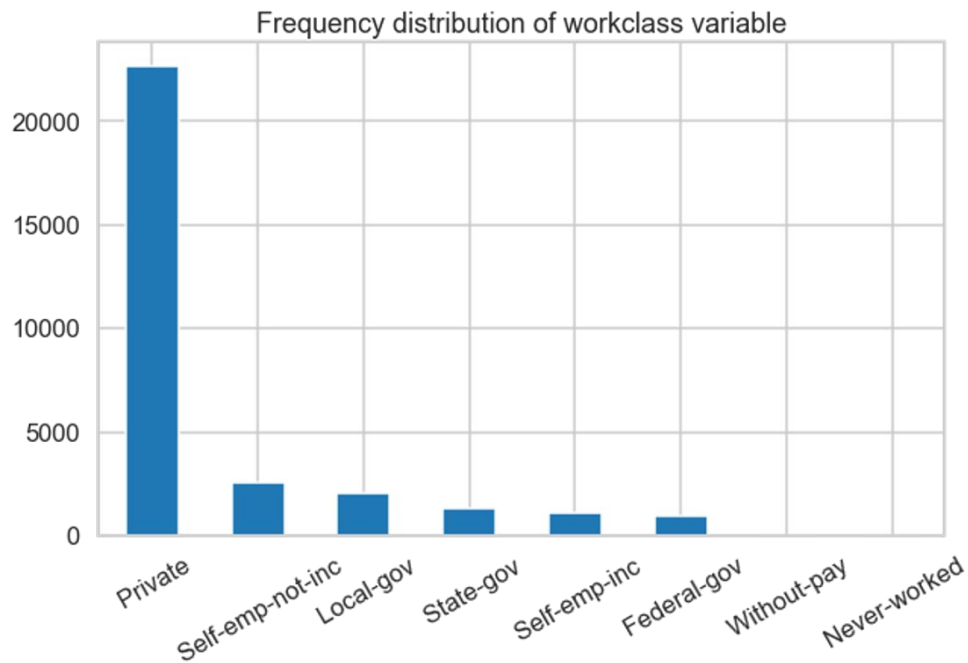
Choose the model with the best F1 score

# Data → ExploDA



The percentage of the population in dire need of the check  
coincides with age

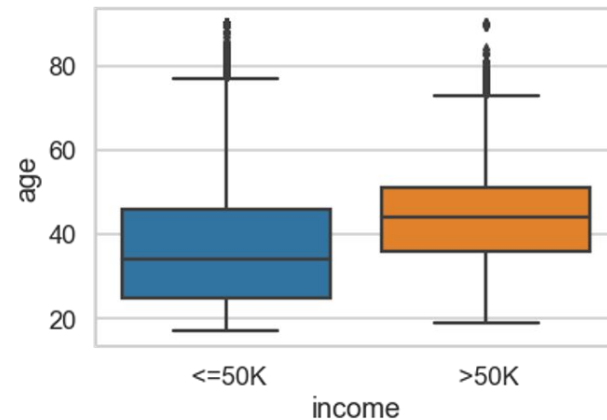
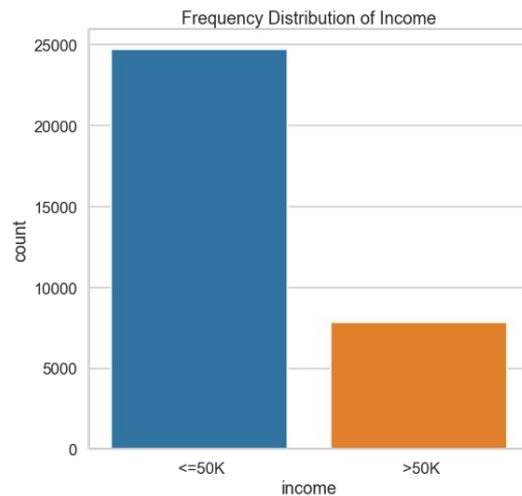
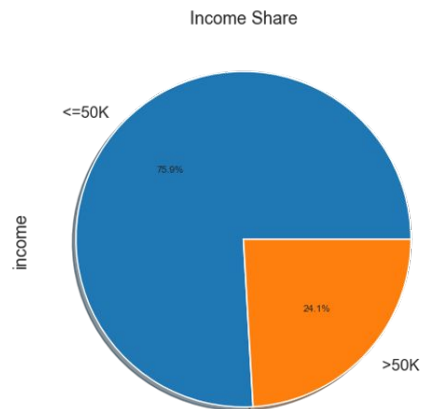
# Data → EDA



A focus on private industry gives the highest opportunity to address the 75% of the population that makes less than 50K



# Data



The percentage of the population in dire need of the check coincides with age

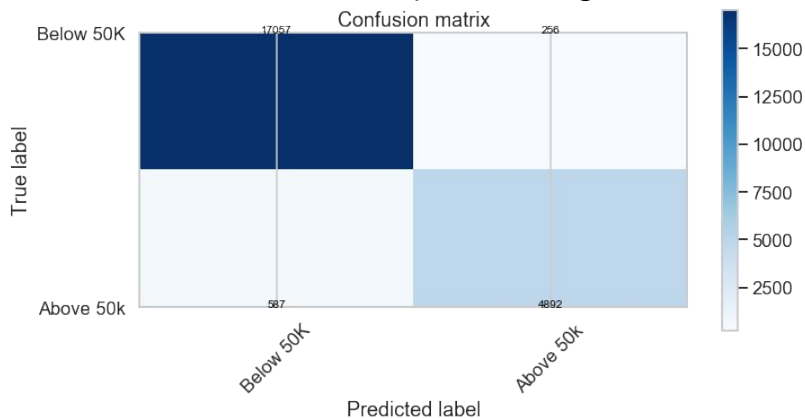
# Methods

- Data Cleaning Preprocess
  - Dropped fnlwgt, race, country, sex, education from model
  - Checked for missing data and place holders
- Feature engineering for categorical values
- Shotgun approach to choose a couple models and compare their precision
  - Used SMOTE (**Synthetic Minority Oversampling Technique**) due to the severe class imbalance to improve model performance
- Use feature scoring from models to determine variables with the highest influence on classification

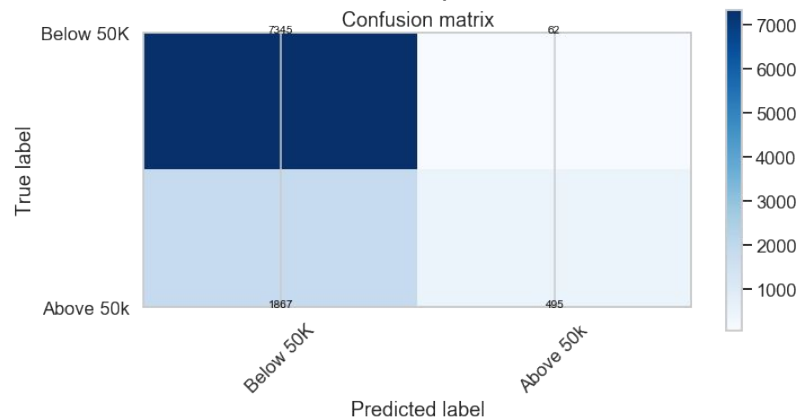
# Results (Training Data / Test Data)

Model	Precision	Accuracy	Recall	F1 Score	ROC_AUC
Random Forest <b>70cv_score:</b> 0.70 / 0.71	<b>70:</b> 0.95 / 0.89	<b>70:</b> 0.96 / 0.80	<b>70:</b> 0.89 / 0.21	<b>70:</b> 0.92 / 0.34	<b>70:</b> 0.94 / 0.60

70/30 Train-Test Split: Training Data



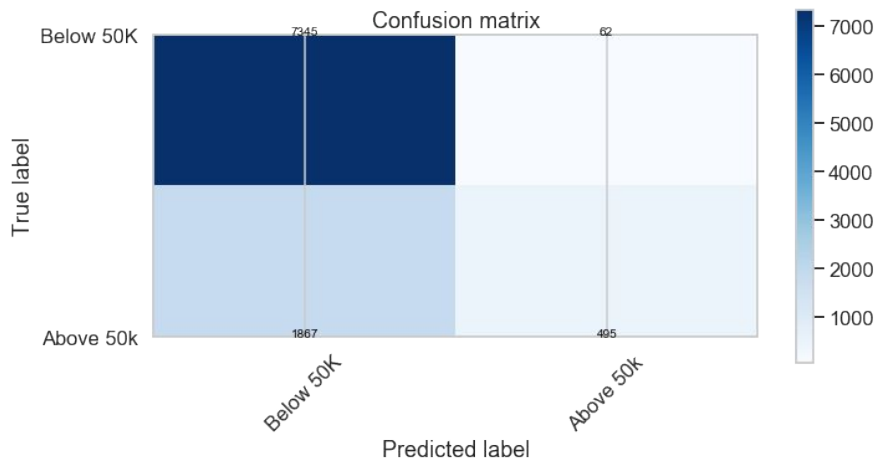
70/30 Train-Test Split: Test Data



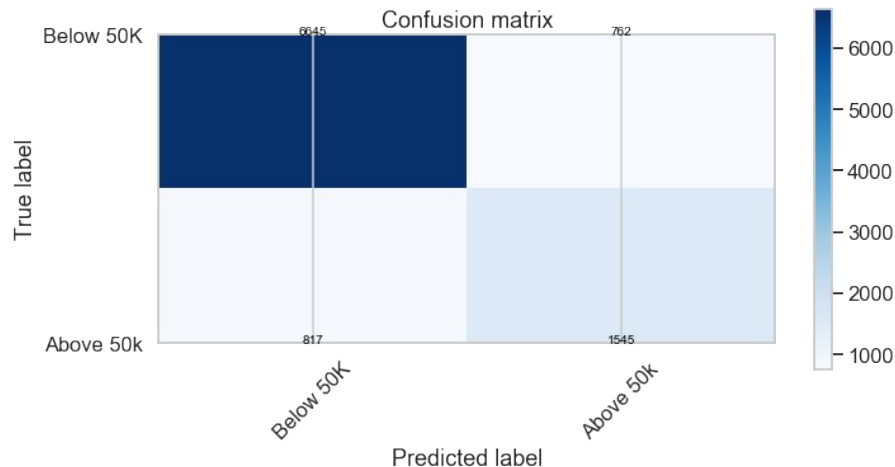
# Results (Test Data)

Model	Precision	Accuracy	Recall	F1 Score	ROC_AUC
Random Forest	<b>70:</b> 0.89 <b>70_smt:</b> 0.67	<b>70:</b> 0.80 <b>70_smt:</b> 0.84	<b>70:</b> 0.21 <b>70_smt:</b> 0.65	<b>70:</b> 0.34 <b>70_smt:</b> 0.66	<b>70:</b> 0.60 <b>70_smt:</b> 0.78

70/30 Train-Test Split: Test



70/30 Train-Test Split: Test w/ SMOTE



# Results (Test Data)

Model	Precision	Accuracy	Recall	F1 Score	ROC_AUC
Random Forest (cv_score =0.70)	<b>70:</b> 0.93 <b>70_smt:</b> 0.67 <b>80:</b> 0.92 <b>80_smt:</b> 0.66	<b>70:</b> 0.79 <b>70_smt:</b> 0.84 <b>80:</b> 0.79 <b>80_smt:</b> 0.84	<b>70:</b> 0.16 <b>70_smt:</b> 0.65 <b>80:</b> 0.18 <b>80_smt:</b> 0.66	<b>70:</b> 0.27 <b>70_smt:</b> 0.66 <b>80:</b> 0.30 <b>80_smt:</b> 0.66	<b>70:</b> 0.57 <b>70_smt:</b> 0.78 <b>80:</b> 0.59 <b>80_smt:</b> 0.78
KNN [15] (cv_score =0.69)	<b>70:</b> 0.79 <b>70_smt:</b> 0.67 <b>80:</b> 0.77 <b>80_smt:</b> 0.66	<b>70:</b> 0.86 <b>70_smt:</b> 0.85 <b>80:</b> 0.85 <b>80_smt:</b> 0.84	<b>70:</b> 0.55 <b>70_smt:</b> 0.72 <b>80:</b> 0.55 <b>80_smt:</b> 0.73	<b>70:</b> 0.65 <b>70_smt:</b> 0.70 <b>80:</b> 0.65 <b>80_smt:</b> 0.70	<b>70:</b> 0.75 <b>70_smt:</b> 0.80 <b>80:</b> 0.75 <b>80_smt:</b> 0.80
Decision Tree [11] (cv_score =0.71)	<b>70_smt:</b> 0.58 <b>80_smt:</b> 0.58	<b>70_smt:</b> 0.81 <b>80_smt:</b> 0.81	<b>70_smt:</b> 0.78 <b>80_smt:</b> 0.79	<b>70_smt:</b> 0.67 <b>80_smt:</b> 0.67	<b>70_smt:</b> 0.80 <b>80_smt:</b> 0.80
LogReg (cv_score = 0.73)	<b>70_smt:</b> 0.65 <b>80_smt:</b> 0.64	<b>70_smt:</b> 0.84 <b>80_smt:</b> 0.83	<b>70_smt:</b> 0.71 <b>80_smt:</b> 0.71	<b>70_smt:</b> 0.68 <b>80_smt:</b> 0.67	<b>70_smt:</b> 0.80 <b>80_smt:</b> 0.58

# Conclusions

Our best performing model was the KNN model with a CV score of .70 , a F1 score of .70,

Candidate income profile heavily influenced from our random forest model:

- Age → focus on citizens < 36
- Marital Status → focus on single
- Education → focus on education < 10th grade (those who don't have a high school diploma)

\*Workclass/occupation did not rank high for the 1994 data set as an income determinant

- Does not provide accurate representation of factors that would be relevant to today's job market

# Future Recommendations

- Feed next rounds of data into model to train, collect information relevant to industries of relevant today -->clearer
- Bin wage classes differently to assign different stimulus amounts based on wage classes
- Alternate analysis can also be done with respect to wage class → different factors might distinguish those in the 150k and up class vs 75k-90k → zipcode could segment the data → identify risk of being able to weather pandemic (vaccine selection)
- Though not done in this study might be interesting to explore wage disparities by community / culture → if funds are dispense at a state level, could prioritize support based on tiered approach [zipcode]

# Thank You!

**Email:** [scottlajohnson92@gmail.com](mailto:scottlajohnson92@gmail.com)

**GitHub:** [@scottlou](https://github.com/scottlou)

**LinkedIn:** [www.linkedin.com/in/scott-johnson-432b83163/](https://www.linkedin.com/in/scott-johnson-432b83163/)