

Predicting new locations for Whole Foods

GA Data Science Project

Markus Huber

Data

- IRS data 2012:
 - 73 features
 - 27718 zip codes
- Whole Foods location:
383 locations

Challenges

- Few number of positive labels and no negative labels
- Different feature scales
- “Prediction”

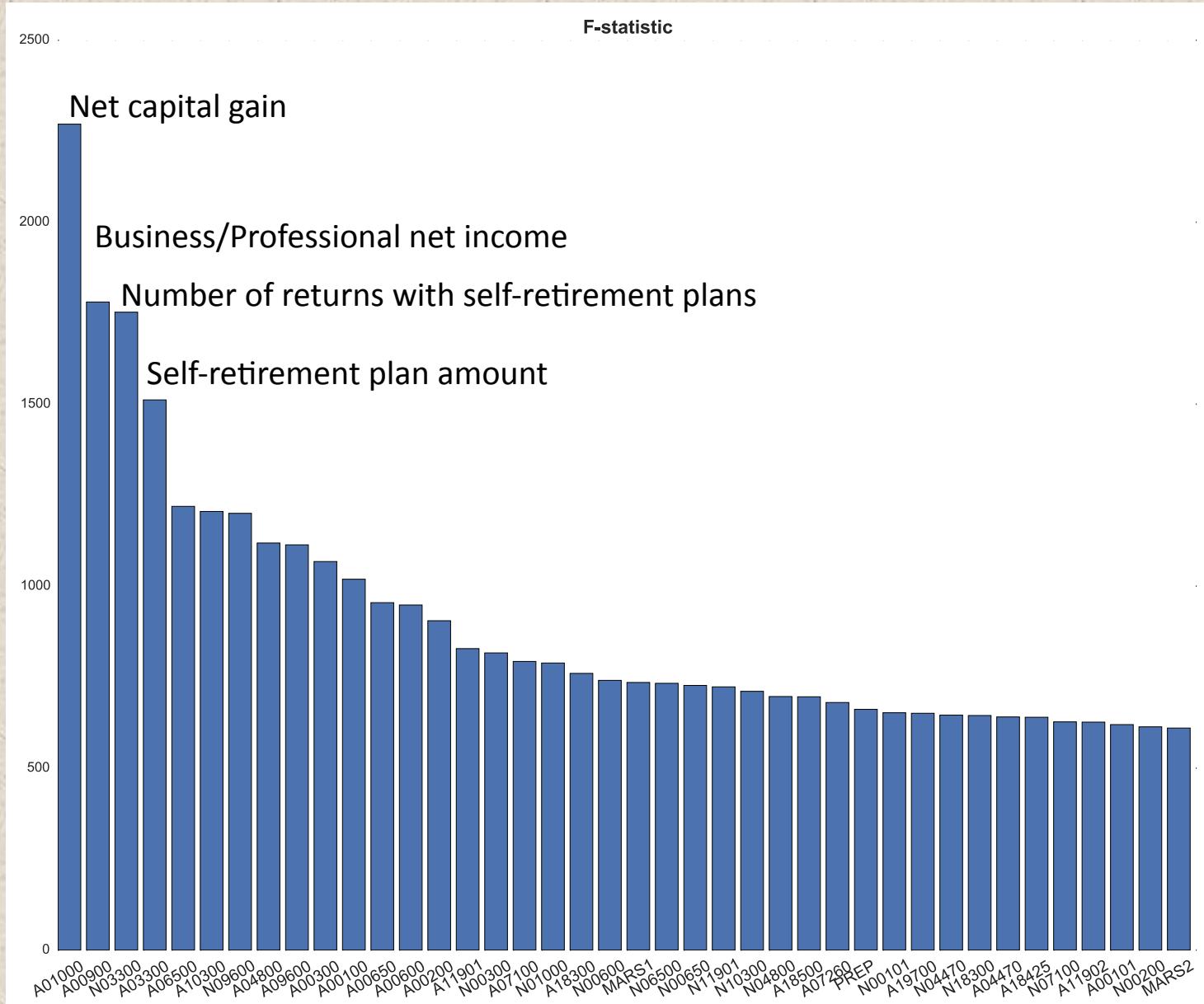
Approaches

- Dimensionality reduction
- Supervised learning

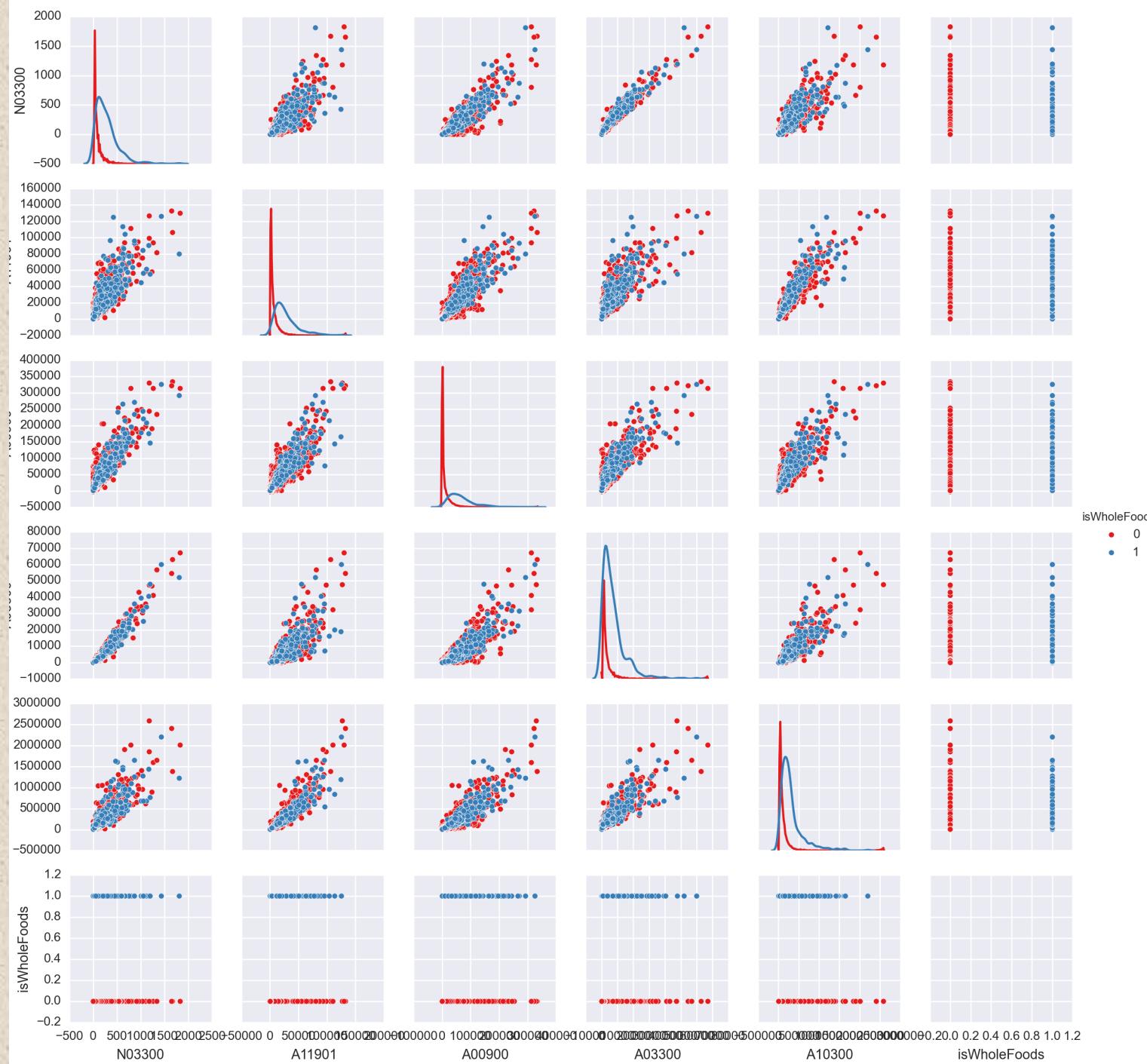
IRS data: 73 features

N1	Number of returns	N18500	Number of returns with real estate taxes
MARS1	Number of single returns	A18500	Real estate taxes amount
MARS2	Number of joint returns	N18300	Number of returns with taxes paid
MARS4	Number of head of household returns	A18300	Taxes paid amount
PREP	Number of returns with paid preparer's signature	N19300	Number of returns with mortgage interest paid
N2	Number of exemptions	A19300	Mortgage interest paid amount
NUMDEP	Number of dependents	N19700	Number of returns with contributions
A00100	Adjust gross income (AGI) [2]	A19700	Contributions amount
N00200	Number of returns with salaries and wages	N04800	Number of returns with taxable income
A00200	Salaries and wages amount	A04800	Taxable income amount
N00300	Number of returns with taxable interest	N09600	Number of returns with alternative minimum tax
A00300	Taxable interest amount	A09600	Alternative minimum tax amount
N00600	Number of returns with ordinary dividends	N07100	Number of returns with total tax credits
A00600	Ordinary dividends amount	A07100	Total tax credits amount
N00650	Number of returns with qualified dividends	N07180	Number of returns with child and dependent care credit
A00650	Qualified dividends amount [3]	A07180	Child and dependent care credit amount
N00900	Number of returns with business or professional net income (less loss)	N07220	Number of returns with child tax credit
A00900	Business or professional net income (less loss) amount	A07220	Child tax credit amount
SCHF	Number of farm returns	N07260	Number of returns with residential energy tax credit
N01000	Number of returns with net capital gain (less loss)	A07260	Residential energy tax credit amount
A01000	Net capital gain (less loss) amount	N11070	Number of returns with additional child tax credit
N01400	Number of returns with taxable individual retirement arrangements distributions	A11070	Additional child tax credit amount
A01400	Taxable individual retirement arrangements distributions amount	N59660	Number of returns with earned income credit
N01700	Number of returns with taxable pensions and annuities	A59660	Earned income credit amount [5]
A01700	Taxable pensions and annuities amount	N59720	Number of returns with excess earned income credit
N02300	Number of returns with unemployment compensation	A59720	Excess earned income credit (refundable) amount [6]
A02300	Unemployment compensation amount [4]	N06500	Number of returns with income tax
N02500	Number of returns with taxable Social Security benefits	A06500	Income tax amount [7]
A02500	Taxable Social Security benefits amount	N10300	Number of returns with tax liability
N03300	Number of returns with self-employment retirement plans	A10300	Total tax liability amount [8]
A03300	Self-employment retirement plans amount	N11901	Number of returns with tax due at time of filing
N04470/N00101	Number of returns with itemized deductions	A11901	Tax due at time of filing amount [9]
		N11902	Number of returns with overpayments refunded

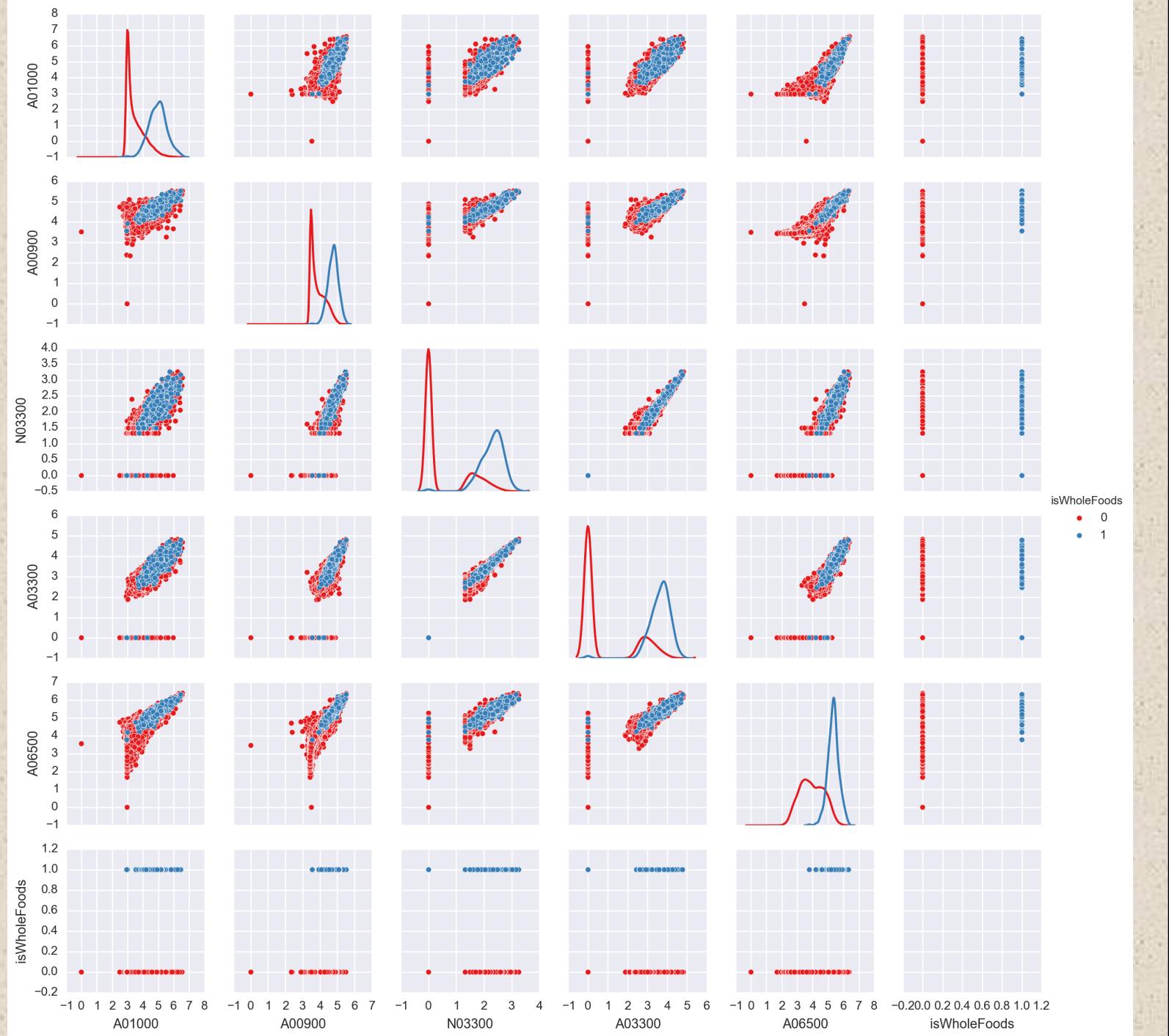
F-statistic



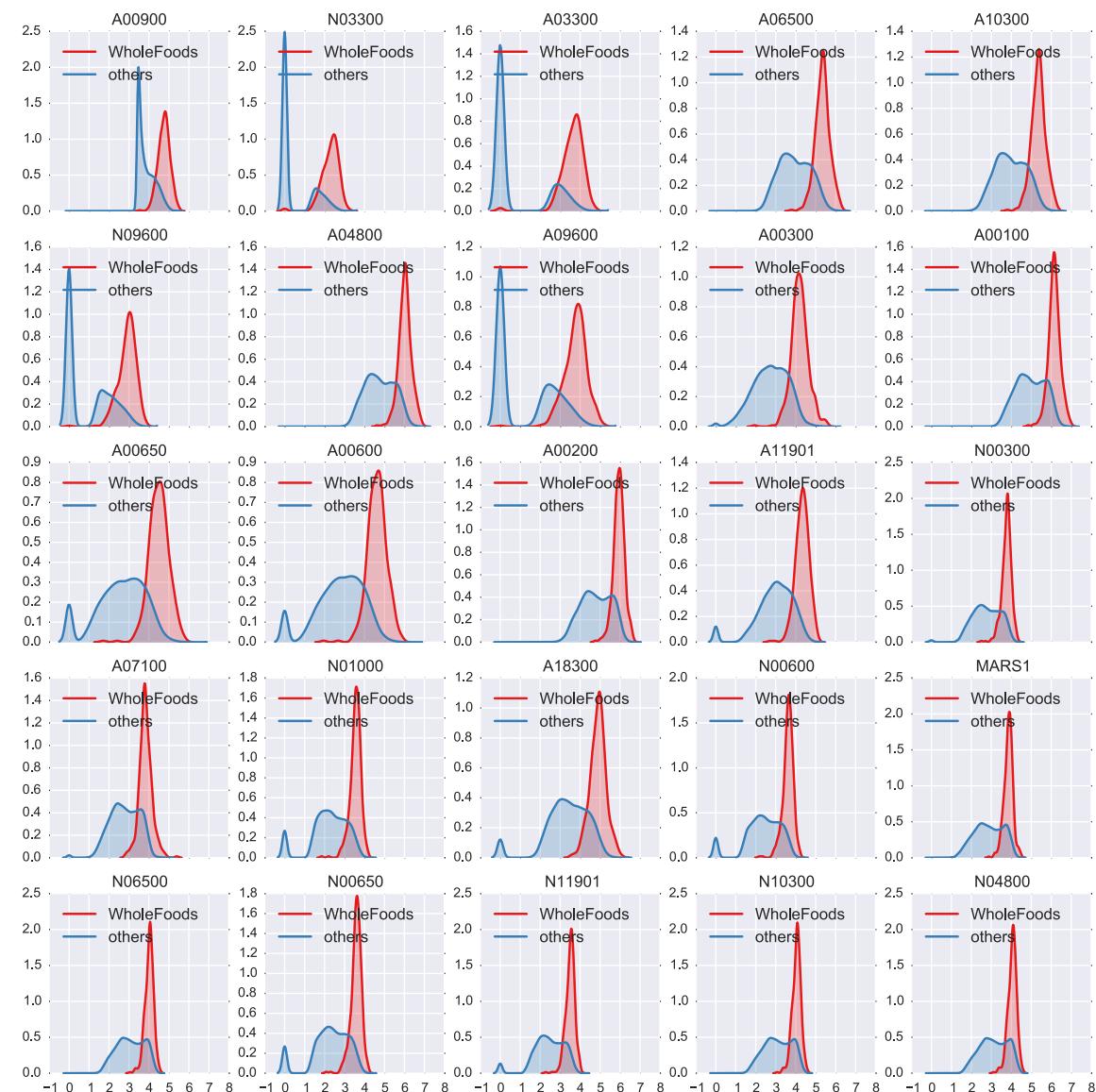
Original Units



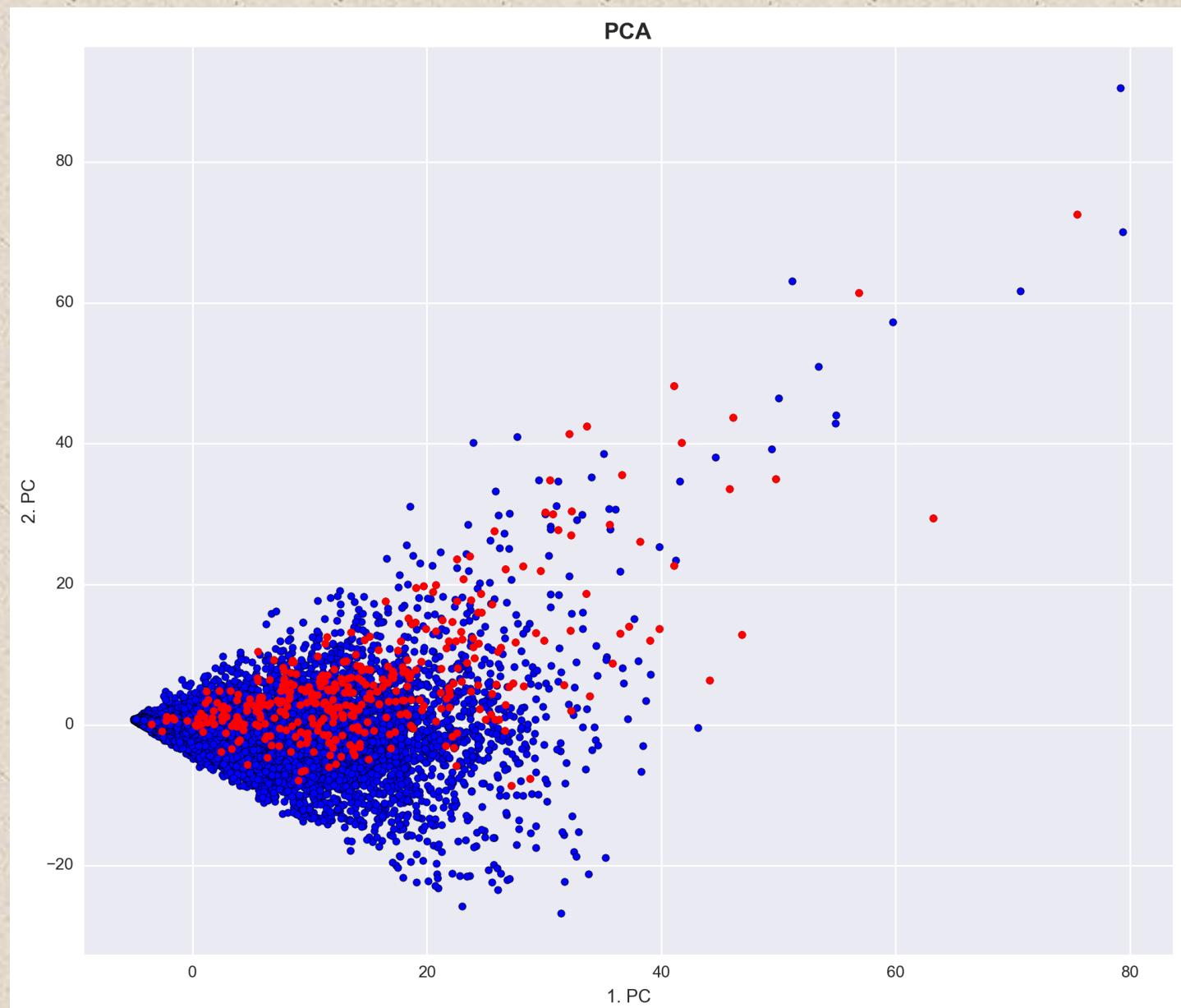
Log-transformed



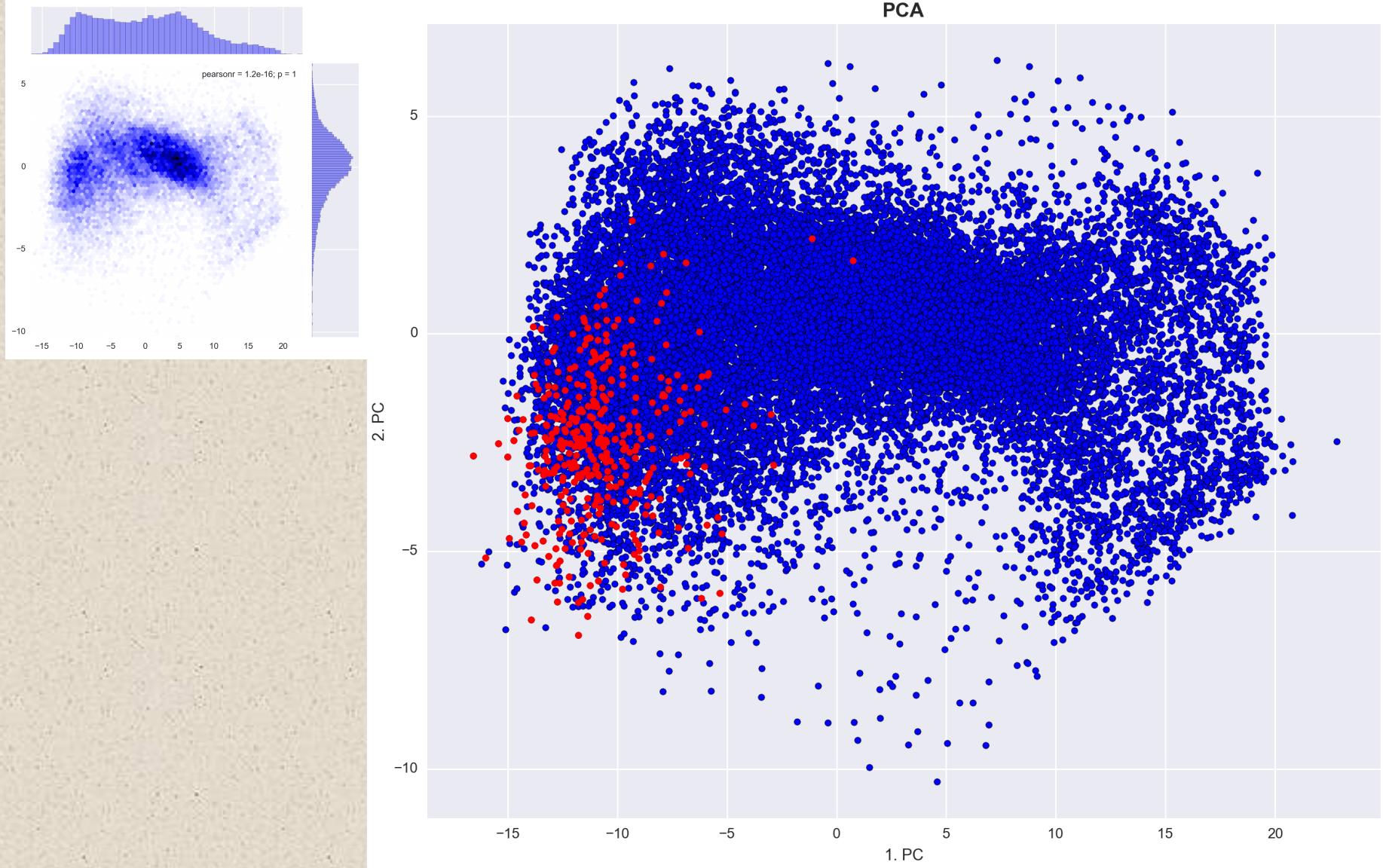
Histograms



PCA:



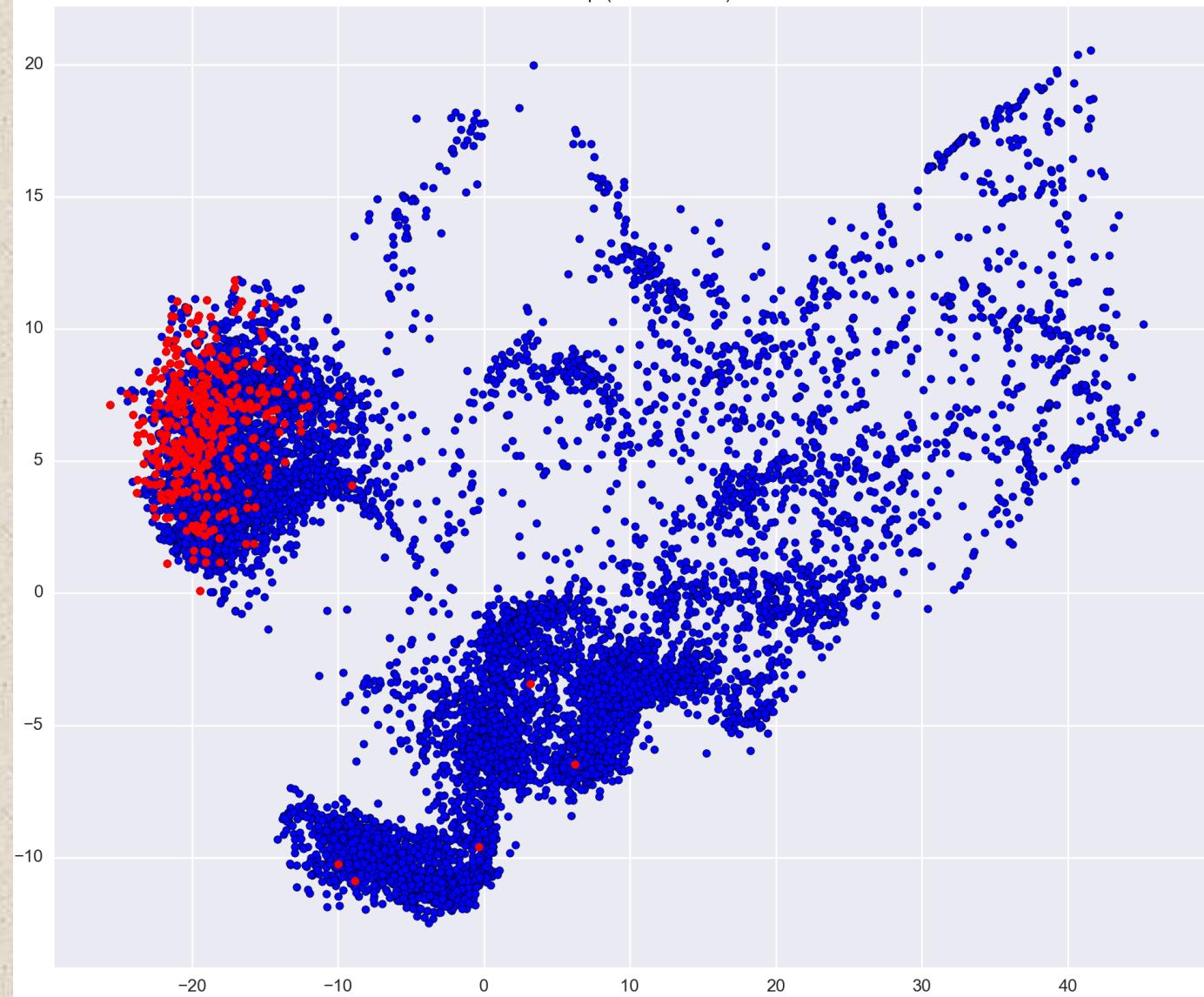
PCA:



Non-linear dimensionality reduction: Isomap

intrinsic geometry of the data manifold (subsample)

Isomap (4.8e+02 sec)

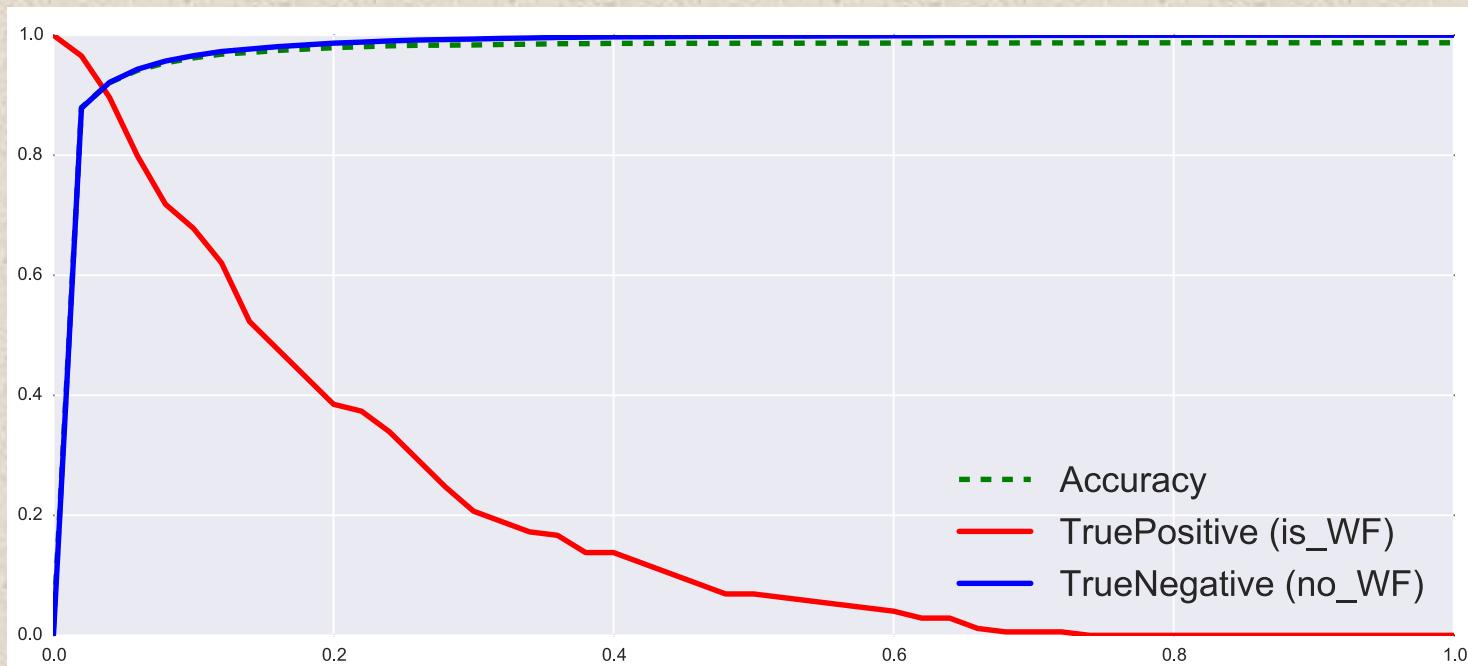
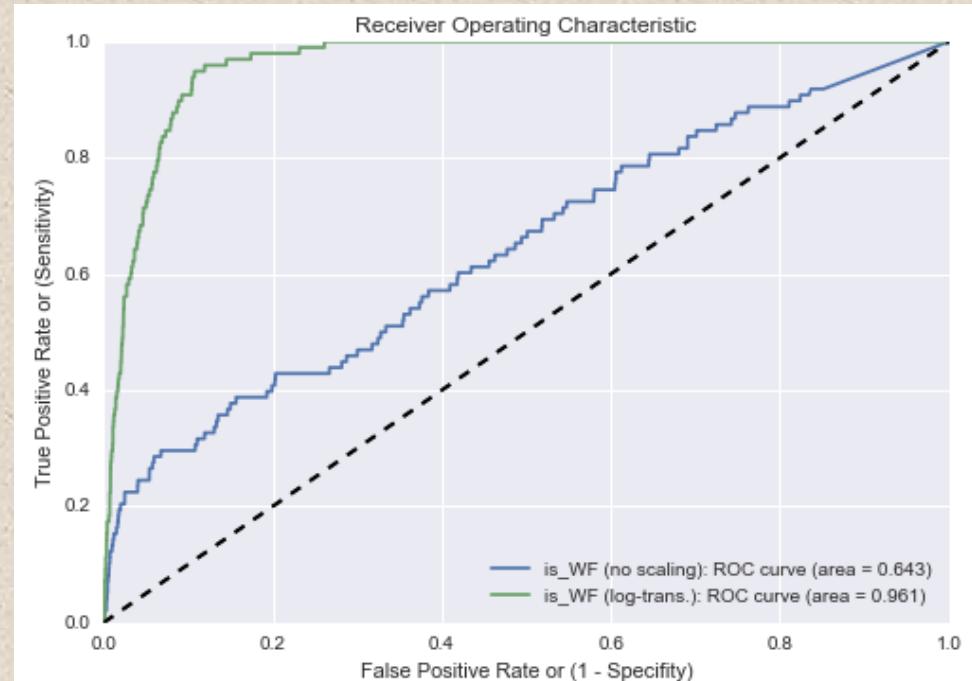


Model: Logistic Regression

- Log-transformed features
- Cross-validation
- Optimize for True Positives
- ROC curves, TruePositive curve, LR thresholds
- Consider False Positives as new locations for new locations for Whole Foods

Logistic Regression:

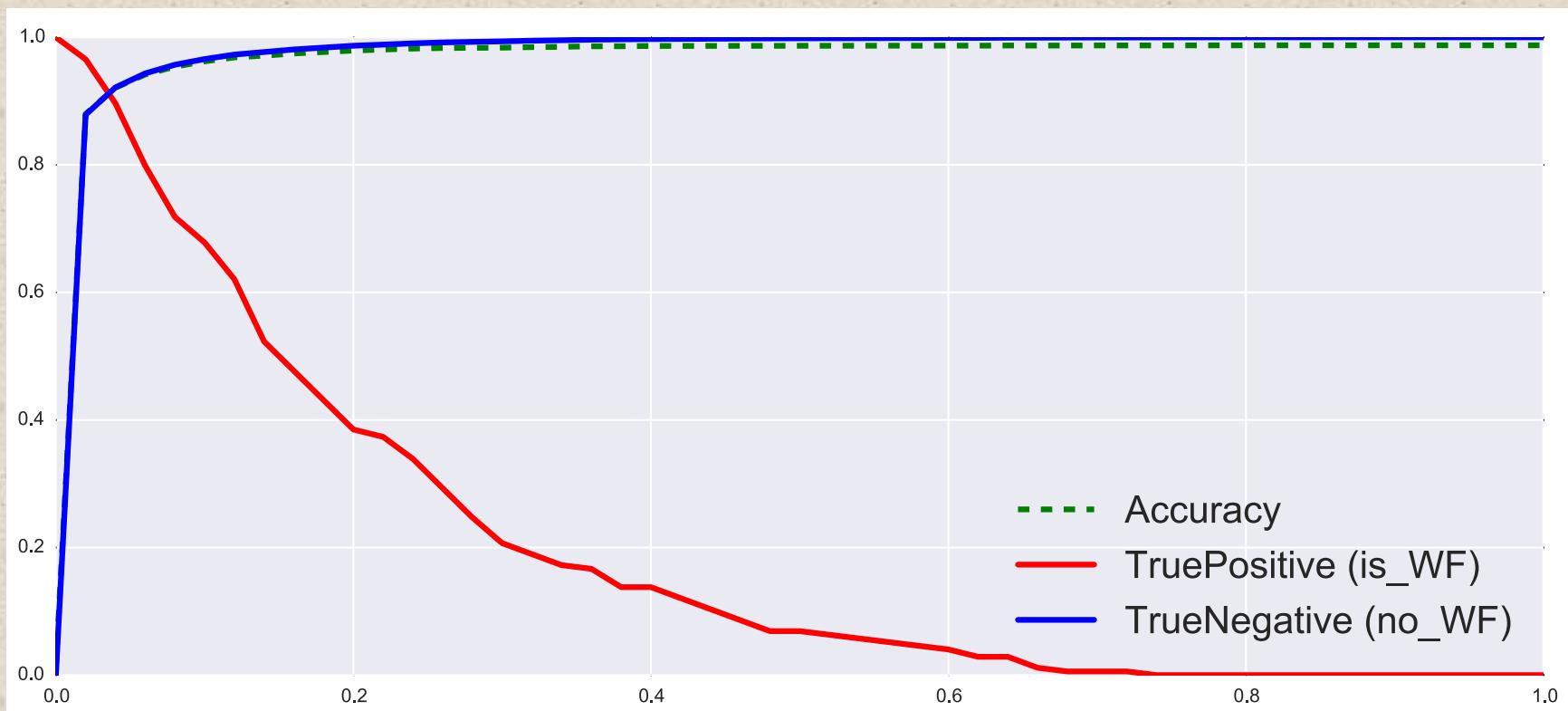
		predicted	
		no_WF	is_WF
actual	no_WF	8094	124
	is_WF	83	15



Logistic Regression Thresholds:



TP: 97%
TN: 88%



Summary:

- Model can reduce potential zip codes by 88%.
- Prediction power of the data may be reached.
- More features may be necessary to exclude more zip codes.