

HUMAN SKILL THROUGH A DIGITAL LENS: EVALUATING ANALYSTS WITH MACHINE LEARNING

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

^aUniversity of Maryland

^bSaint Louis University

^cGeorgia State University

University of Louisville
Feb 2024

MOTIVATION

Human skills

- 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?

OVERVIEW

- ➊ Design a human-friendly AI model for financial data.
 - Integrate domain knowledge.
 - Facilitate local non-linear interactions.
 - Convert images to numerical data.
- ➋ Evaluate analysts' skills from a machine perspective.
 - Machine vs. human assess human skills differently in important dimensions
 - Answer the puzzle of post-revision drift
- ➌ Extract valuable information from individual and collective analyst forecasts.
 - Generate significant abnormal returns from machine-selected analysts
 - Create a “smart” analyst consensus that better proxies for earning news before earnings announcements than the traditional analyst consensus

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

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 $\text{Star}_{ijt} = 0$.

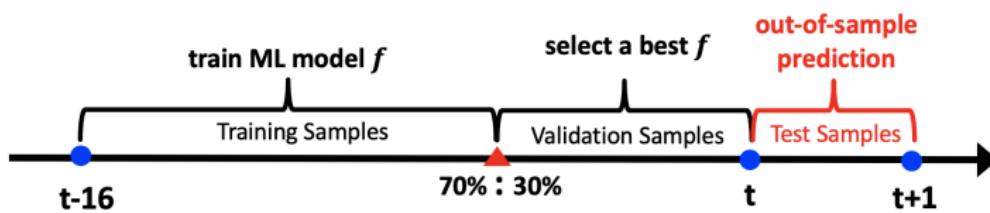
CONSTRUCT TARGET VARIABLE

| 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



ISSUES OF MACHINE LEARNING MODELS IN FINANCE

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

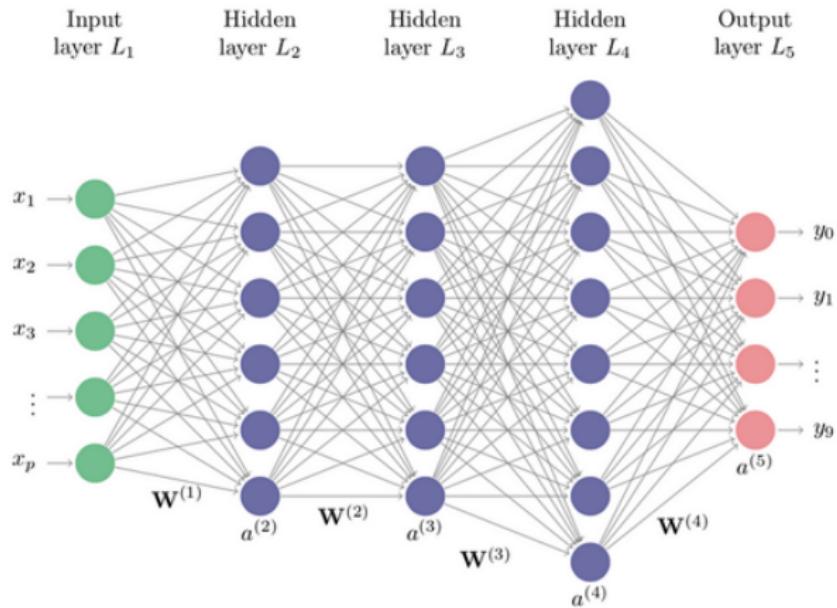
Quants have tried for decades with limited success at their biggest challenge



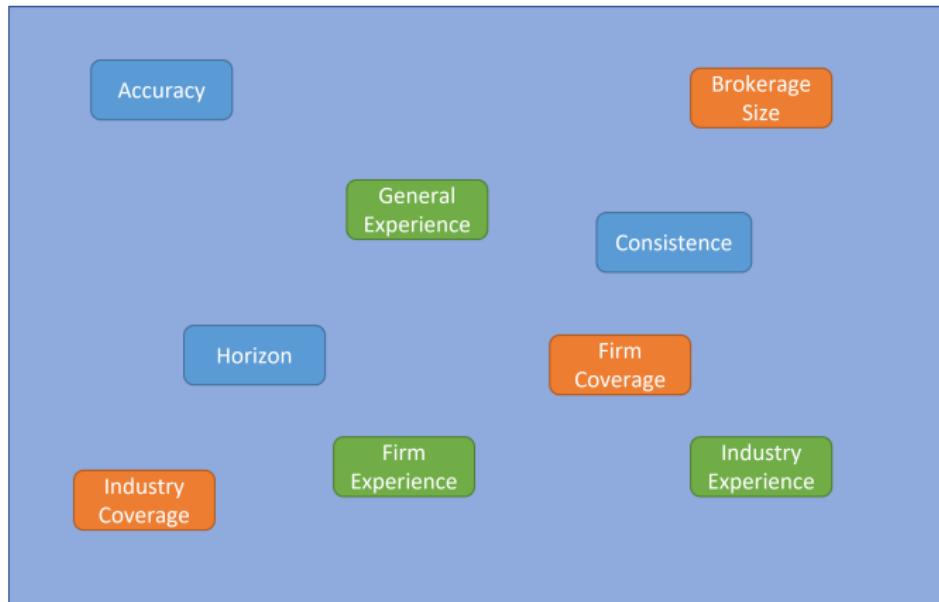
Market data is smaller in size and “noisier” than language and other data, making it harder to use it to explain or predict market moves

EXAMPLE OF ML: FEED FORWARD NEURAL NETWORKS

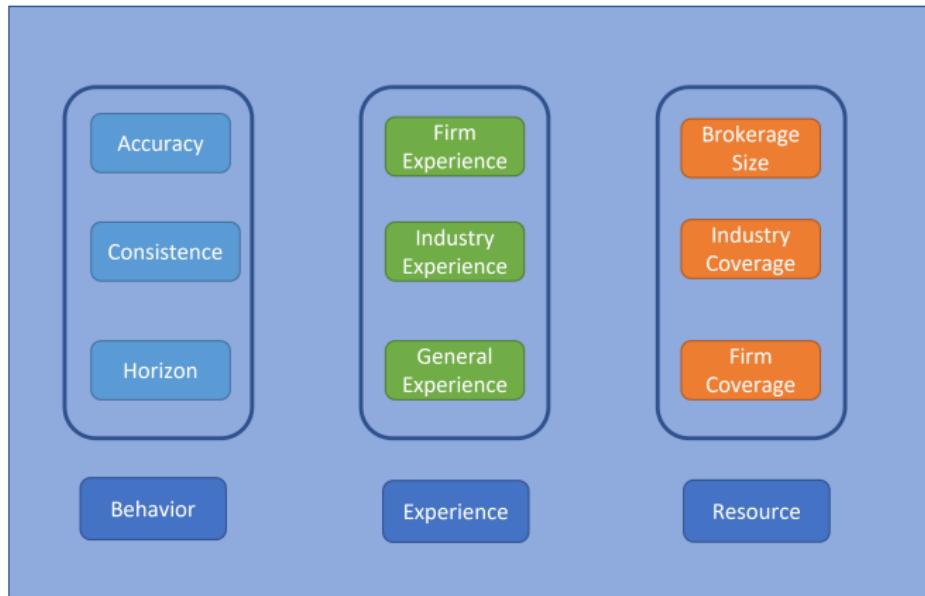
Neural Networks



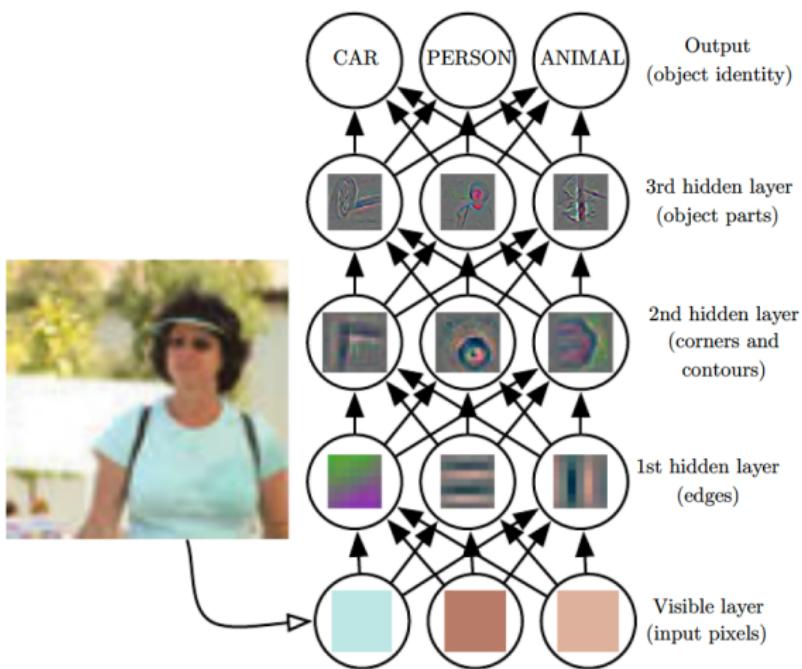
EXAMPLE OF ML: CONVOLUTIONAL NEURAL NETWORKS (CNN)



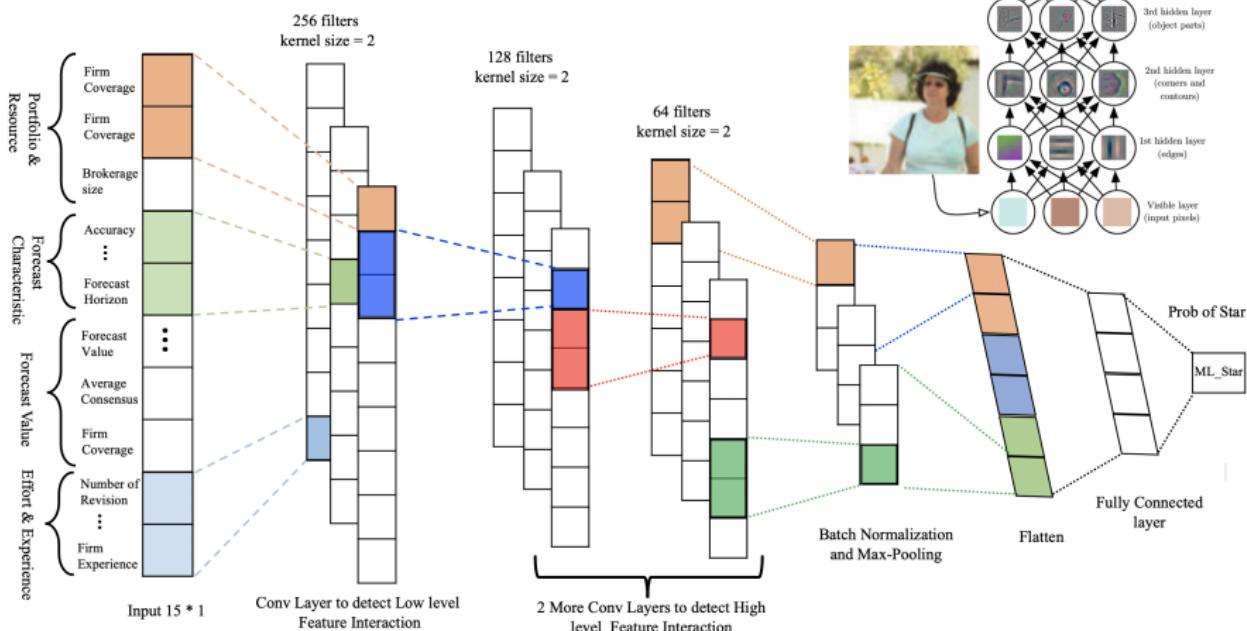
EXAMPLE OF ML: CONVOLUTIONAL NEURAL NETWORKS (CNN)



EXAMPLE OF ML: CONVOLUTIONAL NEURAL NETWORKS (CNN)



CNN ARCHITECTURE



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)

FEATURE SELECTION

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

| Feature | Accuracy | Precision | Recall | F1 Score |
|------------------------|----------|-----------|---------|----------|
| [Analyst] | 61.57% | 62.89% | 83.06% | 71.10% |
| [Analyst, Firm] | 60.80% | 61.60% | 83.33% | 70.49% |
| [Analyst, Macro] | 60.30% | 60.57% | 87.70% | 71.26% |
| [Analyst, Firm, Macro] | 59.49% | 60.74% | 82.92% | 69.85% |
| [Firm] | 56.33% | 56.38% | 99.68% | 72.02% |
| [Macro] | 56.38% | 56.38% | 100.00% | 72.10% |
| [Firm, Macro] | 56.29% | 56.35% | 99.66% | 71.99% |

BRINGING IN DOMAIN KNOWLEDGE

Group Analyst features in four categories based on literature:

- **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

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% |

IMPORTANCE OF THE ORDER OF FEATURE CATEGORIES IN CNN



- The order of different feature groups matters

| Feature | Accuracy | Precision | Recall | F1 Score |
|--|----------|-----------|--------|----------|
| [Portfolio&Resource, Effort&Experience, ForecastChar, ForecastValue] | 69.79% | 71.83% | 76.35% | 74.02% |
| [Portfolio&Resource, ForecastChar, Effort&Experience, ForecastValue] | 69.58% | 70.23% | 79.92% | 74.76% |
| [Effort&Experience, ForecastValue, ForecastChar, Portfolio&Resource] | 69.57% | 71.58% | 76.30% | 73.87% |
| ... | | | | |
| [ForecastChar, ForecastValue, Portfolio&Resource, Effort&Experience] | 68.61% | 69.49% | 79.01% | 73.94% |
| [ForecastChar, Effort&Experience, ForecastValue, Portfolio&Resource] | 68.00% | 69.33% | 77.56% | 73.21% |
| [Effort&Experience, ForecastValue, Portfolio&Resource, ForecastChar] | 67.27% | 67.10% | 82.27% | 73.92% |

ORDER WITHIN EACH FEATURE CATEGORY



- The order of different variables within each group matters

| Features | Accuracy | Precision | Recall | F1 Score |
|--|----------|-----------|--------|----------|
| <code>IndustryCoverage, FirmCoverage, BrokerageSize, Accuracy, ForecastHorizon, Consistency, meanest.ibes, consensus_avg, ForecastValue, ReportCashflow, GeneralExperience, IndustryExperience, ReportRevenue, NumberofRevision, FirmExperience</code> | 70.33% | 71.92% | 78.33% | 74.97% |
| <code>FirmCoverage, BrokerageSize, IndustryCoverage, Consistency, ForecastHorizon, Accuracy, consensus_avg, ForecastValue, meanest.ibes, FirmExperience, GeneralExperience, ReportRevenue, ReportCashflow, IndustryExperience, NumberofRevision</code> | 70.06% | 71.64% | 78.15% | 74.74% |
| <code>BrokerageSize, FirmCoverage, IndustryCoverage, Accuracy, Consistency, ForecastHorizon, ForecastValue, meanest.ibes, consensus_avg, FirmExperience, IndustryExperience, NumberofRevision, ReportCashflow, ReportRevenue, GeneralExperience</code> | 69.96% | 71.55% | 78.09% | 74.66% |
| *** | | | | |
| <code>FirmCoverage, IndustryCoverage, BrokerageSize, Accuracy, ForecastHorizon, Consistency, meanest.ibes, consensus_avg, ForecastValue, ReportCashflow, ReportRevenue, IndustryExperience, NumberofRevision, FirmExperience, GeneralExperience</code> | 69.55% | 71.32% | 77.50% | 74.20% |
| <code>IndustryCoverage, BrokerageSize, FirmCoverage, ForecastHorizon, Accuracy, Consistency, consensus_avg, meanest.ibes, ForecastValue, NumberofRevision, IndustryExperience, ReportRevenue, FirmExperience, ReportCashflow, GeneralExperience</code> | 69.49% | 71.52% | 76.75% | 74.04% |
| <code>FirmCoverage, BrokerageSize, IndustryCoverage, Consistency, ForecastHorizon, Accuracy, meanest.ibes, consensus_avg, ForecastValue, ReportRevenue, NumberofRevision, ReportCashflow, GeneralExperience, IndustryExperience, FirmExperience</code> | 69.43% | 71.09% | 77.62% | 74.20% |

MODEL COMPARISON

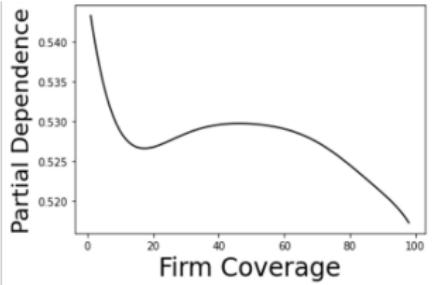
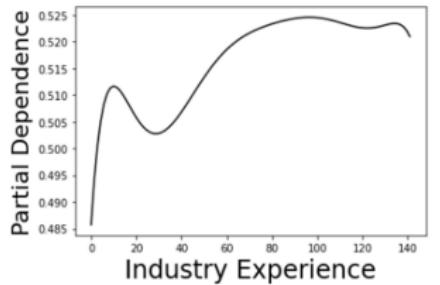
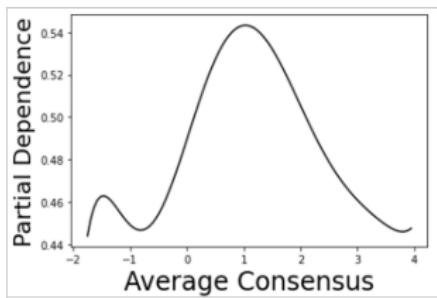
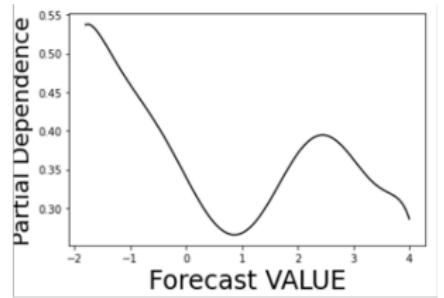
- Non-linear models outperform
- Convolutional Neural Networks (CNN), which proceeds from low-dimensional interactions of features to high-dimensional interactions, excels

| | 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 | 59.81% | 58.16% | 90.86% | 70.93% |
| | Convolutional Neural Network | 70.33% | 71.92% | 78.33% | 74.97% |

FEATURE IMPORTANCE

ML-Star*All-Star*

PARTIAL DEPENDENCE



FORECAST ACCURACY

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

| Variables | (1) | (2) | (3) | (4) |
|-------------------|----------------------|---------------------|----------------------|---------------------|
| | <i>Star</i> | | | |
| <i>ML-Star</i> | 0.381*** (123.66) | 0.382*** (81.42) | 0.380*** (123.65) | 0.380*** (81.02) |
| <i>Prior Star</i> | | | 0.018*** (19.72) | 0.018*** (16.41) |
| <i>All-Star</i> | | | 0.009*** (6.29) | 0.006*** (3.06) |
| Year-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>Prior 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) |
| Year-Quarter 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

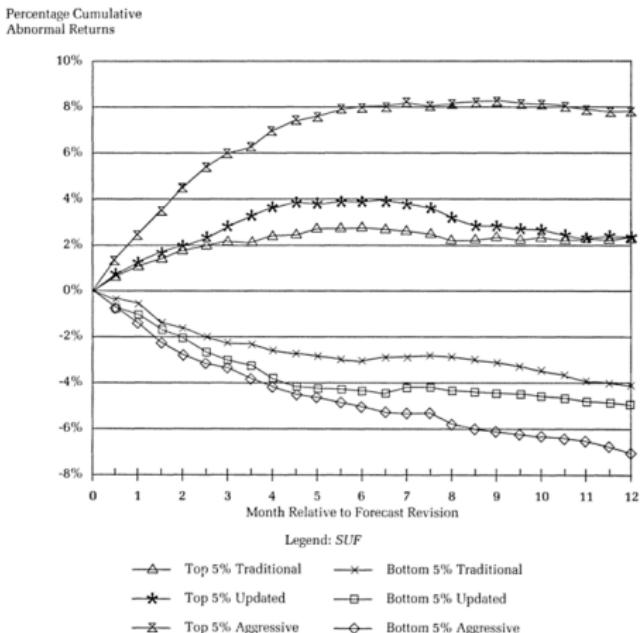
Analyst skill can be more accurately predicted by machines when

- Firm information is more transparent
- The economy is in a normal state

| | <i>ML-Star on Analyst Forecast Accuracy</i> | | | |
|--|---|----------|-----------|---------|
| | High | Low | Diff | t-Stat |
| <i>Bid Ask Spread</i> | 0.364*** | 0.415*** | -0.051*** | (-5.70) |
| <i>Adj probability of informed trading</i> | 0.389*** | 0.431*** | -0.042*** | (-3.08) |
| <i>Flesch-Kincaid Grade Level</i> | 0.398*** | 0.383*** | 0.015*** | (2.28) |
| <i>Accruals Quarlity</i> | 0.386*** | 0.409*** | -0.023*** | (-2.64) |
| <i>Earning Quality</i> | 0.416*** | 0.384*** | 0.032*** | (4.49) |
| <i>Cashflow Volatility</i> | 0.365*** | 0.403*** | -0.038*** | (-6.15) |
| <i>Return Volatility</i> | 0.362*** | 0.417*** | -0.055*** | (-6.64) |
| <i>Firm Age</i> | 0.397*** | 0.380*** | 0.017*** | (2.56) |
| <i>NBER Crisis Dummy</i> | 0.366*** | 0.393*** | -0.027** | (-1.96) |

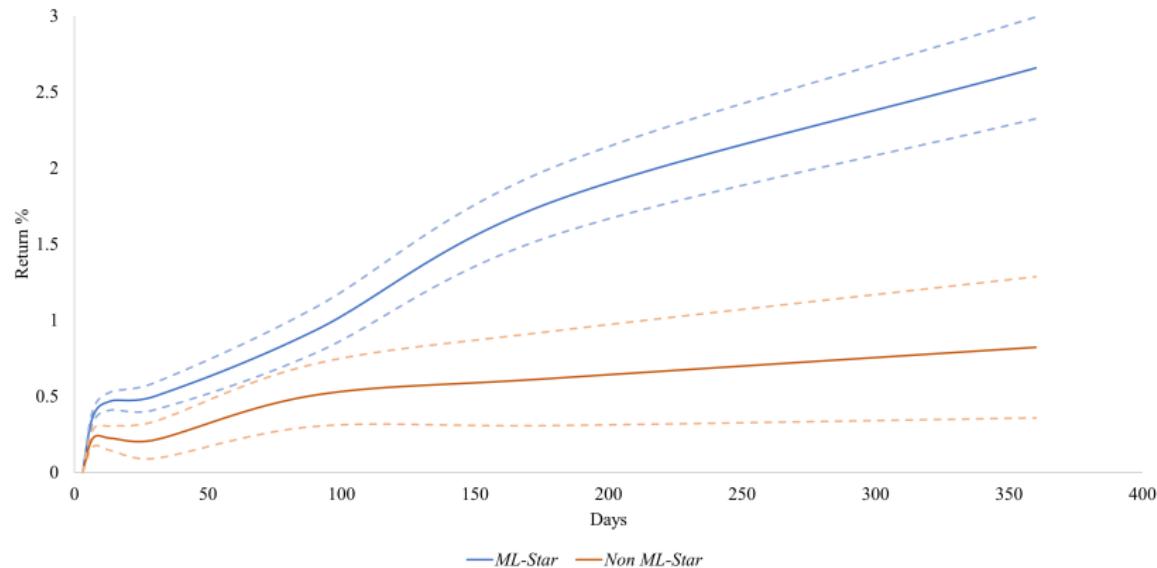
POST REVISION DRIFT

- The phenomenon of delayed stock price reactions to analyst forecast revisions, is a well-documented market anomaly. (Stickel (1999))



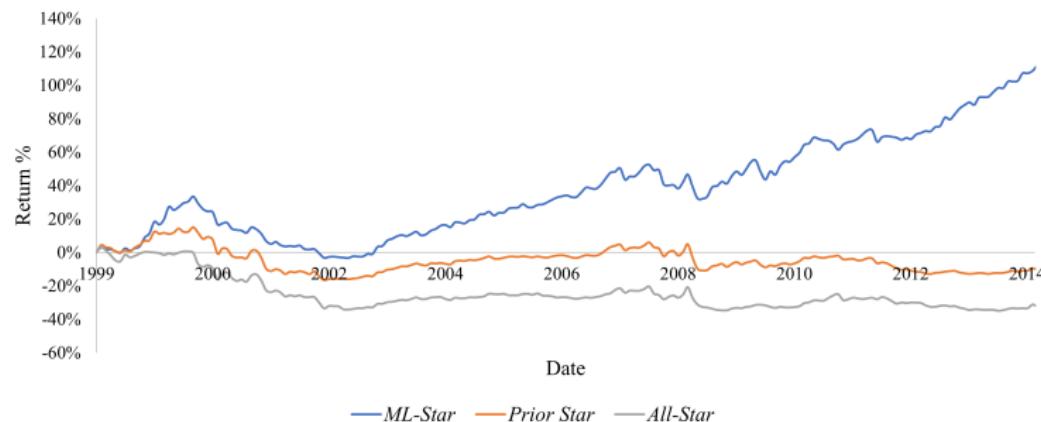
POST REVISION DRIFT

- 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

- ML consensus provide additional predicting future earnings

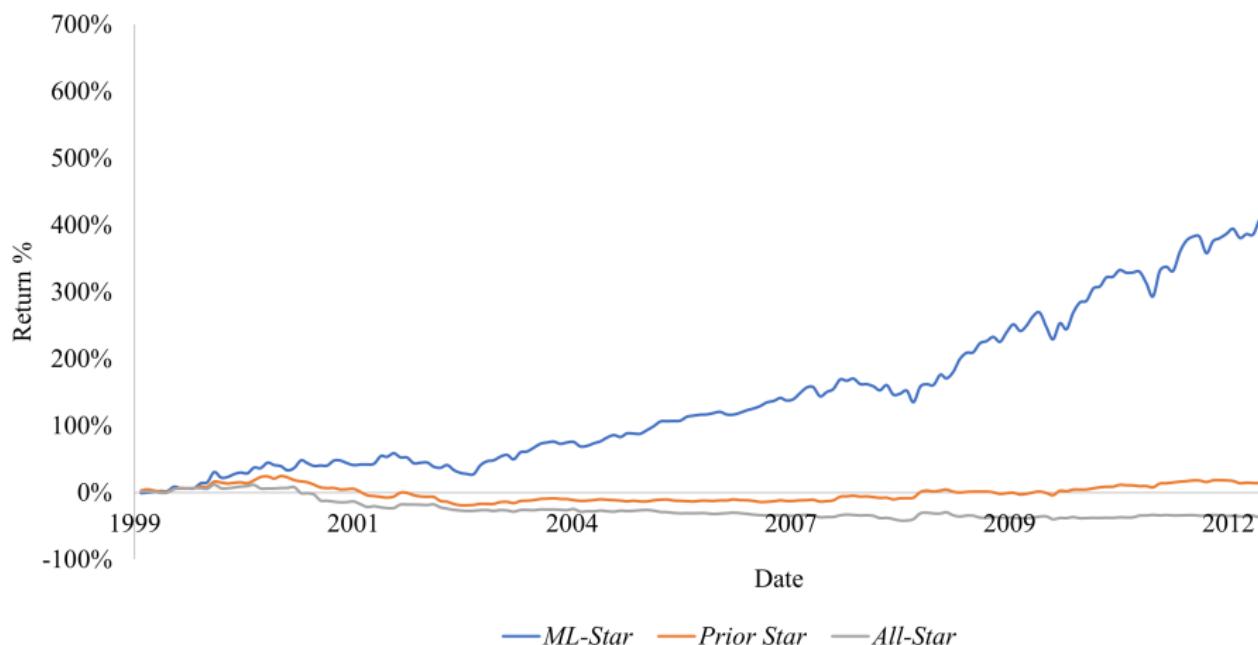
| Dependent Variable | (1) | (2) | (3) | (4) |
|---------------------------------|---------------------|---------------------|---------------------|---------------------|
| | Earnings | | | |
| <i>ML-Consensus - Consensus</i> | 2.142*** (4.32) | 2.034*** (4.93) | 2.133*** (4.25) | 2.058*** (4.73) |
| <i>Consensus</i> | 1.059*** (48.97) | 1.091*** (25.18) | 1.064*** (41.69) | 1.112*** (22.46) |
| <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) |
| Year-Quarter FE | Yes | Yes | Yes | Yes |
| Firm FE | No | Yes | No | Yes |
| Observations | 203,759 | 203,118 | 156,635 | 156,158 |
| Adj. R-squared | 0.771 | 0.791 | 0.790 | 0.808 |

MARKET EXPECTATION

- ML consensus predicts returns around earnings announcements

| Variables | (1) CAR [-1, 1] | (2) CAR [2, 7] | (3) CAR [8, 14] |
|---------------------------------|-----------------------|----------------------|----------------------|
| <i>ML-Consensus - 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) |
| Year-Quarter FE | Yes | Yes | Yes |
| Firm 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



FACTOR ANALYSIS

- Machine learning-based strategy captures unique insights not fully integrated into the market

| Variables | (1) | (2) | (3) | (4) | (5) |
|-----------------|--------------------------------|---------------------|----------------------|---------------------|---------------------|
| | Algorithm Return Based on PEAD | | | | |
| <i>Mkt-RF</i> | 0.479*** (13.60) | 0.491*** (12.93) | 0.534*** (12.75) | 0.495*** (11.57) | 0.544*** (11.59) |
| <i>SMB</i> | 0.004 (0.08) | -0.002 (-0.04) | 0.031 (0.58) | 0.018 (0.35) | 0.025 (0.48) |
| <i>HML</i> | -0.113** (-2.46) | -0.103** (-2.16) | -0.230*** (-3.50) | | |
| <i>Mom</i> | | 0.025 (0.84) | | | |
| <i>RMW</i> | | | 0.122* (1.69) | | |
| <i>CMA</i> | | | 0.179* (1.95) | | |
| <i>R_IA</i> | | | | -0.018 (-0.24) | |
| <i>R_ROE</i> | | | | 0.059 (0.93) | |
| <i>MGMT</i> | | | | | 0.018 (0.31) |
| <i>PERF</i> | | | | | 0.095** (2.52) |
| <i>Constant</i> | 0.007*** (4.62) | 0.007*** (4.50) | 0.006*** (3.79) | 0.007*** (4.17) | 0.006*** (3.24) |
| Observations | 236 | 236 | 236 | 236 | 216 |
| adj R-squared | 0.468 | 0.468 | 0.478 | 0.440 | 0.460 |

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