

CAN MACHINES UNDERSTAND HUMAN DECISIONS: DISSECTING STOCK FORECASTING SKILL

Sean Cao^a Norman Guo^b Houping Xiao^c Baozhong Yang^c

^aUniversity of Maryland

^bSaint Louis University

^cGeorgia State University

Friday Seminar
April 2023

MOTIVATION

Human decisions

- Important to the economy and markets
- Difficult to analyze
 - Rational paradigm (e.g., Muth 1961; Lucas 1987)
 - Cognitive heuristics (e.g., Kahneman and Tversky 1979)

Machines have succeeded in many tasks

- Image recognition
- Natural language processing
- Game playing
- Automatic driving

Can machines understand and evaluate human skills?



MOTIVATION

AI Can Write a Song, but It Can't Beat the Market

Quants have tried for decades with limited success at their biggest challenge



Market data is smaller in size and “noisier” than language and other data, making it harder to use it to explain or predict market moves

RESEARCH OBJECTIVES

- ① Design a machine-friendly AI model for financial data.
- ② Explore how do machines evaluate analyst skills.
- ③ Extract valuable information from individual and collective analyst forecasts?.

RESEARCH OBJECTIVES

Why the analyst setting

- Analysts are important financial intermediaries
- Earnings forecasts are measurable individual opinions
- Observable features from analysts, firms and economy
- Past realized earnings can serve as benchmark for evaluation of performance
 - Manual labelling is labor-intensive
 - Learning from labels: Which analyst has information or is more skilled

Challenges

- Each analyst's private information and expertise
- High-dimensional, nonlinear interactions

WHY DO WE USE MACHINE LEARNING

Traditional Econometrics, e.g., OLS

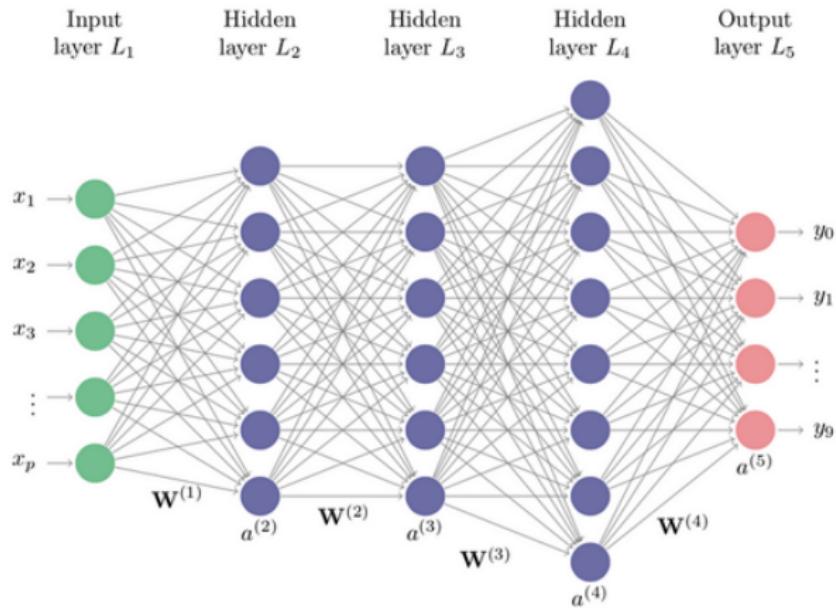
- Have difficulty dealing with a large number of variables
- Cannot handle complicated nonlinear relations
- Optimized for in-sample interpretation, not out-of-sample prediction

Machine Learning Methods, e.g., Neural Networks

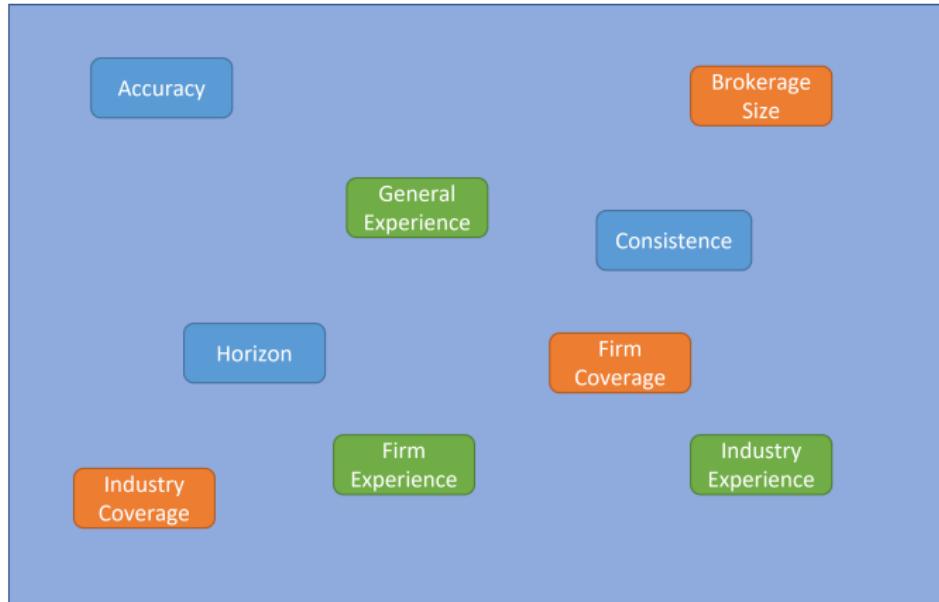
- Built-in dimension reduction to focus on more important variables
- Incorporate highly flexible nonlinear relations
- Model designs are optimized for out-of-sample predictions

EXAMPLE OF ML: FEED FORWARD NEURAL NETWORKS

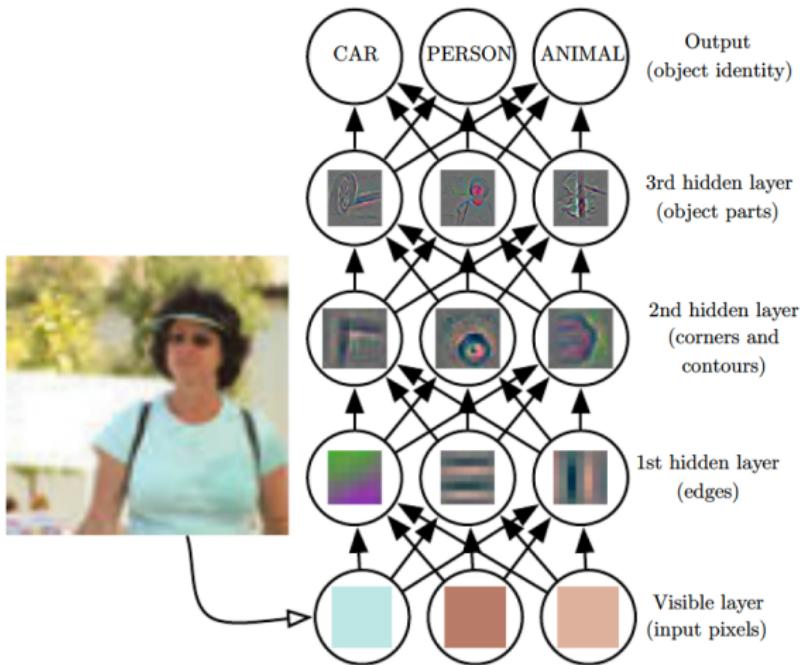
Neural Networks



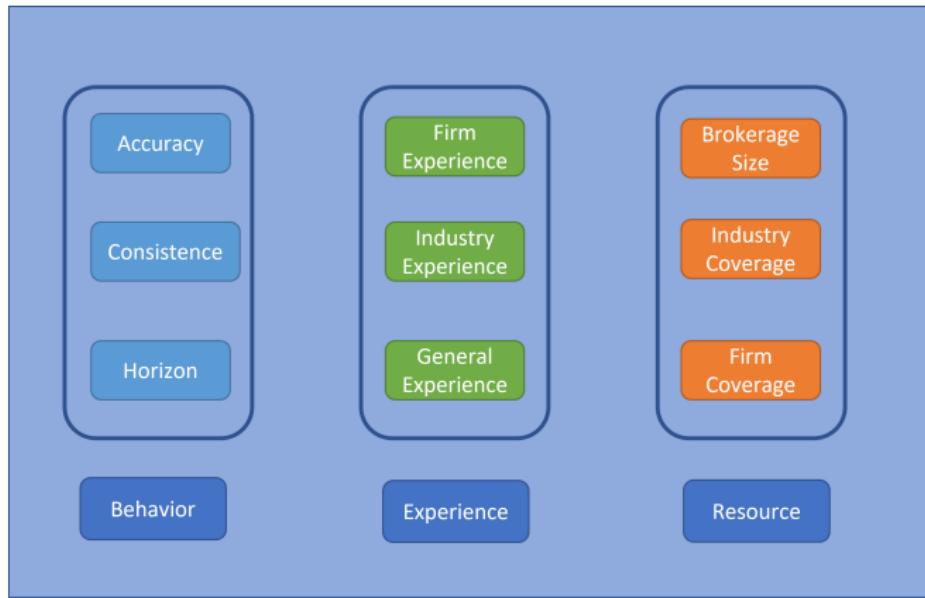
EXAMPLE OF ML: CONVOLUTIONAL NEURAL NETWORKS (CNN)



EXAMPLE OF ML: CONVOLUTIONAL NEURAL NETWORKS (CNN)



EXAMPLE OF ML: CONVOLUTIONAL NEURAL NETWORKS (CNN)



DATA AND FEATURES

Data Sources: IBES, Compustat, CRSP, Fed St. Louis, Thomson 13F, etc.

- **Analyst-level features, $A_{i,j,t}$**
 - 15 features
 - including Firm Experience, Forecast Horizon, Effort, Consensus (IBES), etc
- **Macro-level features, T_t**
 - 12 features
 - including Inflation, Oil Prices, Term Spread, Default Spread, VIX, etc
- **Firm-level features, $F_{j,t}$**
 - 40 features
 - including Size, Book to Market, Momentum, Accruals, Profit Margin, Asset Liquidity, Closed Price, Turnover, Institutional Ownership, etc

Note: analyst i , firm j , and time t

CONSTRUCT TARGET VARIABLE

$$\text{Star}_{i,j,t+1} = f(A_{i,j,t}, F_{j,t}, T_t) + \epsilon_{i,j,t+1}$$

for analyst i , firm j , and time t .

- **Classification:** $\text{Star}_{ijt} = 1$ if the absolute forecast error of analyst i in the quarter t is lower than median of all analysts covering the firm j ; otherwise $\text{Star}_{ijt} = 0$.

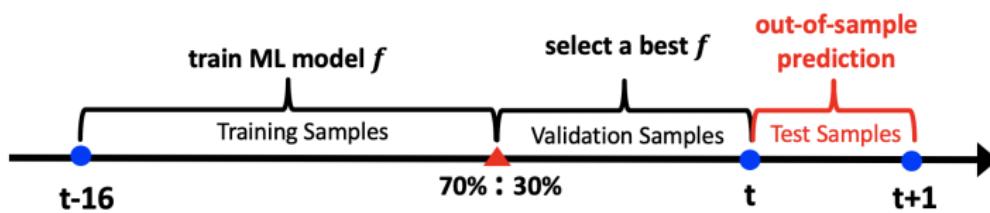
CONSTRUCT TARGET VARIABLE

Analysts	FE
John	0.15
Mary	0.2
Sarah	0.3
Lenard	-0.1
Brooke	0.5
Clifford	0.45
Emerson	-0.2
Olive	-0.4
Shelia	-0.35

CONSTRUCT TARGET VARIABLE

Analysts	Abs FE	Star
Lenard	0.1	1
John	0.15	1
Mary	0.2	1
Emerson	0.2	1
Sarah	0.3	0
Shelia	0.35	0
Olive	0.4	0
Clifford	0.45	0
Brooke	0.5	0

TIMELINE



EXPLANATION OF MACHINE LEARNING METRICS

Metrics used to evaluate ML models:

- **Accuracy:** $\frac{\text{True Positives}}{\text{Total Sample}}$
- **Precision:** $\frac{\text{True Positive}}{\text{True Positive} + \text{False Positive}}$, measures Type I error
- **Recall:** $\frac{\text{True Positive}}{\text{True Positive} + \text{False Negative}}$, measures Type II error
- **F1 Score:** Average of precision and recall (harmonic average)

MODEL COMPARISON

- Non-linear models outperform
- Convolutional Neural Networks (CNN), which proceeds from low-dimensional interactions of features to high-dimensional interactions, excels

	Models	Accuracy	Precision	Recall	F1 Score
Linear	Logistic Regression	53.81%	54.33%	90.23%	67.84%
	Logistic LASSO	55.49%	55.90%	82.93%	66.78%
Non-Linear	Gradient Boost	58.14%	57.70%	83.95%	68.40%
	Neural Network	62.27%	62.02%	75.72%	69.7%
	Convolutional Neural Network	69.03%	69.10%	77.04%	72.86%

FEATURE SELECTION

- Not always “the more the better”
- Analyst features are the most important ones

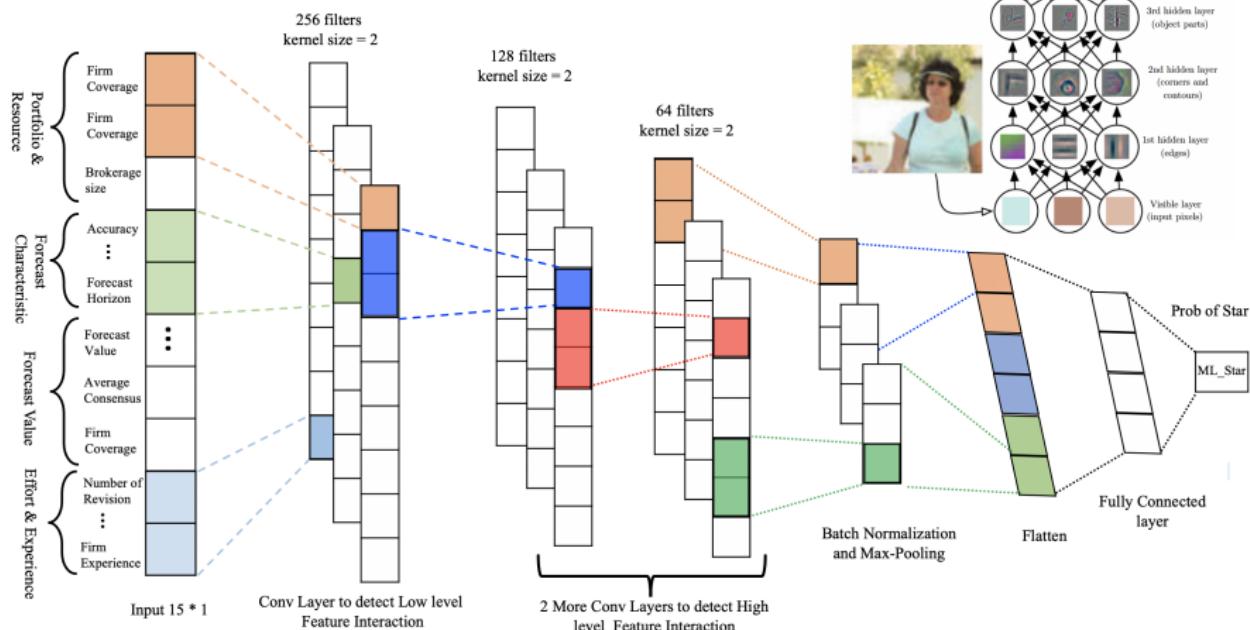
Feature	Accuracy	Precision	Recall	F1 Score
[Analyst, Firm, Macro]	68.60%	68.37%	77.45%	72.63%
[Analyst, Macro]	68.56%	68.53%	77.29%	72.65%
[Analyst, Firm]	68.62%	68.44%	77.30%	72.60%
[Firm, Macro]	53.93%	53.96%	99.66%	70.01%
[Analyst]	69.03%	69.10%	77.04%	72.86%
[Firm]	53.93%	53.94%	99.77%	70.03%
[Macro]	54.16%	54.16%	99.97%	70.26%

GROUPING ANALYST FEATURES

Analyst features come naturally in 4 groups

- **Forecast Values:** Forecast, Consensus from I/B/E/S, Average Consensus
- **Forecast Characteristics:** Accuracy, Consistency, Horizon
- **Effort & Experience:** Number of Revisions, Whether Report Revenue Forecast, Whether Report Cash flow Forecast, General Experience, Industry Experience, Firm Experience
- **Portfolio & Resource:** Analyst Firm Coverage, Analyst Industry Coverage, Brokerage Size

CNN ARCHITECTURE



COMPARISON WITH CASES WITH RANDOM ORDERS

- Random orders of variables do not work well with CNN
- Grouping of features are important!

Feature	Accuracy	Precision	Recall	F1 Score
<code>[FirmExperience, FirmCoverage, Accuracy, ReportCashflow, Consistency, IndustryExperience, BrokerageSize, IndustryCoverage, ForecastHorizon, ReportRevenue consensus_avg, ForecastValue, NumberofRevision, meanest_ibes, GeneralExperience]</code>	68.83%	71.58%	74.62%	73.05%
<code>[GeneralExperience, IndustryCoverage, FirmExperience, consensus_avg, ReportCashflow, Accuracy, ForecastValue, ReportRevenue, BrokerageSize, ForecastHorizon, IndustryExperience, meanest_ibes, NumberofRevision, Consistency, FirmCoverage]</code>	66.92%	69.02%	75.57%	72.13%
<code>[ForecastHorizon, NumberofRevision, Accuracy, ReportCashflow, ReportRevenue, consensus_avg, FirmExperience, IndustryCoverage, ForecastValue, meanest_ibes, Consistency, IndustryExperience, GeneralExperience, FirmCoverage, BrokerageSize]</code>	66.16%	66.84%	80.01%	72.83%
<code>[...]</code>				
<code>[ForecastValue, ReportCashflow, IndustryExperience, GeneralExperience, FirmCoverage, Consistency, consensus_avg, NumberofRevision, ForecastHorizon, IndustryCoverage, meanest_ibes, BrokerageSize, FirmExperience, ReportRevenue, Accuracy]</code>	57.14%	57.95%	89.00%	70.19%
<code>[ForecastHorizon, consensus_avg, GeneralExperience, FirmExperience, meanest_ibes, IndustryCoverage, IndustryExperience, ForecastValue, Consistency, Accuracy, ReportRevenue, NumberofRevision, FirmCoverage, BrokerageSize, ReportCashflow]</code>	57.04%	58.09%	87.22%	69.72%
<code>[ForecastHorizon, FirmCoverage, Accuracy, ReportRevenue, BrokerageSize, IndustryCoverage, ForecastValue, GeneralExperience, consensus_avg, Consistency, IndustryExperience, ReportCashflow, FirmExperience, NumberofRevision, meanest_ibes]</code>	56.83%	57.97%	86.87%	69.52%

FEATURE IMPORTANCE

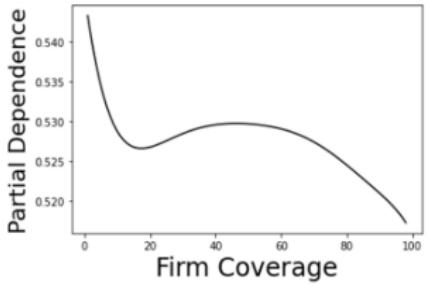
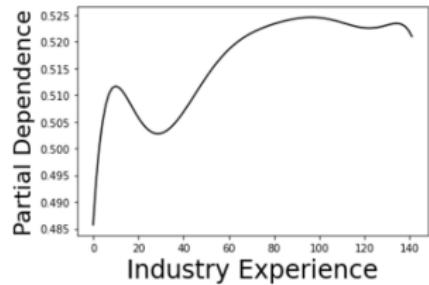
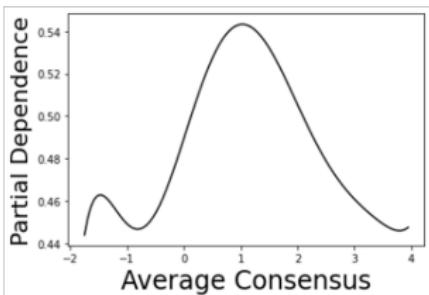
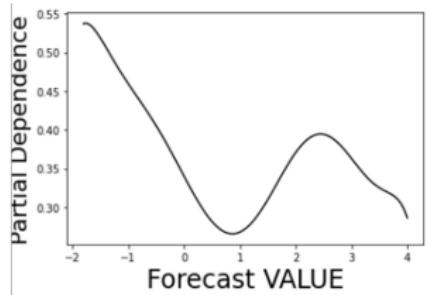
Importance from Linear Regression



Importance from Shapely Values



PARTIAL DEPENDENCE



FORECAST ACCURACY

- ML predicted star analysts outperform historically accurate analysts and (human-labeled) all star analysts

Variables	(1)	(2)	(3) <i>Star</i>	(4)
<i>ML-Star</i>	0.381*** (123.66)	0.382*** (81.42)	0.380*** (123.65)	0.380*** (81.02)
<i>Historical Star</i>			0.018*** (19.72)	0.018*** (16.41)
<i>All-Star</i>			0.009*** (6.29)	0.006*** (3.06)
Quarter FE	No	Yes	No	Yes
Firm FE	Yes	Yes	Yes	Yes
Observations	1,488,430	1,488,430	1,488,430	1,488,430
R-squared	0.145	0.145	0.145	0.145

FORECAST PERSISTENCE

- The predictive power of the ML-Star is persistent

Variables	(1)	(2)	(3)	(4)	(5)
	<i>Star</i>				
	1 Qtr	2 Qtr	3 Qtr	4 Qtr	8 Qtr
<i>ML-Star</i>	0.056*** (29.43)	0.042*** (25.39)	0.038*** (24.08)	0.035*** (20.95)	0.025*** (17.03)
<i>Historical Star</i>	0.036*** (23.98)	0.032*** (20.21)	0.028*** (18.99)	0.026*** (18.14)	0.022*** (13.11)
<i>All-Star</i>	-0.001 (-0.23)	-0.003 (-0.99)	-0.006 (-1.55)	-0.007 (-1.63)	-0.008* (-1.90)
Time FE	Yes	Yes	Yes	Yes	Yes
Fund FE	Yes	Yes	Yes	Yes	Yes
Observations	1,308,358	1,172,893	1,054,670	951,168	640,661
R-squared	0.014	0.013	0.013	0.013	0.014

FORECAST ACCURACY: SUBSAMPLE ANALYSIS

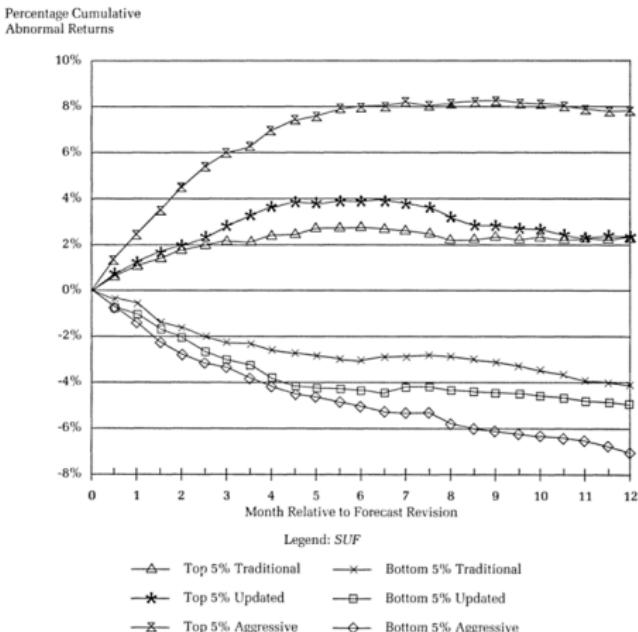
Analyst skill can be more accurately predicted by machines when

- Firm information is more transparent
- Analyst is more experienced, more focused, and has more resources
- The economy is in a normal state

	High	Low	Diff	t-Stat
Information Asymmetry				
Bid Ask Spread	0.364***	0.415***	-0.051	(-5.70)
Adj probability of informed trading	0.389***	0.431***	-0.042	(-3.08)
Return Volatility	0.362***	0.417***	-0.055	(-6.64)
Cashflow Volatility	0.365***	0.403***	-0.038	(-6.15)
Earning Quality	0.416***	0.384***	0.032	(4.49)
Firm Age	0.397***	0.380***	0.017	(2.56)
Analyst Characteristics				
General Experience	0.395***	0.383***	0.012	(4.70)
Analyst Firm Coverage	0.385***	0.393***	-0.008	(-2.86)
Brokerage Size	0.395***	0.383***	0.012	(4.49)
Market Condition				
NBER Crisis Dummy	0.366***	0.393***	-0.027	(-1.96)

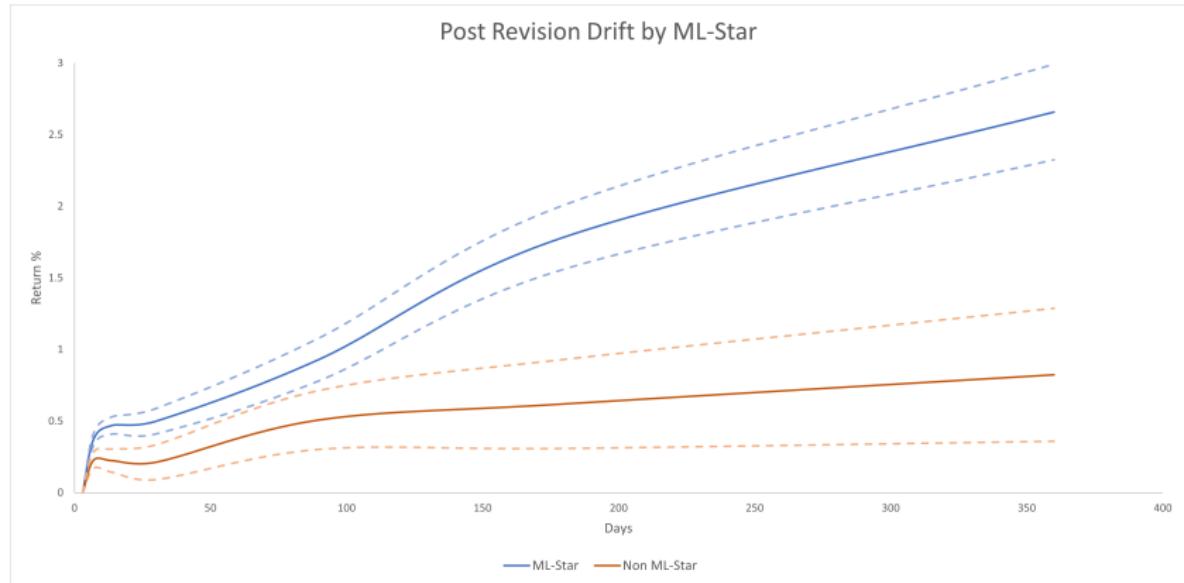
POST REVISION DRIFT

- The phenomenon of delayed stock price reactions to analyst forecast revisions, is a well-documented market anomaly.



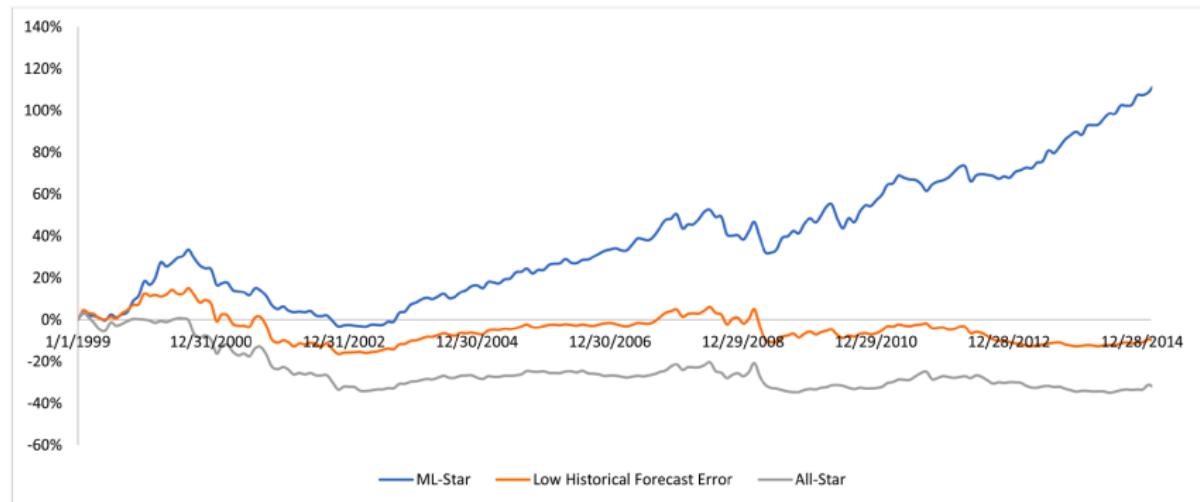
POST REVISION DRIFT

- ML predicted star analysts explain the bulk of post analyst revision drifts.



TRADING STRATEGY RETURNS: POST ANALYST REVISION DRIFT

Long positive revision, short negative revision



GENERATE CLOUD WISDOM

- **ML Earnings Consensus:** We compute the ML consensus as the average of predicted ML-Star analysts' forecasts

EARNING FORECAST

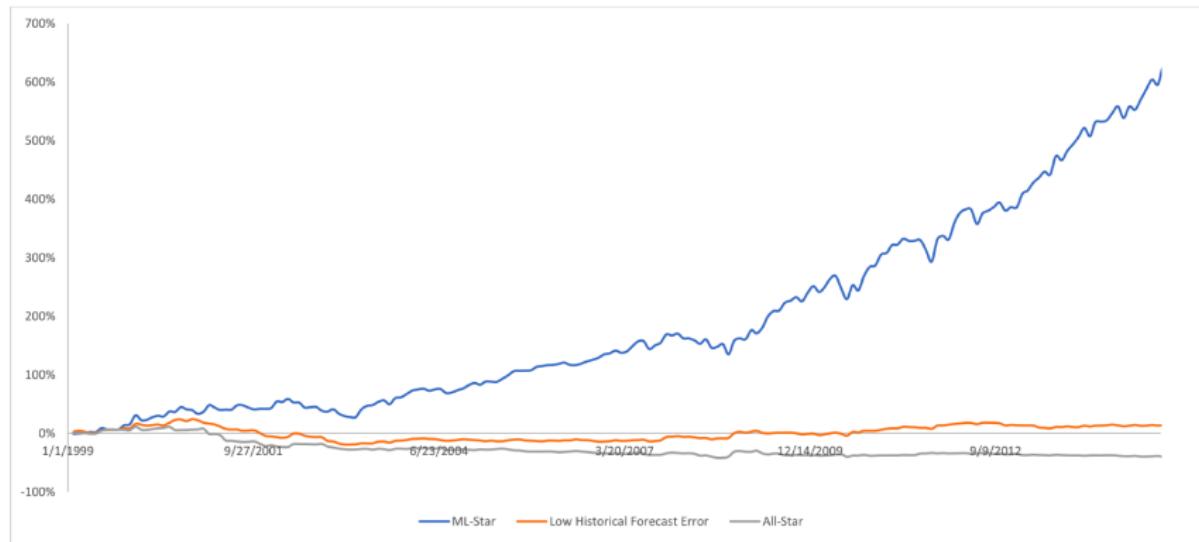
Dependent Variable	(1)	(2)	(3)	(4)
	<i>Earnings</i>			
<i>Consensus_ML - Consensus</i>	2.133*** (4.25)	2.142*** (4.32)	2.058*** (4.73)	2.034*** (4.93)
<i>Consensus</i>	1.064*** (41.69)	1.059*** (48.97)	1.112*** (22.46)	1.091*** (25.18)
<i>Liquidity</i>	-0.003 (-1.32)		0.004** (2.55)	
<i>Momentum</i>	0.030*** (6.32)		0.012** (2.13)	
<i>Log_Size</i>	-0.004 (-0.80)		-0.007 (-0.73)	
<i>Book to Market</i>	-0.010* (-1.97)		0.028* (1.90)	
<i>Coverage</i>	0.001** (2.24)		-0.001** (-2.23)	
Quarter FE	Yes	Yes	Yes	Yes
Fund FE	No	No	Yes	Yes
Observations	156,635	203,759	156,158	203,118
Adj R-squared	0.790	0.771	0.808	0.791

MARKET EXPECTATION

- ML consensus predicts returns around earnings announcements

Variables	(1) CAR [-1, +1]	(2) CAR [+2, +7]	(3) CAR [+8, +14]
<i>Consensus_ML - Consensus</i>	0.019** (2.60)	-0.001 (-0.16)	-0.003 (-0.51)
<i>Consensus</i>	0.002*** (2.70)	0.003*** (3.02)	0.002*** (3.46)
<i>Liquidity</i>	-0.001 (-1.46)	-0.001* (-1.96)	0.000 (0.54)
<i>Momentum</i>	0.001 (0.86)	-0.004** (-1.99)	-0.002 (-1.24)
<i>Log_Size</i>	-0.012*** (-11.45)	-0.007*** (-5.98)	-0.007*** (-6.86)
<i>Book to Market</i>	-0.001 (-1.26)	-0.001 (-0.84)	-0.002* (-1.82)
<i>Coverage</i>	-0.000 (-0.81)	0.000 (0.82)	0.000 (0.06)
Quarter FE	Yes	Yes	Yes
Fund FE	Yes	Yes	Yes
Observations	154,783	154,767	154,662
adj R-squared	0.0293	0.0410	0.0446

TRADING STRATEGY RETURNS: POST EARNING DRIFT



CONCLUSION

- **A ML measure of analyst skill**

- A persistent skill measure that outperforms human-labeled star analysts and historically accurate analysts in future analyst forecasts
- Explains the post-revision drift anomaly for analysts
- Skill prediction is more accurate in a transparent information environment

- **A new earnings expectation measure from ML analyst consensus**

- Better predicts earnings surprise
- Predicts stock returns around earnings announcements
- Generates profitable trading strategies for investors
- AI provides significant incremental information to common consensus

- **Methodological contribution**

- Feature and model selection in Machine Learning
- CNN can capture subtle variable interactions by grouping and ordering of features
- Interpretation of non-linear relations in deep-learning models
- A new ML method to aggregate information from heterogeneous agents: Applicable to general settings, e.g., online forums, political opinions, and macroeconomic outlooks