

GA DSI-10 CAPSTONE PROJECT

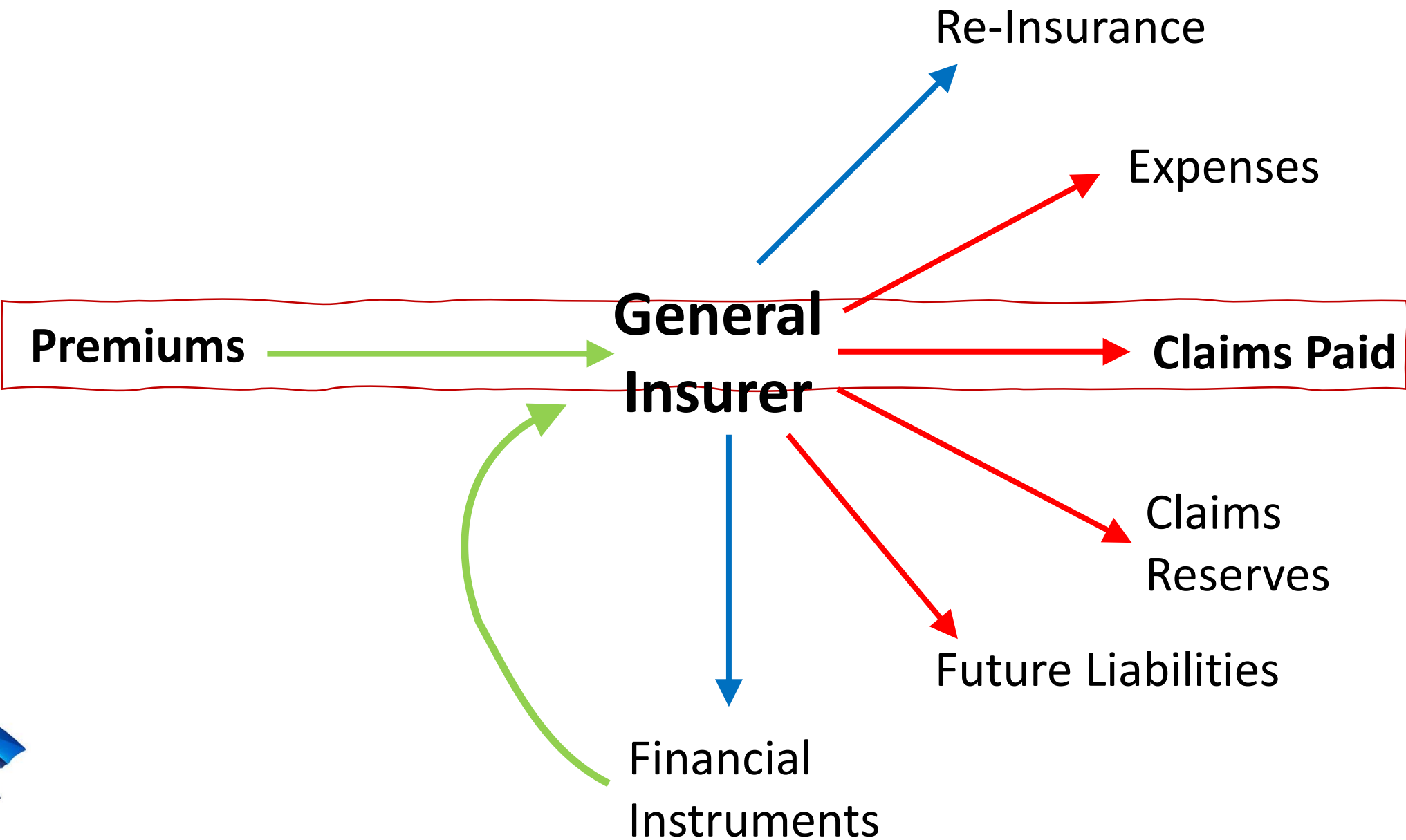
Predicting Whether Insurance
Underwriting Gain Will Be Negative

Irwin Wei

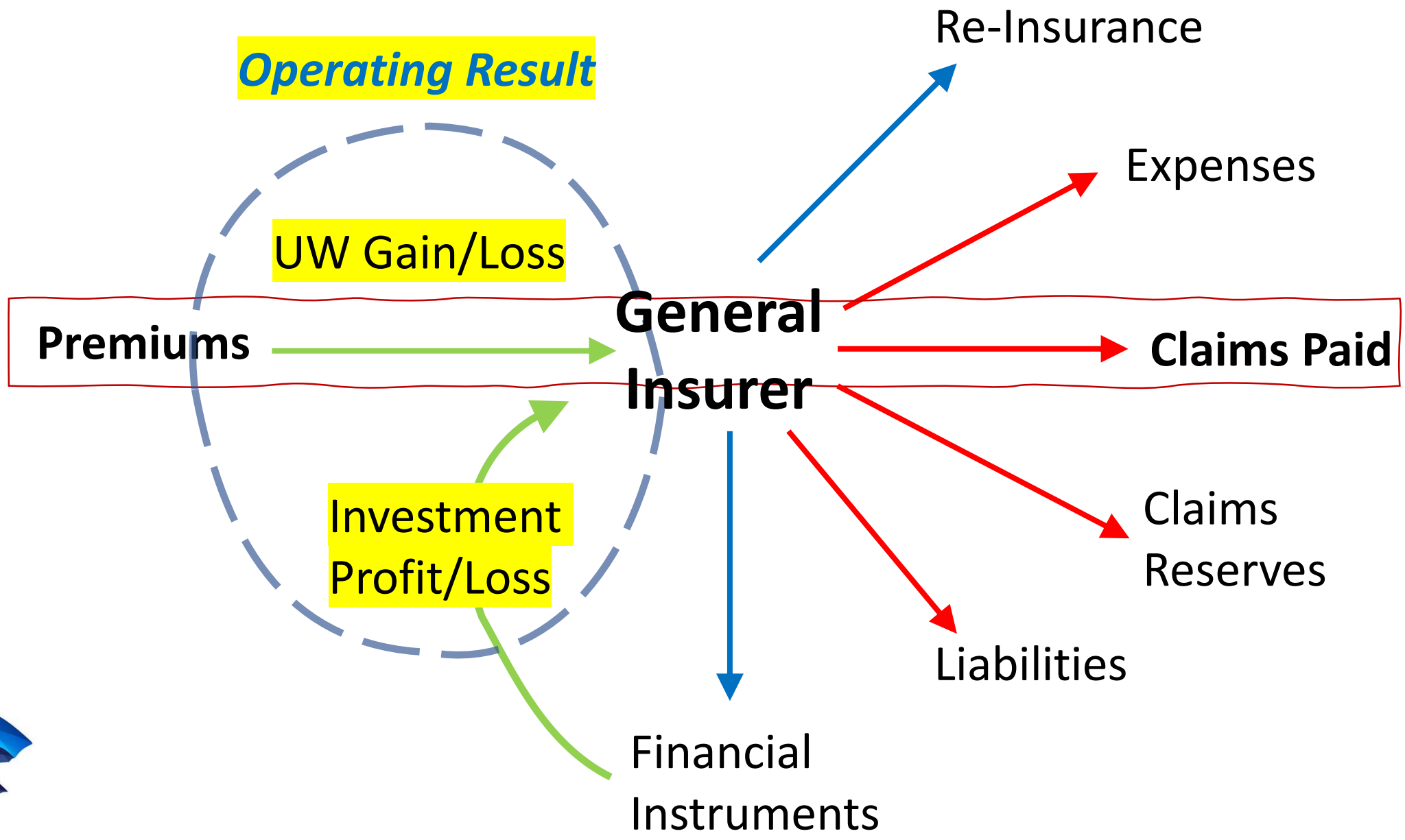
5 December 2019



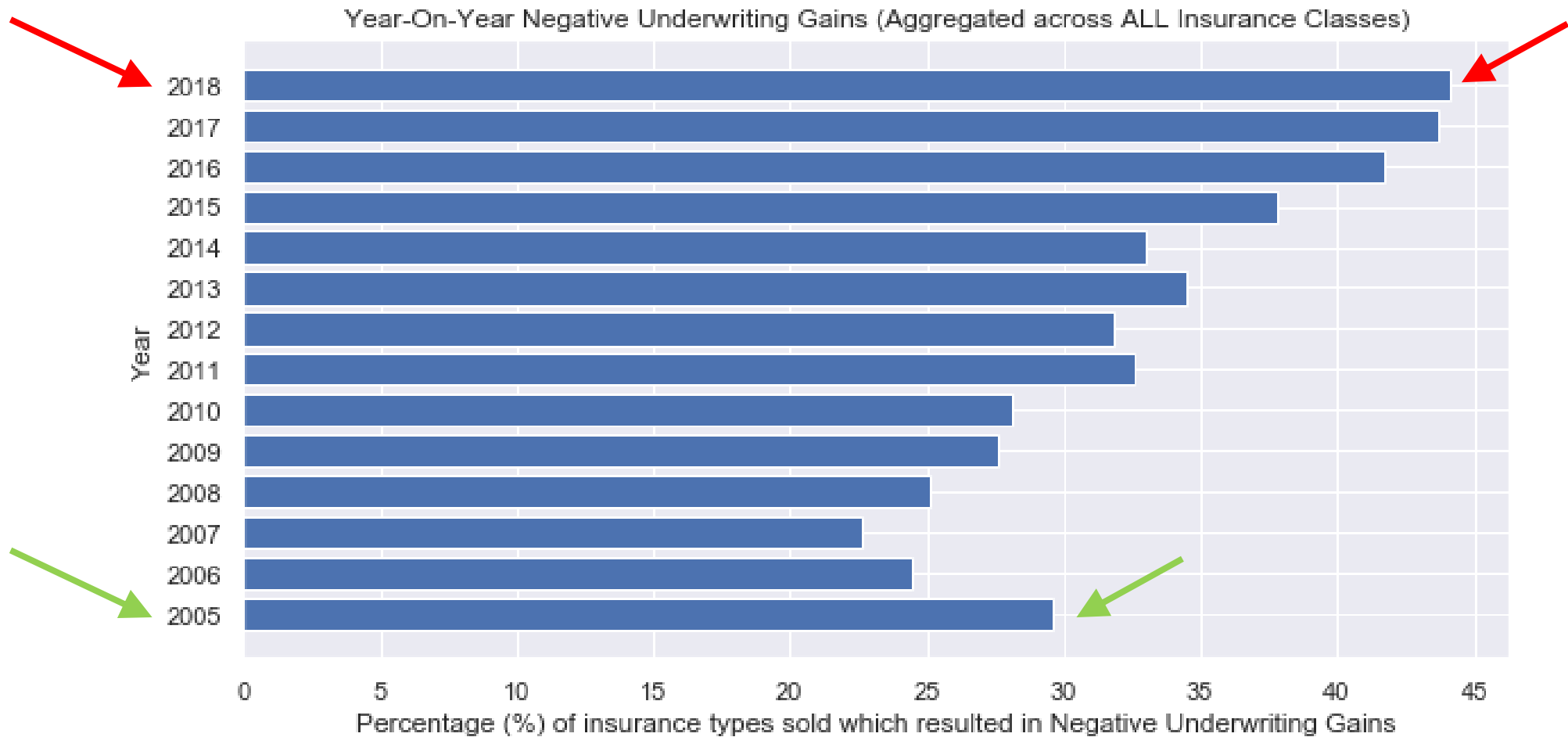
INSURANCE MONEY: INs & OUTs



INSURANCE MONEY: INs & OUTs



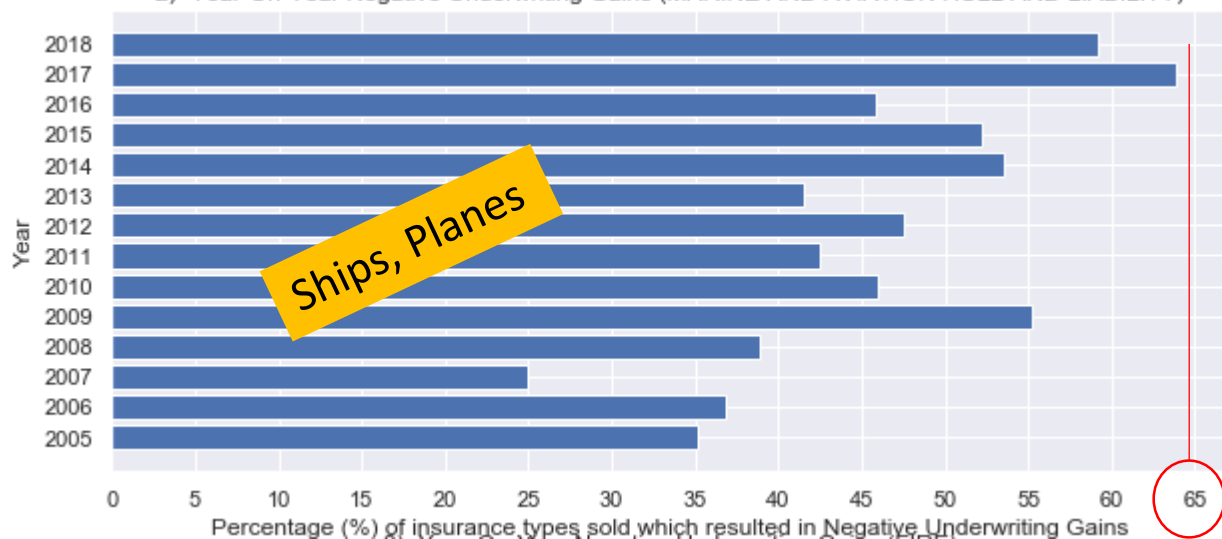
The Rising Trend of Underwriting Losses



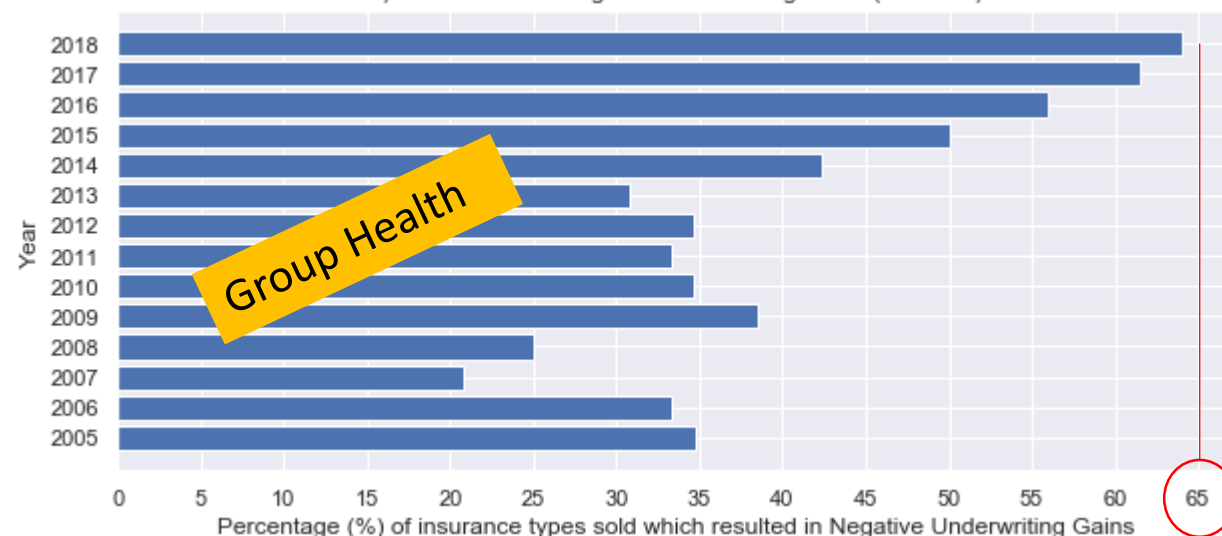


Losing streak ...

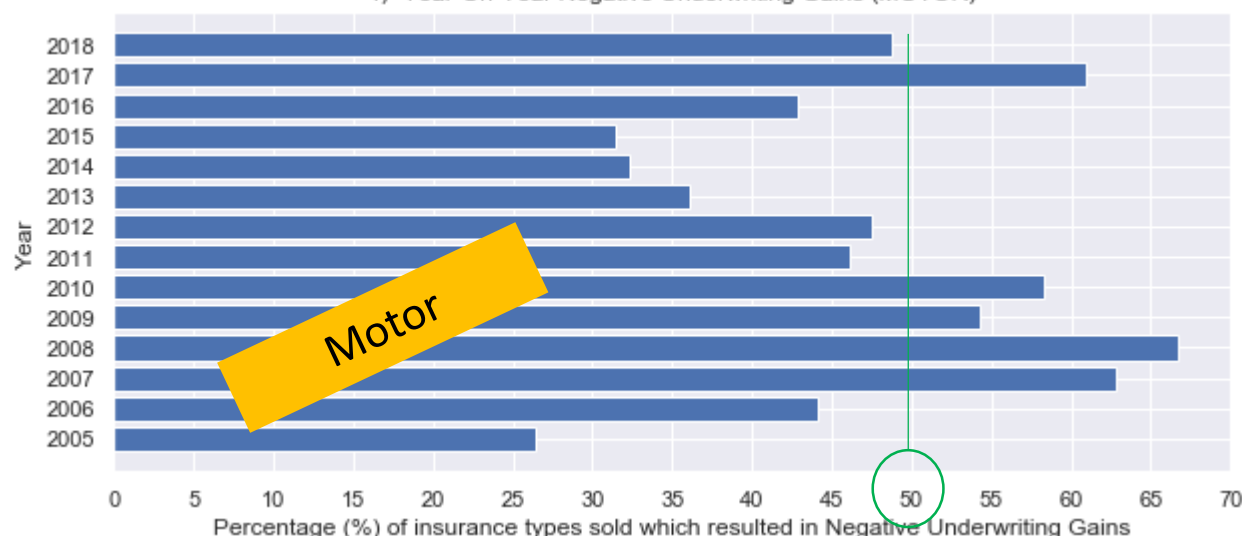
2) Year-On-Year Negative Underwriting Gains (MARINE AND AVIATION HULL AND LIABILITY)



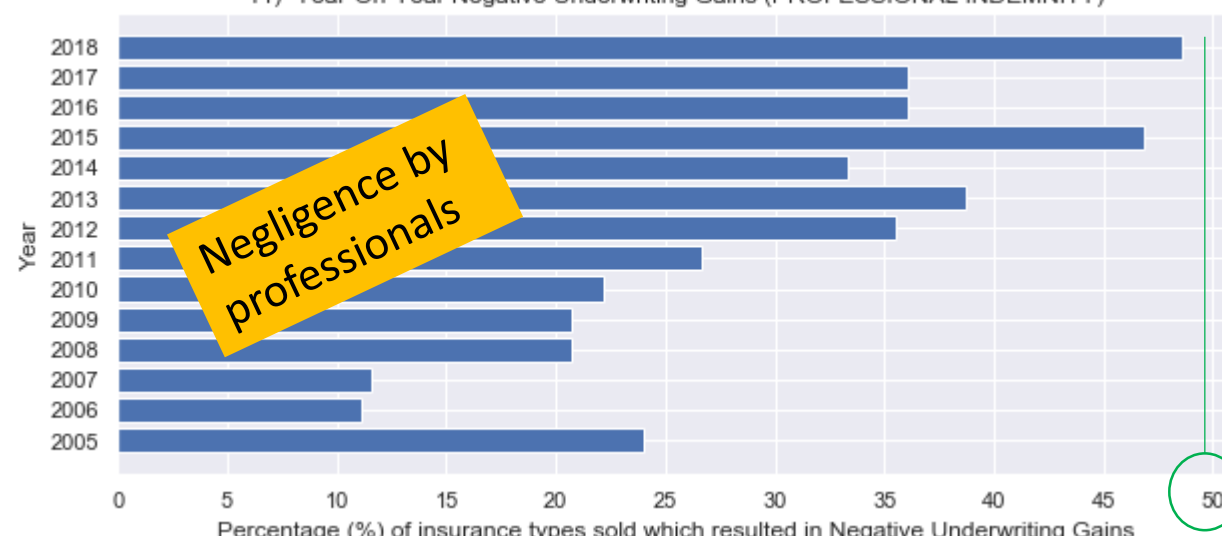
7) Year-On-Year Negative Underwriting Gains (HEALTH)



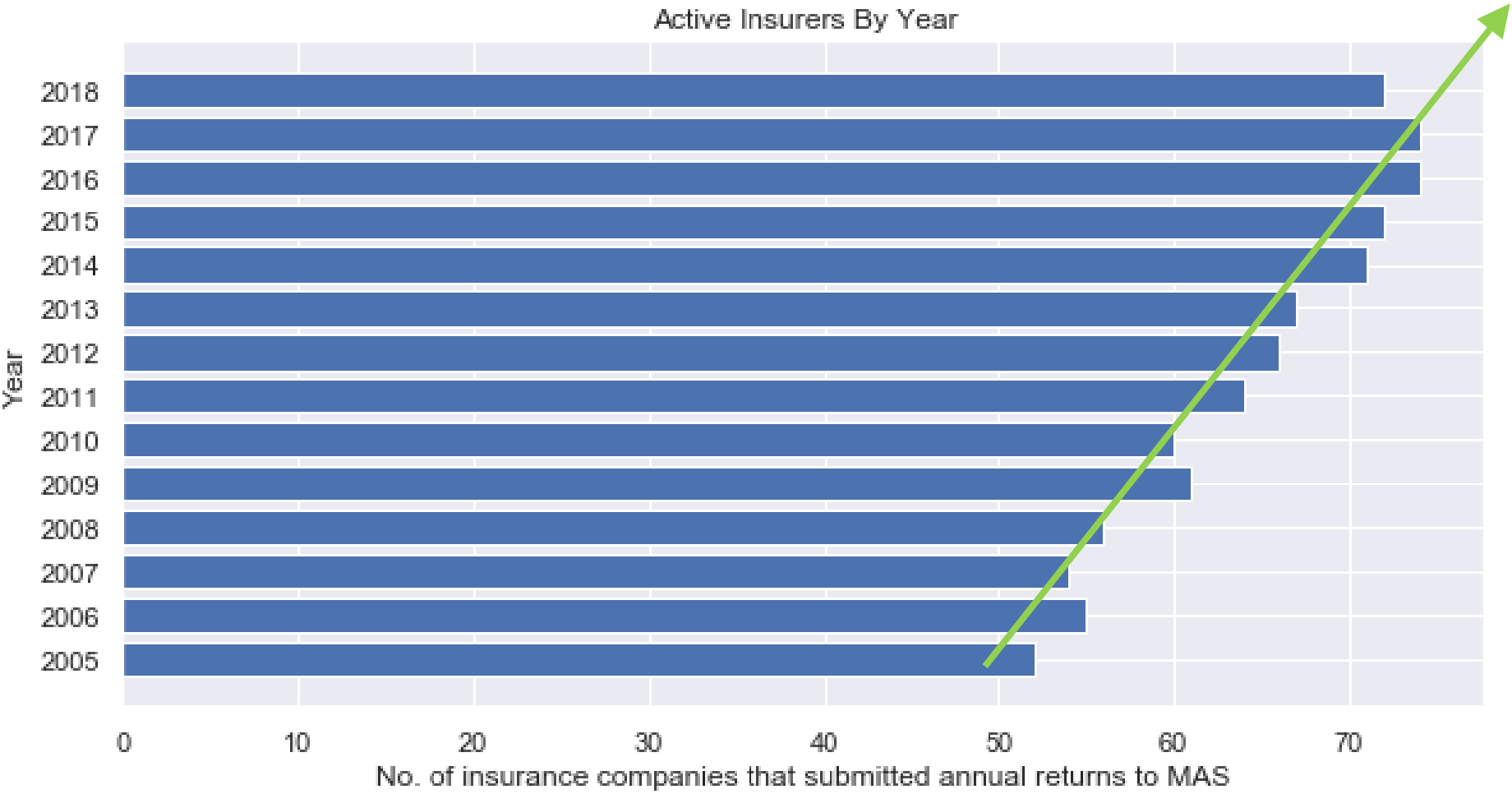
4) Year-On-Year Negative Underwriting Gains (MOTOR)



11) Year-On-Year Negative Underwriting Gains (PROFESSIONAL INDEMNITY)



... and yet... More Insurers!



Problem Statement

To implement a proof-of-concept, employing only publicly available data, to predict whether the underwriting performance of any given insurance class is likely to result in an underwriting loss at the end of the current 12-month reporting period.



Value Proposition

With such predictions, underwriters can place more focus on certain insurance classes and review their underwriting approach, and/or take necessary risk management measures such as re-insuring more.



Caveat


Publicly available data is very coarse. Only shows the start/end state of a 12-month period.

Insurers have much higher quality, high-resolution data concerning underwriting and claims details.



Description	Row No.	Marine and Aviation - Cargo	Marine and Aviation - Hull and Liability	Fire	Motor	Work Injury Compensation	Personal Accident	Health	Misc - Public Liability	Misc - Bonds	Misc - Engineering / CAR / EAR	Misc - Professional Indemnity	Misc - Credit / Political Risk	Misc - Others	Misc - Sub-Total	Total
A. PREMIUMS																
Gross premiums																
Direct business	1	1,459,528	2,647,014	36,747,999	192,134,197	53,595,401	73,154,783	37,503,540	14,951,033	5,198,553	1,822,746	12,455,454	9,333,364	19,854,998	83,213,873	480,458,413
Reinsurance business accepted -																
In Singapore	2	0	346,573	1,123,873	0	0	0	0	832,885	144,034	1,521,033	158,102	0	2,868,321	5,452,184	6,825,768
From other ASEAN countries	3	0	0	9,904	0	0	0	0	0	0	0	0	0	0	0	9,904
From other countries	4	367	407,239	579,129	40,365	11,260	23,248	0	3,274	32,268	860	2,850	1,361	4,758	45,511	1,107,048
Total (2 to 4)	5	367	753,792	1,713,005	40,365	11,260	23,248	0	835,959	176,242	1,521,893	160,752	1,361	3,001,079	5,497,675	8,043,683
Reinsurance business ceded -																
In Singapore	6	0	0	367,875	311	590,024	251,265	(43)	0	0	0	0	0	29,967	29,967	1,148,069
To other ASEAN countries	7	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
To other countries	8	1,067,998	3,314,779	26,868,179	17,178,901	6,372,411	8,052,257	1,978,083	2,839,529	4,280,834	3,144,640	3,857,348	9,087,485	11,824,478	34,747,398	97,015,673
Total (6 to 8)	9	1,067,998	3,314,779	27,236,054	17,177,212	6,872,435	8,303,523	1,978,040	2,839,529	4,280,834	3,144,640	3,857,348	9,087,485	11,851,143	34,773,853	98,761,873
Net premiums written (1 + 5 - 9)	10	371,835	68,027	11,264,850	174,967,350	46,734,326	68,874,482	35,523,521	13,947,483	1,678,971	(233)	8,858,854	247,884	11,774,826	33,807,685	368,737,083
Premium liabilities at beginning of period	11	48,383	42,880	6,378,654	65,868,041	15,421,829	15,239,890	11,696,811	3,533,735	773,568	0	4,755,206	12,202	15,090,292	33,153,918	159,990,106
Premium liabilities at end of period	12	267,668	53,593	6,371,841	73,963,301	14,320,105	14,674,483	10,159,239	3,280,513	898,000	0	4,852,170	86,043	13,768,323	33,003,648	144,683,488
Premiums earned during the period (10 + 11 - 12)	13	212,230	76,294	11,309,793	186,613,090	47,835,950	67,439,867	37,267,103	12,339,865	918,479	(233)	8,481,896	173,882	12,326,868	34,598,585	365,043,711
B. CLAIMS																
Gross claims settled																
Direct business	14	24,008	595,197	10,807,119	109,041,158	23,334,384	18,620,483	19,036,153	1,380,907	877,343		36,648	1,871,299	483,084	3,826,314	187,718,058
Reinsurance business accepted -																
In Singapore	15	0	0	0	0	0	0	0	0	11,783	421,783	9,890	0	0	443,256	443,256
From other ASEAN countries	16	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
From other countries	17	0	0	5,844	0	0	0	0	0	0	0	0	0	0	0	5,844
Total (15 to 17)	18	0	0	5,844	0	0	0	0	0	11,783	421,783	9,890	0	0	443,256	449,100
Recoveries from reinsurance business ceded -																
In Singapore	19	0	0	82,531	825	0	0	0	0	(2,198)	0	0	0	0	80,333	81,331
To other ASEAN countries	20	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
To other countries	21	18,008	592,992	10,112,850	57,625	1,454,808	485,182	107,177	317,959	989,281	453,431	381,422	497,307	12,358	2,543,372	15,421,638
Total (19 to 21)	22	18,008	592,992	10,195,187	58,450	1,454,808	485,182	107,177	317,959	987,083	453,431	381,422	497,307	12,358	2,543,349	15,482,959
Net claims settled (14 + 18 - 22)	23	6,000	2,205	617,578	108,982,678	21,879,576	18,135,301	18,928,975	942,948	23,360		1,582,549	173,882	3,813,314	173,718,967	
Claims liabilities at end of period	24	29,599	103,277	1,487,839	121,486,439	41,762,917	13,818,106	8,373,443	8,892,899	485,885		8,193,841	28,414	3,711,323	208,637,689	
Claims liabilities at beginning of period	25	8,417	57,931	1,043,039	119,545,941	34,818,724	14,584,125	6,582,127	6,403,887	898,715	11,375	8,827,795	101,673	5,890,230	21,843,485	198,464,759
Net claims incurred (23 + 24 - 25)	26	27,184	47,551	1,082,179	119,933,174	26,823,869	15,369,302	20,746,290	3,431,960	(100,937)	19,451	825,729	136,648	1,493,305	5,900,381	182,803,910
C. MANAGEMENT EXPENSES																
Management Expenses	27	299,507	690,073	7,811,890	42,572,178	10,867,675	18,585,198	9,509,853	3,180,037	1,323,067		637,498	2,557,753	1,092,570	4,593,191	104,477,367
D. DISTRIBUTION EXPENSES																
Commissions	28	254,422	516,923	2,702,179	29,623,007	5,975,029	17,899,719	6,926,483	2,942,921	693,285	252,573	1,719,314	1,538,145	2,412,588	9,808,828	73,526,566
Reinsurance commissions	29	885,047	900,721	4,821,557	261,283	938,193	979,590	149,138	874,595	1,817,275	1,072,492	252,014	2,901,892	3,182,739	9,700,789	18,018,282
Net commissions incurred (28 - 29)	30	(410,625)	(383,798)	(2,219,378)	29,261,725	5,036,836	18,920,129	6,777,345	2,068,326	(823,990)	(819,819)	1,467,300	(1,313,548)	(770,151)	(91,942)	54,698,303
Other distribution expenses	31	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
E. UNDERWRITING RESULTS																
Underwriting gain / (loss) (13 - 26 - 27 - 30 - 31)	32	299,164	(275,533)	4,855,332	4,045,013	3,108,239	18,585,058	239,834	3,349,382	618,319	182,733	3,511,117	(536,011)	7,010,503	14,118,023	42,753,134
F. NET INVESTMENT INCOME																
	33	9,493	2,273	388,033	4,487,223	1,193,003	1,707,130	906,934	300,515	27,543	0	221,036	6,302	388,921	846,343	6,438,417

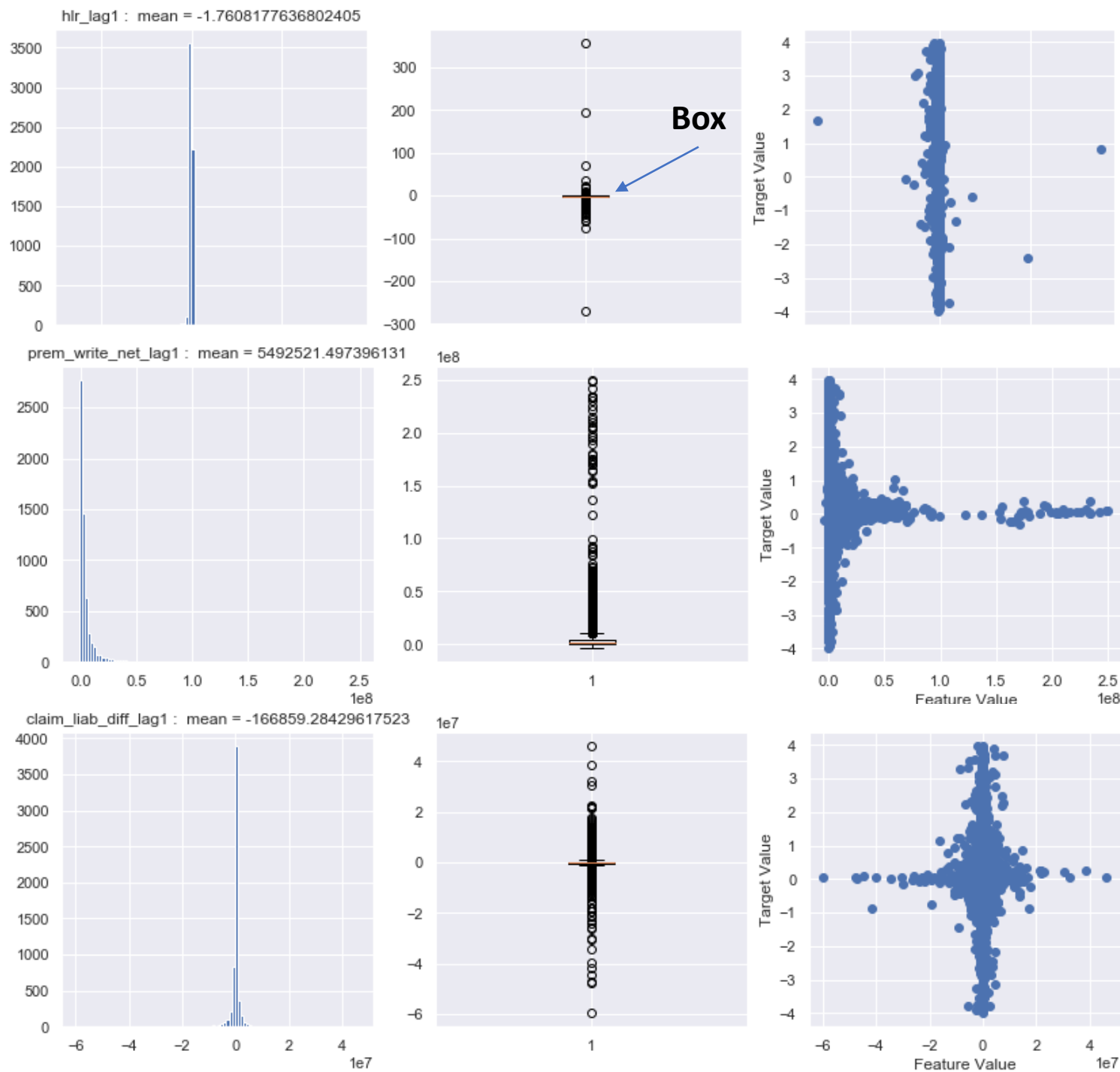
Annual Returns
Submitted to MAS –
Table “FORM 6 (SIF)”



pdfminer
pyPDF4
tabula

[illegible]

Non-Ideal data



~ 9,000 rows



FEATURE ENGINEERING

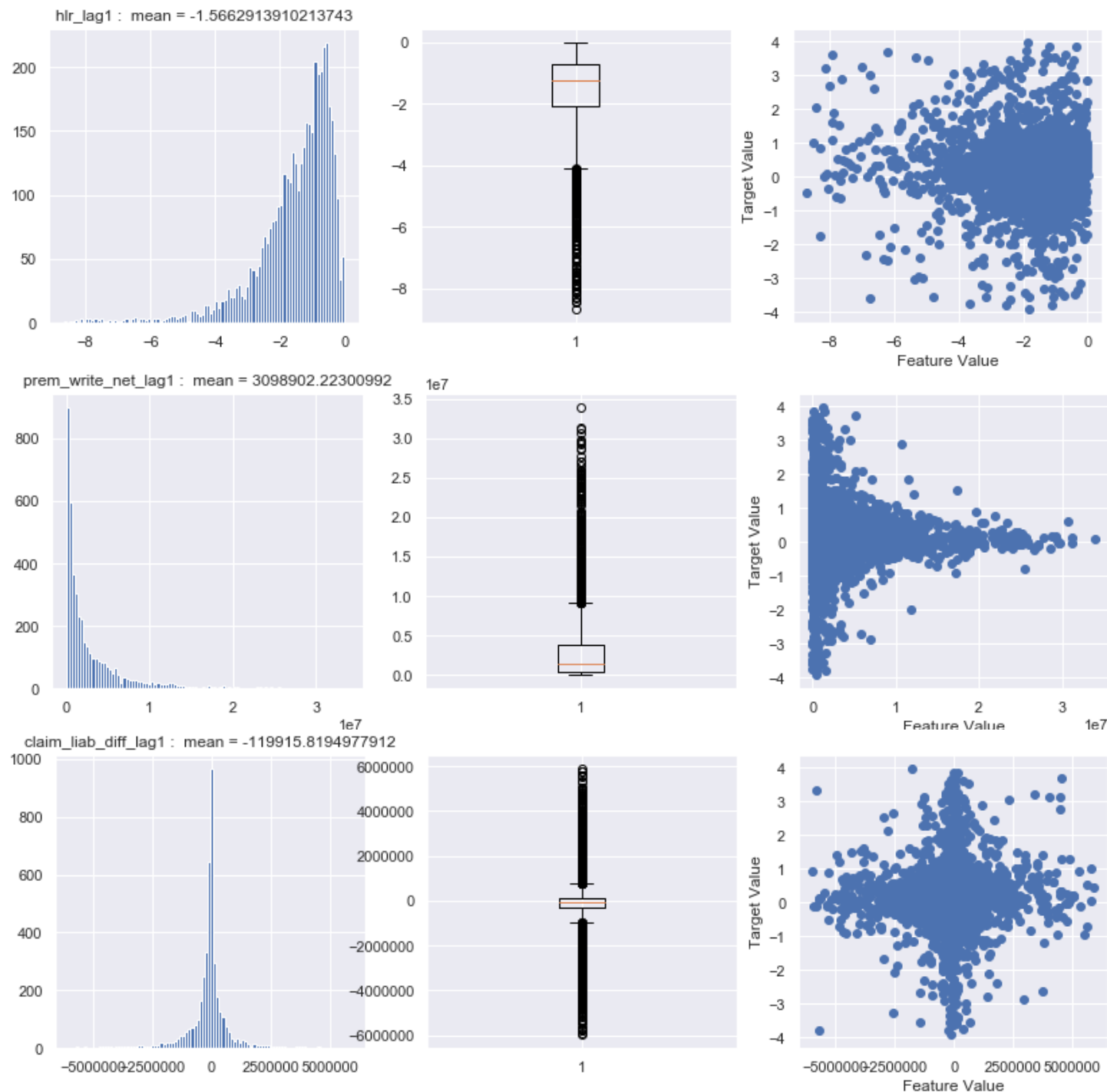
- Calculated ratios in order to normalize data
- Used exponential weighted means on time-lag data (up to 5 years' history)
- Reduced to 5 features, then built them up again
- +/- polarity of features depending on whether it represented incoming/outgoing money



After
trimming
away
Extreme
Outliers

&

Feature
Engineering



~ 4,400 rows

REGRESSION MODELS

- Linear Regression
- Decision Tree
- Random Forest
- Extra Trees
- Ada Boost
- Gradient Boosting



Feature
Re-Engineering



UNSUPERVISED LEARNING

- K-Means
- DBSCAN
- Hierarchical

No meaningful or interpretable clusters...



REGRESSION MODELS

- Linear Regression
- Decision Tree
- Random Forest
- Extra Trees
- Ada Boost
- Gradient Boosting



Feature
Re-Engineering



Best R2 Score => 0.53 (only)



UNSUPERVISED LEARNING

- K-Means
- DBSCAN
- Hierarchical

No meaningful or interpretable clusters...

CLASSIFICATION

Positive Case (1) → Underwriting Gain < 0

Negative Case (0) → Underwriting Gain ≥ 0

```
In [8]: 1 df['classification'].value_counts(normalize=True).sort_index()
```

executed in 34ms, finished 18:49:39 2019-12-04

```
Out[8]: 0    0.680699
```

```
1    0.319301
```

```
Name: classification, dtype: float64
```

Baseline Accuracy is 0.68



CLASSIFICATION MODELS

- Logistic Regression
- Decision Tree
- Random Forest
- Extra Trees
- Ada Boost
- Gradient Boosting
- K Nearest Neighbors
- Support Vector Machine

METRICS:

Recall

Accuracy



CLASSIFICATION MODELS

- Logistic Regression
- Decision Tree
- Random Forest
- Extra Trees
- Ada Boost
- Gradient Boosting
- K Nearest Neighbors
- Support Vector Machine

METRICS:

Recall

Accuracy



6.1 Logistic Regression Classifier

In [13]:

```
1 # lr = LogisticRegression(fit_intercept=False,C=1.0,tol
2 lr = LogisticRegression(solver='lbfgs')
3 lr.fit(X_train,y_train)
4 print('Score(train/test):',lr.score(X_train, y_train),
5 # print('Score (test):\t',lr.score(X_test, y_test))
6 print('\n=== Classification Report =====
7 # predict & evaluate
8 predictions = lr.predict(X_test)
9 print(classification_report(y_test,predictions,target_n
```

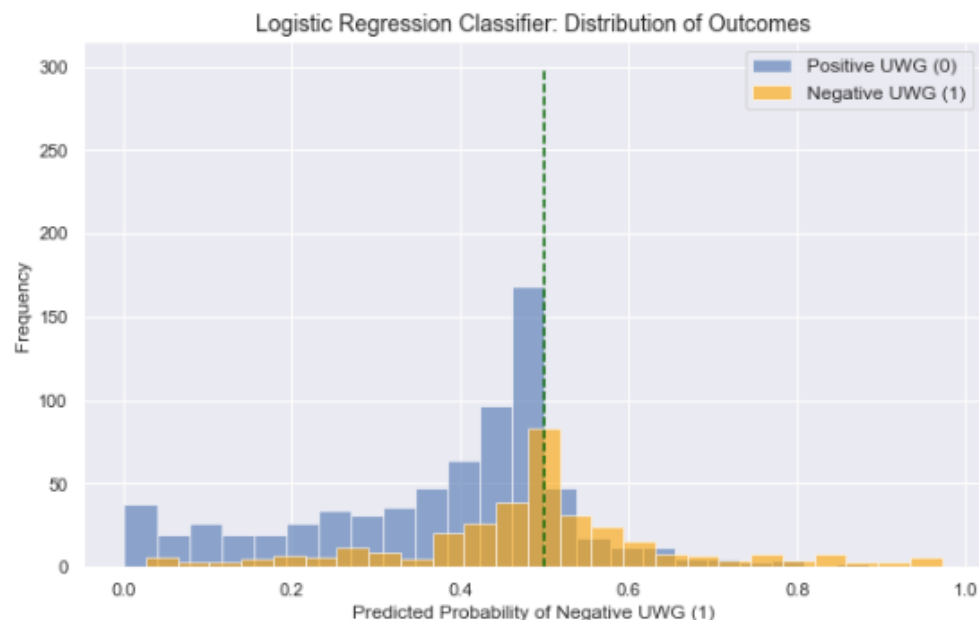
executed in 85ms, finished 16:40:39 2019-12-04

Score(train/test): 0.7239263803680982 , 0.7230910763569457

```
=== Classification Report =====
precision    recall  f1-score   support

Positive UWG (0)    0.76    0.86    0.81     740
Negative UWG (1)    0.59    0.44    0.50     347

   accuracy          0.72     1087
  macro avg          0.68    0.65    0.65     1087
 weighted avg          0.71    0.72    0.71     1087
```



After Tuning

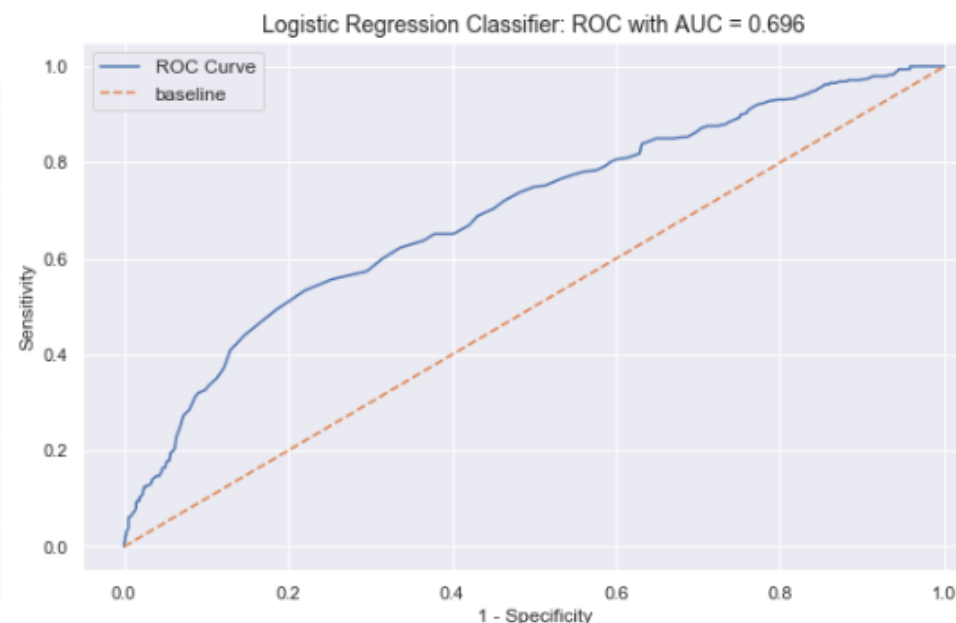
Baseline Accuracy is 0.68

Score(train/test): 0.7220858895705522 , 0.7230910763569457

```
=== Classification Report =====
precision    recall  f1-score   support

Positive UWG (0)    0.77    0.84    0.81     740
Negative UWG (1)    0.58    0.46    0.52     347

   accuracy          0.72     1087
  macro avg          0.68    0.65    0.66     1087
 weighted avg          0.71    0.72    0.71     1087
```



{'C': 1.0,
'fit_intercept': True,
'max_iter': 5000,
'penalty': 'l2', 'solver':
'sag', 'tol': 0.1}

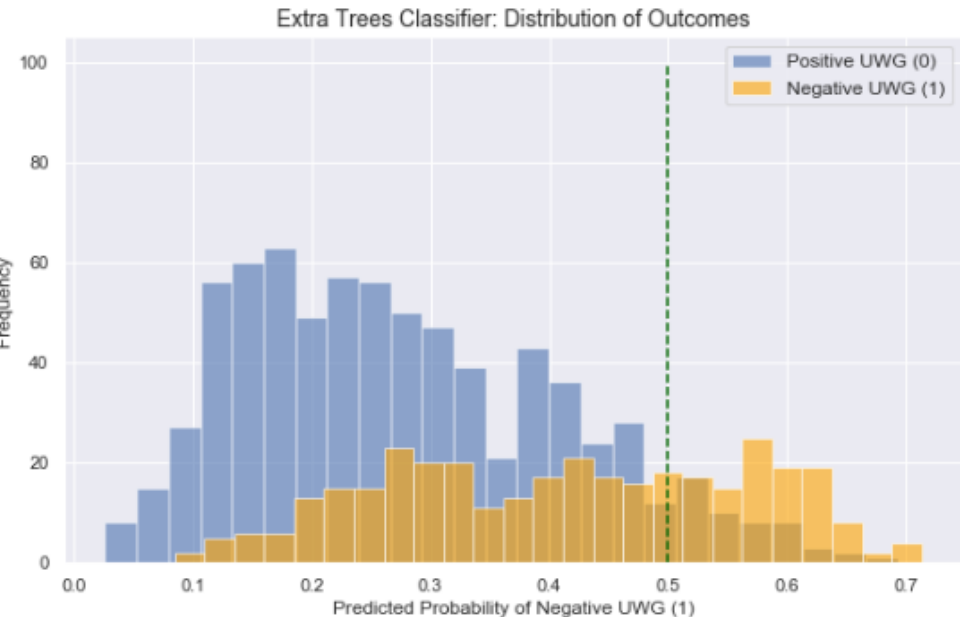
6.4 Extra Trees Classifier

```
In [23]: 1 etcf = ExtraTreesClassifier(bootstrap=True,oob_score=1
2 etcf.fit(X_train,y_train)
3 # Evaluate model.
4 print('Score(train/test):',etcf.score(X_train, y_train)
5 # predict & evaluate
6 predictions = etcf.predict(X_test)
7 print('\n=== Classification Report =====')
8 print(classification_report(y_test,predictions,target_
```

executed in 887ms, finished 16:50:16 2019-12-04

Score(train/test): 1.0 , 0.7332106715731371

=== Classification Report ===				
	precision	recall	f1-score	support
Positive UWG (0)	0.75	0.91	0.82	740
Negative UWG (1)	0.64	0.37	0.47	347
accuracy			0.73	1087
macro avg	0.70	0.64	0.64	1087
weighted avg	0.72	0.73	0.71	1087

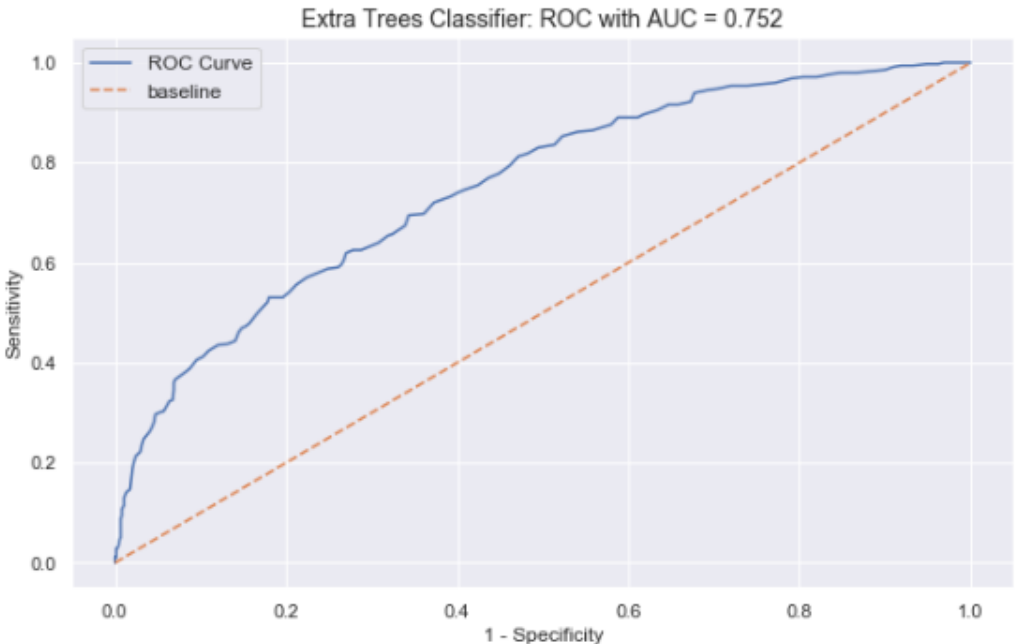


After Tuning

Baseline Accuracy is 0.68

Score(train/test): 0.8076687116564417 , 0.7396504139834407

=== Classification Report ===				
	precision	recall	f1-score	support
Positive UWG (0)	0.75	0.93	0.83	740
Negative UWG (1)	0.70	0.33	0.45	347
accuracy			0.74	1087
macro avg	0.72	0.63	0.64	1087
weighted avg	0.73	0.74	0.71	1087



{'bootstrap': True,
'max_depth': 31,
'min_samples_split':
27, 'n_estimators':
150, 'oob_score': True}

Comparison

Logistic Regressor

```
Score(train/test): 0.7220858895705522 , 0.7230910763569457

=== Classification Report ===
              precision    recall  f1-score   support

Positive UWG (0)       0.77       0.84       0.81       740
Negative UWG (1)       0.58       0.46       0.52       347

   accuracy              0.72       1087
  macro avg              0.68       1087
 weighted avg              0.71       1087
```

Extra Trees

```
Score(train/test): 0.8076687116564417 , 0.7396504139834407

=== Classification Report ===
              precision    recall  f1-score   support

Positive UWG (0)       0.75       0.93       0.83       740
Negative UWG (1)       0.70       0.33       0.45       347

   accuracy              0.74       1087
  macro avg              0.72       1087
 weighted avg              0.73       1087
```



Conclusion

- *Whilst Extra Trees did slightly better in accuracy, Logistic Regression performed much better on Recall ← and that's the most important metric for this problem statement.*
- *Using coarse-grained, low-resolution data from the annual returns alone, the classification model was already able to out-perform the baseline accuracy for classification. Extrapolating further, should insurers' underwriting and claims data be available for analysis, it is certain that far greater insights can be extracted.*

