



Group 18

IPO QUANTITATIVE ANALYSIS AND LISTING PREDICTOR

AKSHITHA, DIVY GUPTA, MOKSHITH, NAGA VENI

INDEX

1. Problem Statement

2. Literature Survey

3. Dataset and Features Preprocessing

4. ML Methodology

5. Performance Metrics &

Deployability of ML Solution

PROBLEM STATEMENT

“To predict IPO listing gains in the Indian stock market by analysing firm fundamentals together with investor demand and market sentiment”



RHP (Red Herring Prospectus):

RHP is an official document released by a company before launching its IPO that provides important information to investors.

GMP (Grey Market Premium) :

GMP refers to the unofficial price at which IPO shares trade in the grey market before the stock is listed on the exchange.



WHY IT MATTERS?

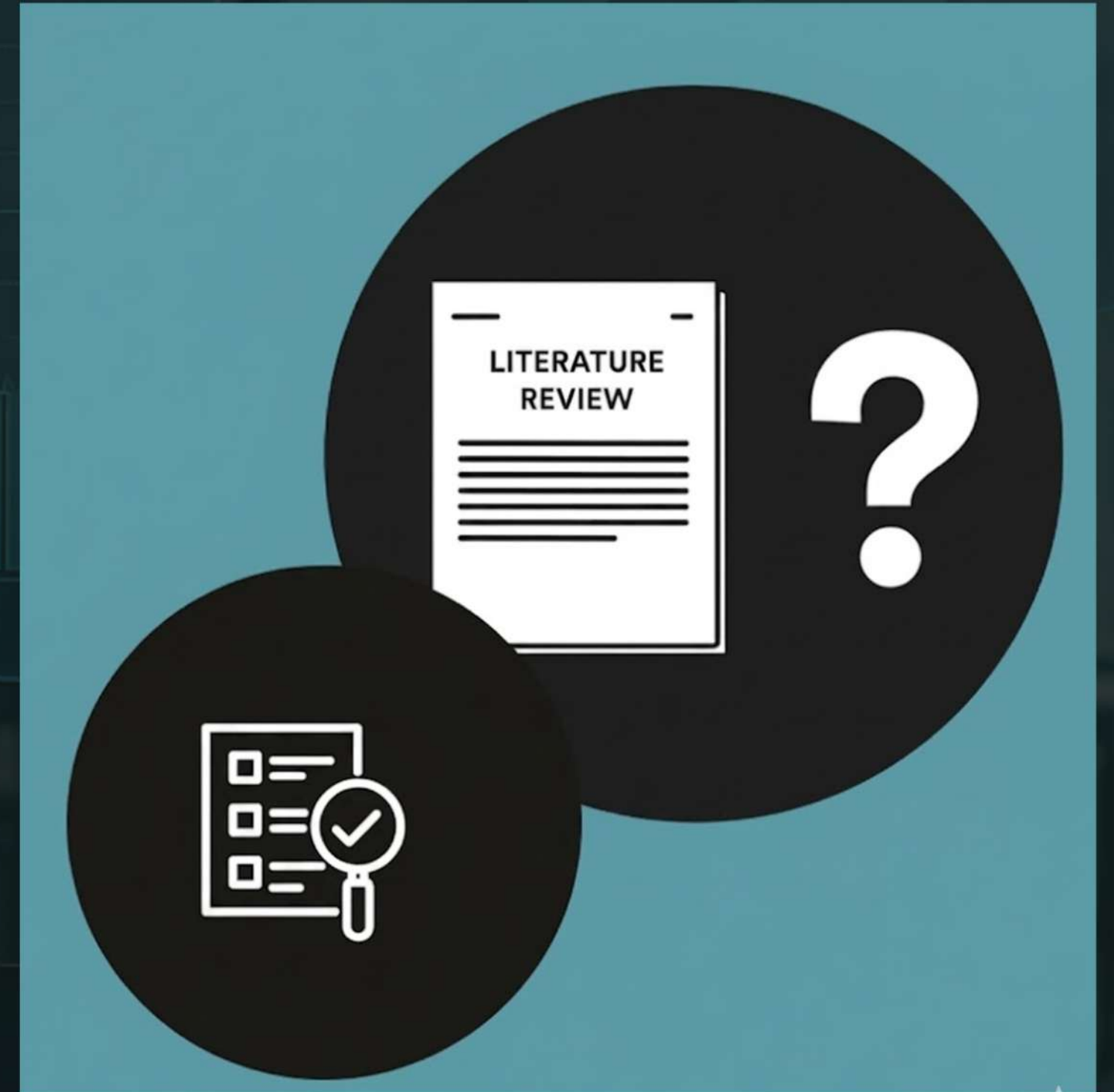
- To improve market efficiency , capital allocation
- Retail investors often rely on hype rather than analysis
- Reduces market mispricing
- Generally IPO markets involve high uncertainty - if prospectus language can predict performance, it reduces uncertainty and improves decision making(words influences the investor behaviour)





Group 18

LITERATURE SURVEY



PHASE 1 AND PHASE 2

PHASE	AUTHOR & YEAR	PAPER CITATIONS	KEY FINDINGS	KEY VALUE
PHASE 1 Underpricing (1990s – Early 2000s)	Ritter (1991)	<u>Ritter, J. R. (1991). The long-run performance of initial public offerings. <i>Journal of Finance</i>, 46(1), 3–27.</u>	IPOs give strong 1st-day returns but underperform long-run vs seasoned firms	Established listing-day prediction is worthwhile
	Ritter & Welch (2002)	<u>Ritter, J. R., & Welch, I. (2002). A review of IPO activity, pricing, and allocations. <i>Journal of Finance</i>, 57(4), 1795–1828.</u>	Standard pricing models (DCF, comparables) cannot explain IPO underpricing	Avg 1st-day return = 18.8% (US, 1980–2001)
PHASE 2 Behavioural & Sentiment (2000 – 2015)	Loughran & Ritter (2004)	<u>Loughran, T., & Ritter, J. R. (2004). Why has IPO underpricing changed over time? <i>Financial Management</i>, 33(3), 5–37.</u>	Psychological biases (anchoring, prospect theory) cause issuers to accept below-market prices	Rational pricing cannot explain underpricing
	Cornelli, Goldreich & Ljungqvist (2006)	<u>Cornelli, F., Goldreich, D., & Ljungqvist, A. (2006). Investor sentiment and pre-IPO markets. <i>Journal of Finance</i>, 61(3), 1187–1216.</u>	Grey market prices predict first-day returns high GMP = investor overexcitement, not fundamental value	GMP contradicts the Efficient Market Hypothesis
	Kim & Ritter (1999)	<u>Kim, M., & Ritter, J. R. (1999). Valuing IPOs. <i>Journal of Financial Economics</i>, 53(3), 409–437.</u>	Financial ratios from the prospectus (P/E, P/Sales, P/Book) explain cross-sectional IPO pricing	Validated RHP fundamentals as predictors
	Shetty et al. (2023)	<u>Shetty, C., Vinish, P., Aluru, S., Pinto, P., & Hawaldar, I. T. (2023). IPO subscription dynamics: A comprehensive inquiry into the Indian stock market. <i>Investment Management and Financial Innovations</i>, 20(4), 400–415.</u>	QIB subscription alone is the strongest predictor of Listing gains in India	QIB explains ~48% of listing-gain variance
	Lowry (2003)	<u>Lowry, M. (2003). Why does IPO volume fluctuate so much? <i>Journal of Financial Economics</i>, 67(1), 3–40.</u>	IPO volume fluctuates with investor sentiment and adverse selection costs	Justifies IPO_Year & market-timing features in model

PHASE 3

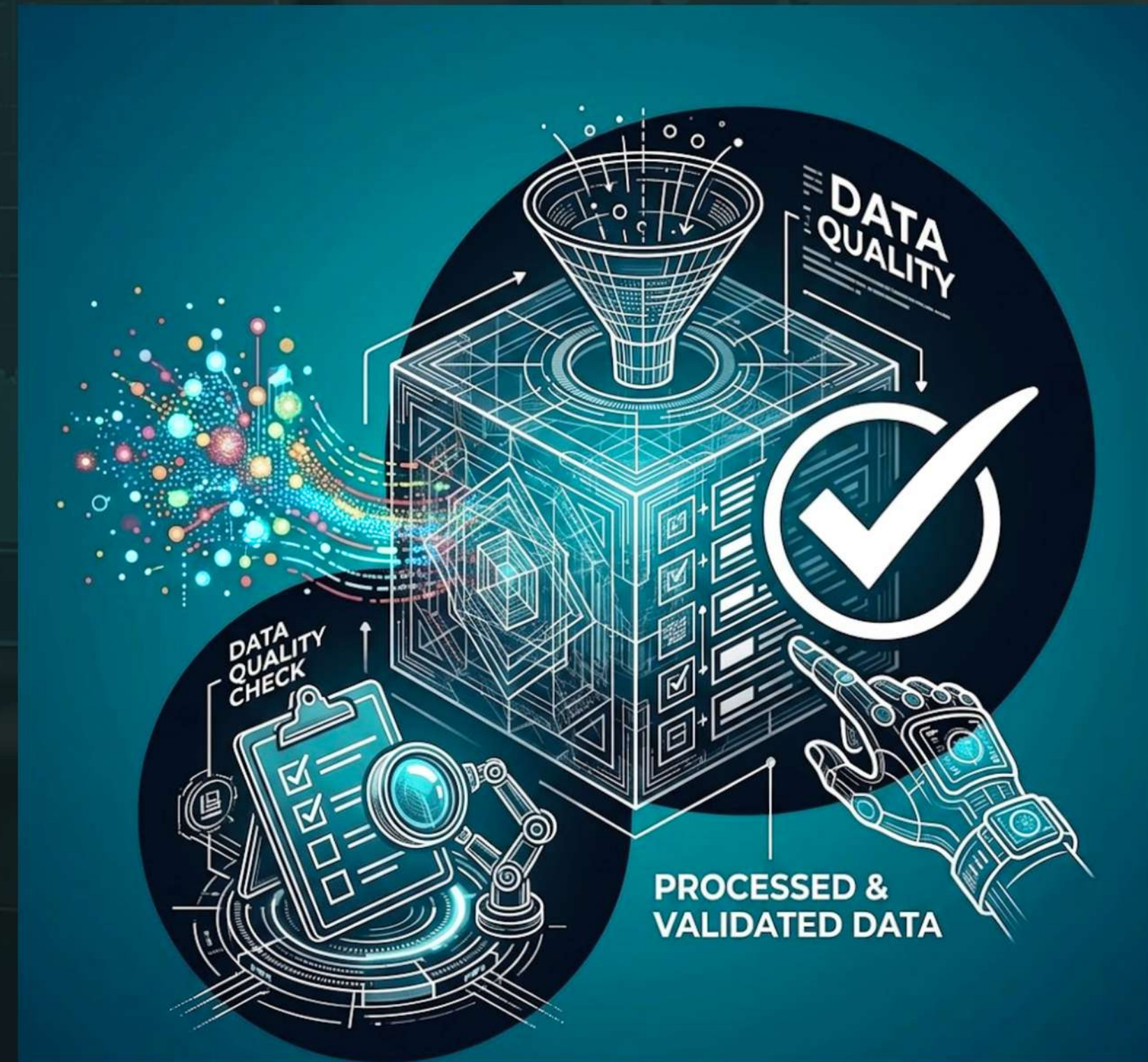
PHASE	AUTHOR & YEAR	PAPER CITATIONS	KEY FINDINGS	KEY VALUE
PHASE 3 Machine Learning Takes Over (2015 – Present)	Gu, Kelly & Xiu (2020)	Gu, S., Kelly, B., & Xiu, D. (2020). Empirical asset pricing via machine learning. Review of Financial Studies, 33(5), 2223–2273.	Tree models & neural nets capture non-linear interactions OLS misses entirely	+6.2% predictive R ² over OLS 0.4% Out of Sample R ²
	Baba & Sevil (2020)	Baba, B., & Sevil, G. (2020). Predicting IPO initial returns using random forest. Borsa Istanbul Review, 20(1), 13–23.	Random Forest outperforms OLS on 1st-day return prediction	RF RMSE 0.148 vs OLS RMSE 0.187
	Katsafados et al. (2023)	Katsafados, A. G., Leledakis, G. N., Pyrgiotakis, E. G., Androutsopoulos, I., Chalkidis, I., & Fergadiotis, E. (2023). Textual information and IPO underpricing: A machine learning approach. Journal of Financial Data Science, 5(2), 100–135.	Prospectus text + financials outperforms either alone; NLP adds measurable signal	Best accuracy 73.6% Text adds +6.1pp
	Alahmadi & Yilmaz (2025)	Alahmadi, M. F., & Yilmaz, M. T. (2025). Prediction of IPO performance from prospectus using multinomial logistic regression, a machine learning model. Data Science in Finance and Economics, 5(1), 105–135.	Fundamentals-only ML (no GMP) still performs well – confirms RHP is a rich signal source	F1 = 0.881 on average class
	Ghosh et al. (2024)	Ghosh, S., Maji, A., Vardhan, N. H., & Naskar, S. K. (2024). Experimenting with multi-modal information to predict success of Indian IPOs. arXiv preprint arXiv:2412.16174.	Multi-modal Indian IPO predictor: GMP + RHP text + macro data combined	GMP: 80.29% Main Board Only 21.29% SME

**THIS STUDY SYNTHESISES ALL THREE PHASES: ML + FUNDAMENTALS
+ GMP + NLP TEXT FOR THE INDIAN MARKET AT SCALE**



Group 18

DATASET & FEATURES PREPROCESSING



DATASET & DATA COLLECTION

1,620

Raw IPOs
collected

820

Final IPOs
with 47 features

2005–2024

Date range
BSE + NSE + SEBI

4

Data
sources merged

SEBI RHP Archive

Pre-issue financials, capital structure, promoter holdings (400–600 page PDFs)

Investorgain / Chittorgarh

Subscription multiples (NII, QIB, Retail) + Grey Market Premium (GMP)

NSE / BSE via Yahoo Finance

Listing price, post-issue share count, trading date

DuckDuckGo+Gemini

IPO listing dates missing from yfinance



DATA PREPROCESSING

PDF Extraction

Gemini 2.5 Flash reads Capital Structure sections from RHPs.
Algebraic derivation: Post-issue shares = Pre-issue + Fresh issue shares.
Listing Market Cap = Listing Price × Post-issue shares.

Sanitisation

Unit misreads capped at ₹50,000 Cr. PE ratios winsorised at 95th %ile. Listing gains above 1,000% clipped at 99th %ile (SME penny stocks).

Peer-Group Imputation

30–60% missingness in SME rows. Peers matched within ±30% anchor metric (revenue, PAT, total assets) in same sector.
Computed on train rows only — applied to cal/test sets via frozen medians.

Train / Cal / Test Split

Chronological split: 70% train | 10% calibration | 20% test.
All encodings, winsorisation bounds, and imputations frozen from training data to prevent look-ahead leakage.

Sector & Underwriter Encoding

Bayesian smoothing: sector encoded as $(n \cdot \mu_{\text{sector}} + k \cdot \mu_{\text{global}}) / (n+k)$, $k=10$. Lead manager 'pop rate' encoded the same way.
Unseen categories default to global training mean.

FEATURE ENGINEERING

Derived Financial Features

- $\text{Fresh_Issue_Shares} = \text{Fresh_Issue_Size} \div \text{Issue_Price}$
- $\text{Post_Issue_Shares} = \text{Pre_Issue_Shares} + \text{Fresh_Issue_Shares}$
- $\text{Listing_Market_Cap} = \text{Listing_Price} \times \text{Post_Issue_Shares}$
- $\text{Free_Float_Pct} = 100 - \text{Promoter_Holding_Post}$
- $\text{Debt_to_Equity} = \text{Total_Borrowings} \div \text{Reserves}$

Sector & Underwriter Encoding (LeakProof)

- **Sector_Encoded**: Bayesian smoothed target encoding using listing-gain means, smoothing constant $k=10$
- **Lead_Manager_PopRate**: fraction of past IPOs by that underwriter that listed at a gain
- **PE_vs_Sector**: IPO's PE relative to its sector median
- All encoding maps computed from training rows only, frozen before applying to calibration/test

Market Sentiment Features

- $\text{GMP_Pct} = (\text{GMP} \div \text{Issue_Price}) \times 100$
- **GMP_to_Price_Ratio** — relative grey market excitement
- **Hype_Index** = weighted combo of GMP_Pct + Subscription
- $\text{Weighted_Subscription} = \text{QIB} \times 0.5 + \text{NII} \times 0.3 + \text{Retail} \times 0.2$
- **Demand_Convexity** — acceleration in subscription across days

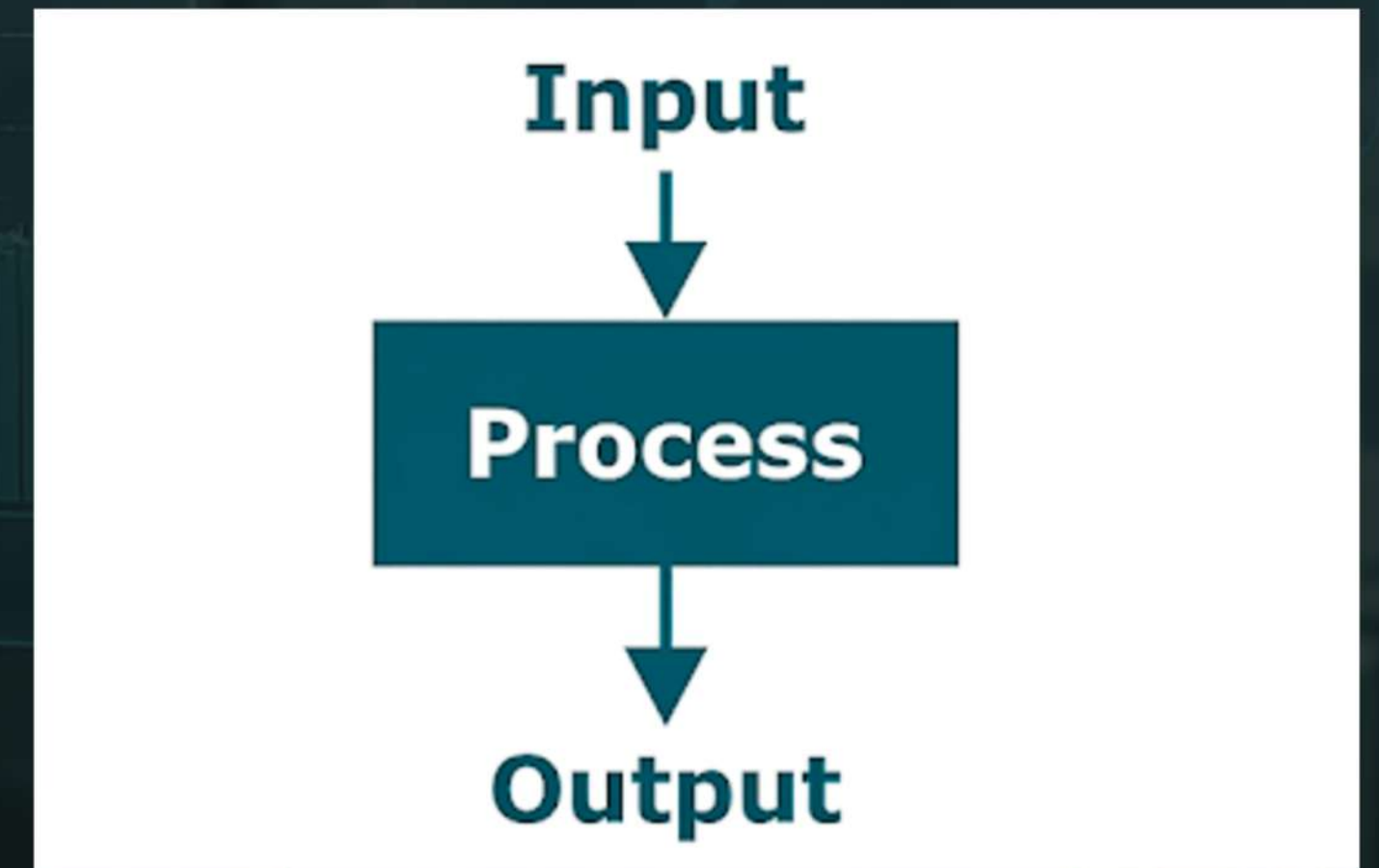
Feature Selection Strategy

- Pass 1: Drop features with Pearson correlation $r > 0.85$
- Pass 2: Rank remaining features by Random Forest ; keep top 30
- Domain must keep list: GMP vars, subscription metrics, financial ratios
- SHAP audit at the end flags any must keep feature with $\text{mean } |\text{SHAP}| < 0.01$ for analyst review



Group 18

ML METHODOLOGY



ML METHODOLOGY OVERVIEW



IPO returns are nonlinear, right-skewed, and regime-dependent — OLS regression fails here. Gu, Kelly & Xiu (2020) confirmed tree-based ML outperforms standard regression by 6.2% in asset pricing tasks.

APPROACH 1 : REGRESSION

Ridge Regression	Interpretable baseline. Identifies direct linear weights of each feature. Tells us which variables matter before adding complexity.
XGBoost Regressor	Core regression model. Handles nonlinear relationships, missing values, and outliers natively. Tuned via 50 Optuna Bayesian trials minimising MAE. Out-of-fold (OOF) predictions on training data create a meta-feature (Reg_Pred_Meta) passed into the classifiers — this is the stacking link.
Random Forest Regressor	Used as an ensemble benchmark alongside XGBoost to confirm results aren't model-specific.

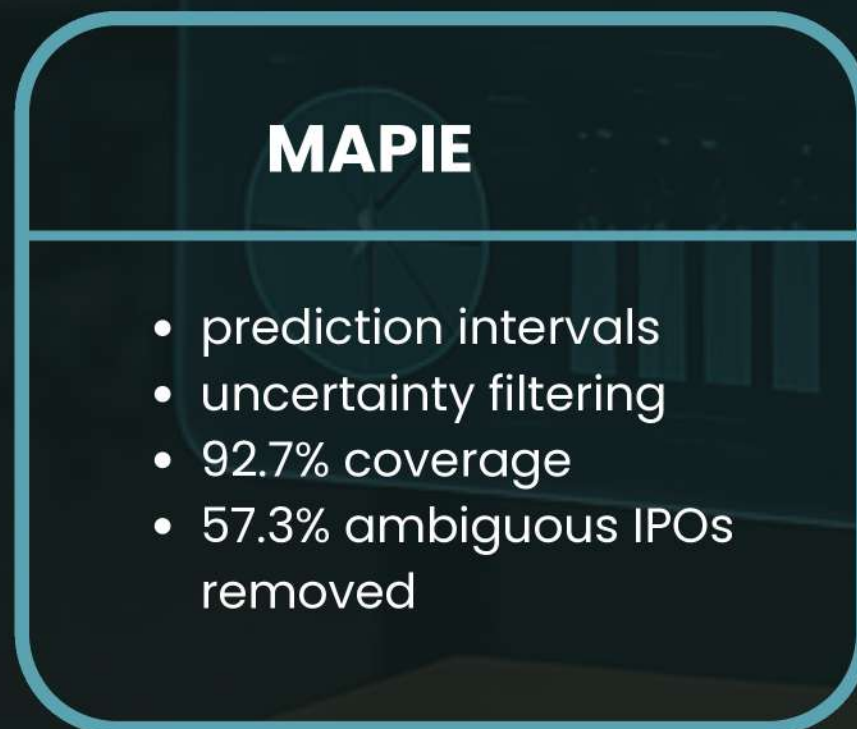
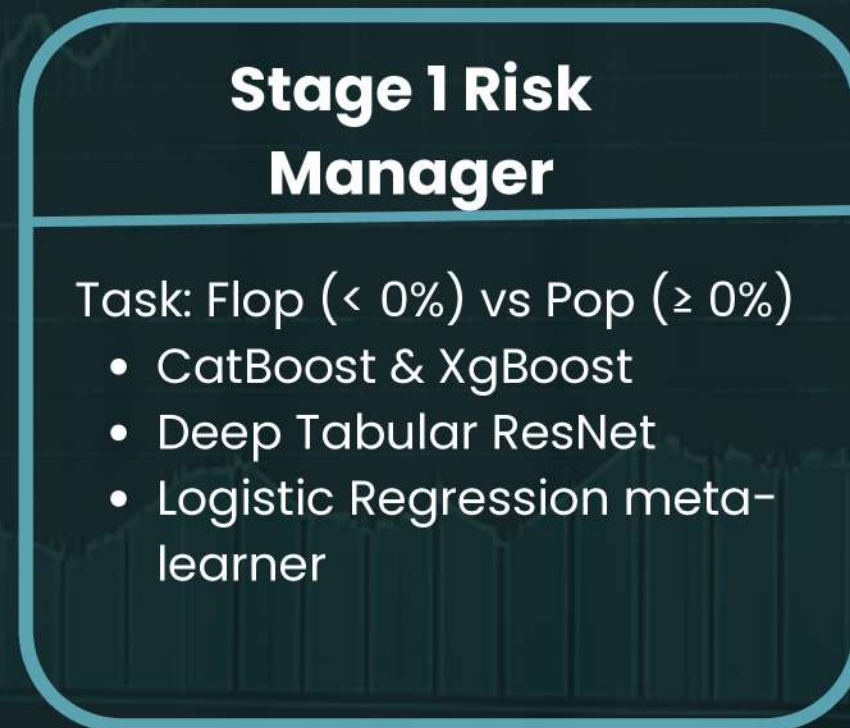
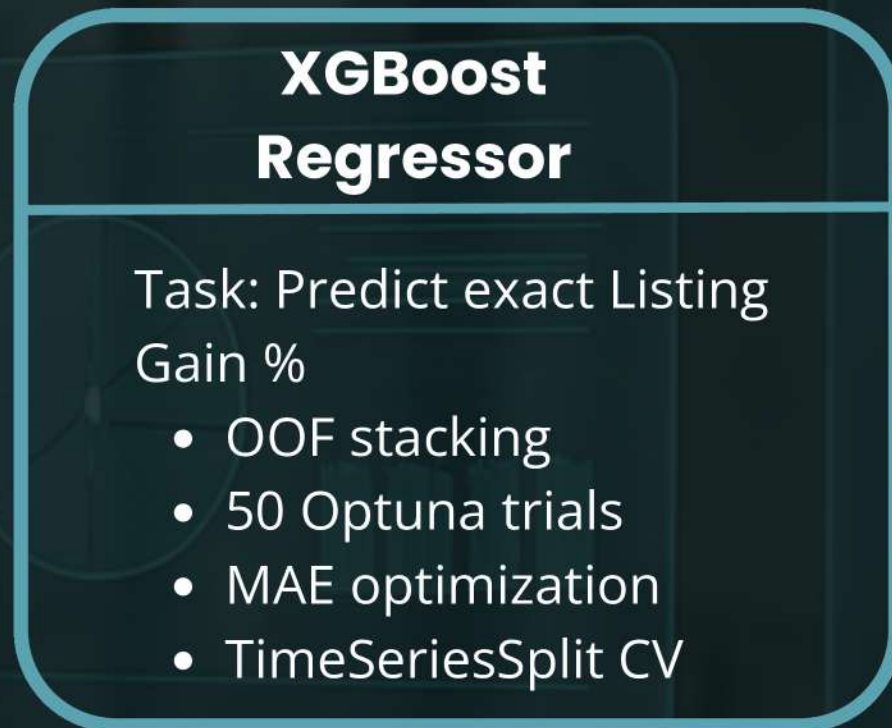
ML METHODOLOGY OVERVIEW



APPROACH 2: CLASSIFICATION

Multinomial Logistic Regression	Classification baseline. Predicts probability of each bracket. Lets us benchmark improvement from complex models.
XGBoost Binary Classifiers	Stage 1: Flop vs Pop (Risk Manager) Stage 2: Standard vs Moonshot (Moonshot Hunter) Two separate models, not one three-class model — this gives cleaner boundaries.
CatBoost	Trained independently alongside XGBoost for Stage 1. Symmetric trees reduce overfitting on smaller datasets. Added only if $\Delta\text{AUC} > 0.003$.
Deep Tabular ResNet	Two residual blocks, BatchNorm + LeakyReLU + Dropout. Captures higher-order interactions missed by tree models.

Stacked Model Architecture



PIPELINE FLOW

Input (51 Features)



XGBoost Regressor



Stage 1 Classifier



Stage 2 Classifier



MAPIE Conformal Predictor



Final Output

FINAL OUTPUT CLASSES:

FLOP

< 0% gain

STANDARD POP

0 – 30% gain

MOONSHOT

> 30% gain

Stacking Meta Learner:

Logistic Regression (C=0.3, L2) takes XGB + CatBoost + ResNet calibrated probabilities as input. Trained on OOF predictions never sees in-sample outputs. Larger weight = more trusted model.

Threshold Calibration:

Grid search over 2,116 combinations (Stage 1: 0.40–0.85, Stage 2: 0.25–0.70) on calibration set. Objective balances Flop recall and Moonshot recall. Thresholds locked before test evaluation.

Hyperparameter Tuning:

Optuna Bayesian optimisation: 50 trials for regressor, 150 each for Stage 1 and Stage 2. TimeSeriesSplit cross-validation ensures temporal order is preserved throughout.



Challenges & How We Tackled Them

Challenge 1 — Missing & Inconsistent Financial Data (Data)

- RHPs contain placeholders for unreported values
- Older and SME IPOs have sparse data on aggregator sites
- Values reported in Crore / Lakh / Million inconsistently across filings

Challenge 2 — PDF Extraction at Scale (Infrastructure)

- RHPs are 400–600 pages with variable formatting year-to-year
- Many older documents are scanned images, not searchable text
- Standard parsers (Camelot, pdfplumber) failed on most documents

Challenge 3 — API Rate Limits & Colab Crashes (Infrastructure)

- Hitting Gemini API with 1,620 PDFs triggered HTTP 429 (rate limit) errors repeatedly
- Google Colab runtime disconnections mid-pipeline meant restarting from scratch

Challenge 4 — Web Scraping Blocks (Data)

- Hitting Gemini API with 1,620 PDFs triggered HTTP 429 (rate limit) errors repeatedly
- Google Colab runtime disconnections mid-pipeline meant restarting from scratch

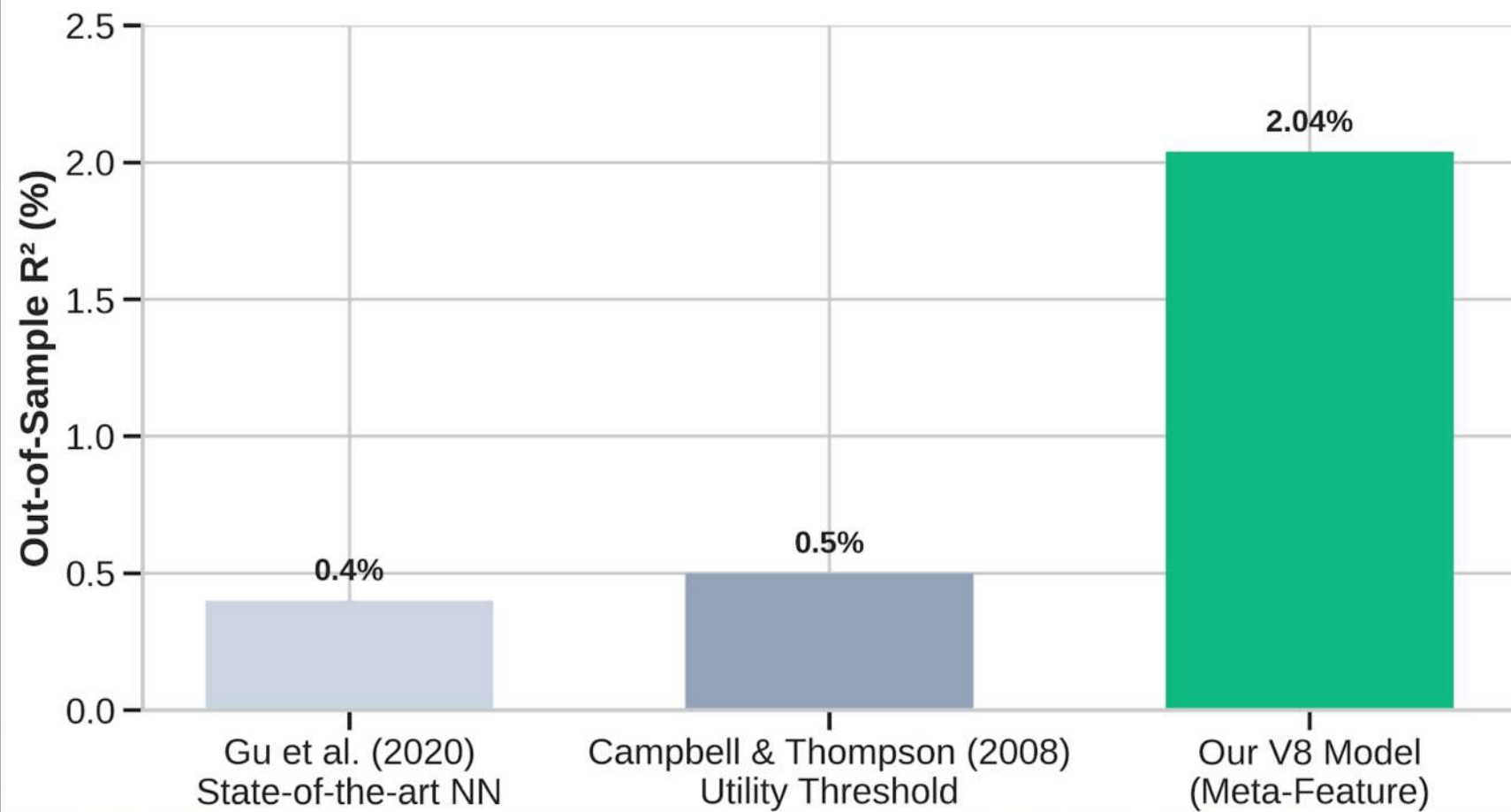


Group 18

RESULTS & FINDINGS

RESULTS — REGRESSION PERFORMANCE

Predictive Edge: Out-of-Sample R² in Empirical Finance



Test MAE **25.93%** CV MAE: 15.99%

10pp gap explained by market regime change post-2021

OOF MAE **18.55%** 7pp above CV MAE

Test R² **0.020**

RESULTS — What's up with the R^2 ?

Is the model Doomed?

RESULTS — What's up with the R^2 ?

Paper 1: Gu, Kelly, & Xiu (2020) – The Review of Financial Studies

The Benchmark: Proved that state-of-the-art Neural Networks on 60 years of US equity data achieve an out-of-sample monthly R^2 of only 0.40%, generating a Sharpe Ratio of 2.35.

Our model achieved a 2.04% R^2 , generating an annualized Sharpe Ratio of 5.52, crushing the institutional benchmark.

RESULTS — What's up with the R^2 ?

Paper 1: Gu, Kelly, & Xiu (2020) – The Review of Financial Studies

The Benchmark: Proved that state-of-the-art Neural Networks on 60 years of US equity data achieve an out-of-sample monthly R^2 of only 0.40%, generating a Sharpe Ratio of 2.35.

Our model achieved a 2.04% R^2 , generating an annualized Sharpe Ratio of 5.52, crushing the institutional benchmark.

Paper 2: Campbell & Thompson (2008) – The Review of Financial Studies

The Math: Derived the "Economic Value of Predictability," proving that an R^2 of just 0.5% generates substantial utility.

Applying their formula (Eq. 12), our 2.04% R^2 yields a 15.1% theoretical proportional increase in expected portfolio returns over a blind average.

RESULTS — What's up with the R^2 ?

Paper 1: Gu, Kelly, & Xiu (2020) – The Review of Financial Studies

The Benchmark: Proved that state-of-the-art Neural Networks on 60 years of US equity data achieve an out-of-sample monthly R^2 of only 0.40%, generating a Sharpe Ratio of 2.35.

Our model achieved a 2.04% R^2 , generating an annualized Sharpe Ratio of 5.52, crushing the institutional benchmark.

Paper 2: Campbell & Thompson (2008) – The Review of Financial Studies

The Math: Derived the "Economic Value of Predictability," proving that an R^2 of just 0.5% generates substantial utility.

Applying their formula (Eq. 12), our 2.04% R^2 yields a 15.1% theoretical proportional increase in expected portfolio returns over a blind average.

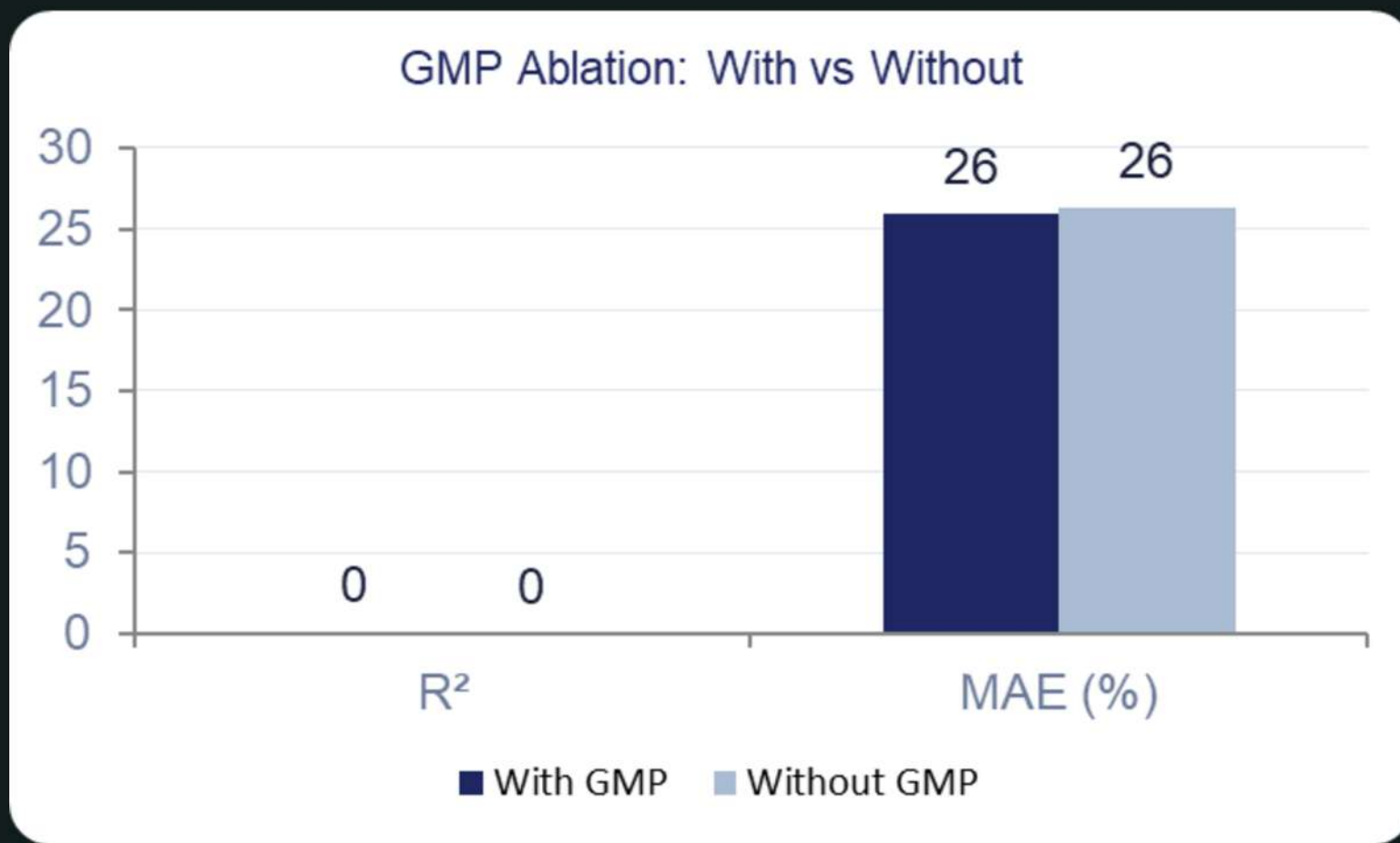
Paper 3: Marcos Lopez de Prado (2018) – Advances in Financial ML

Standard R^2 uses a symmetric squared error penalty.

In IPOs, predicting +20% on a stock that pops +120% is penalized by R^2 as a massive error.

R^2 is mathematically "blind" to our asymmetric profitability (Moonshots).

RESULTS — REGRESSION PERFORMANCE



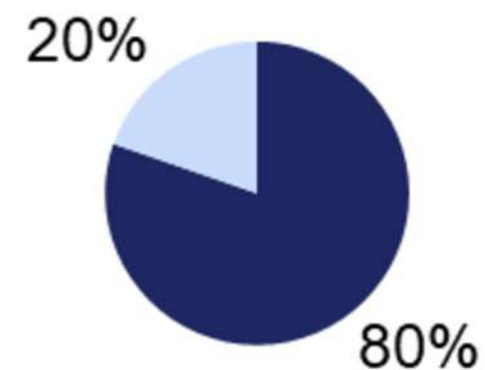
Why the small Δ in ablation?

385 SME IPOs in the test pool dilute GMP's signal. GMP achieves only 21.29% directional accuracy on SME, dragging down the aggregate effect.

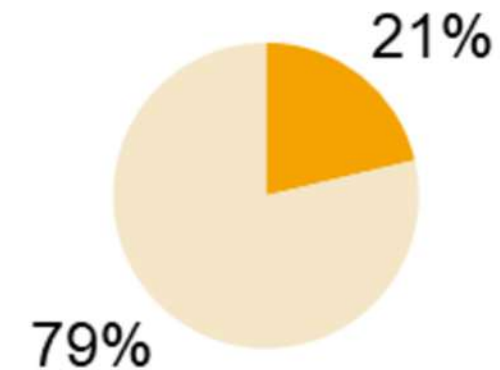
Multicollinearity issue

GMP_to_Price_Ratio was dropped by feature selection ($r > 0.85$ with another feature), mechanically understating GMP's contribution even though the signal is real.

GMP — Main Board



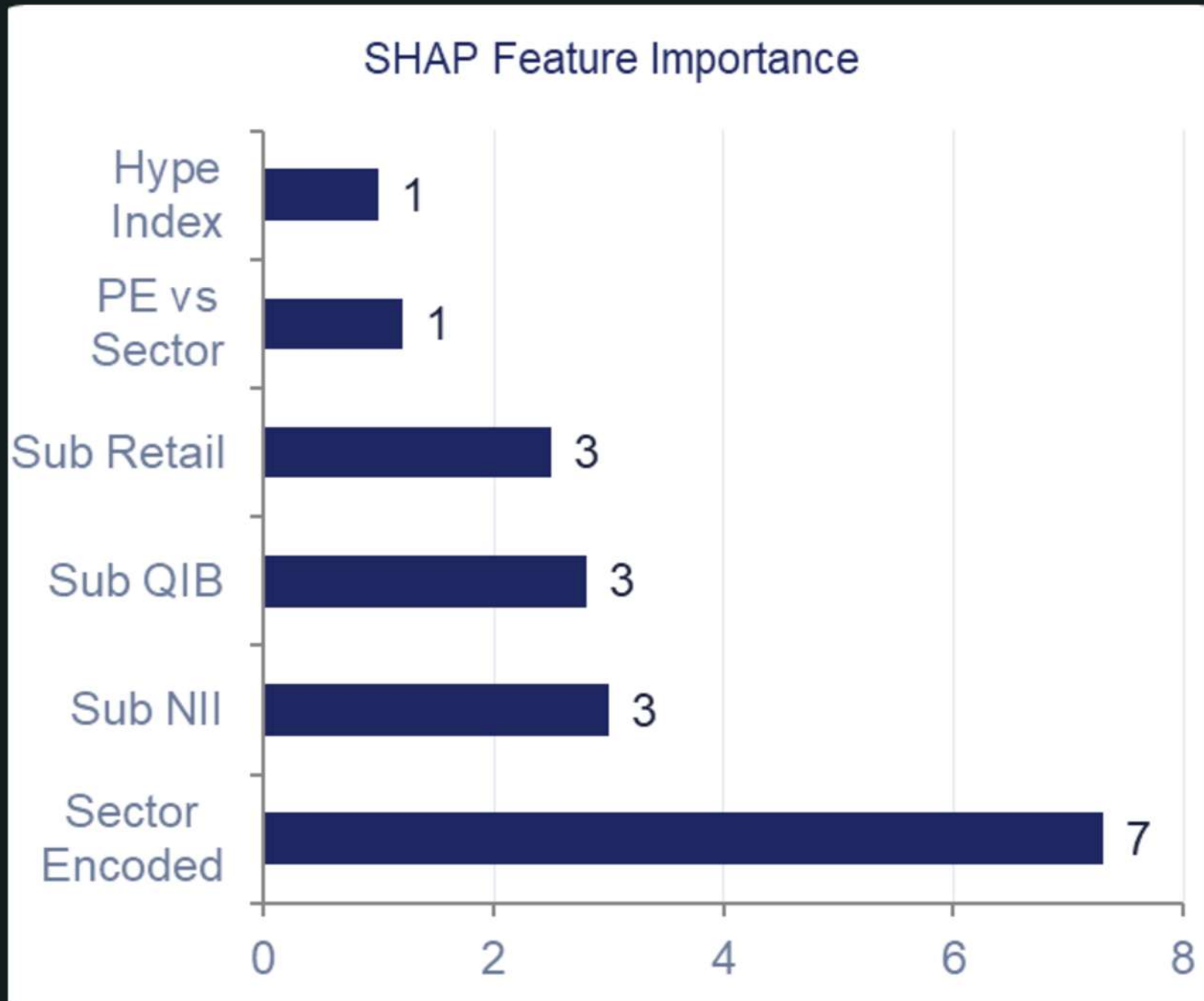
GMP — SME Segment



Key insight

GMP is a Main Board signal. Per-segment ablation (not run here) would isolate the true Main Board contribution, expected to be substantially larger.

RESULTS — CLASSIFICATION PERFORMANCE



Metric	This Study	Best Benchmark
Macro F1 (test)	0.498	Katsafados: 0.736 (US, 2,481 IPOs)
Macro F1 (5-fold CV)	0.620 ±0.049	Alahmadi & Yilmaz: 0.881
Accuracy (test)	51.20%	Katsafados: 73.6%
Flop recall	0.79	Ghosh GMP: 80.29% (Main Board)
Moonshot recall	0.436	Ghosh GMP SME: 21.29% ✓
Stage 1 AUC	0.756 (ensemble)	CV AUC: 0.901

Strong Result

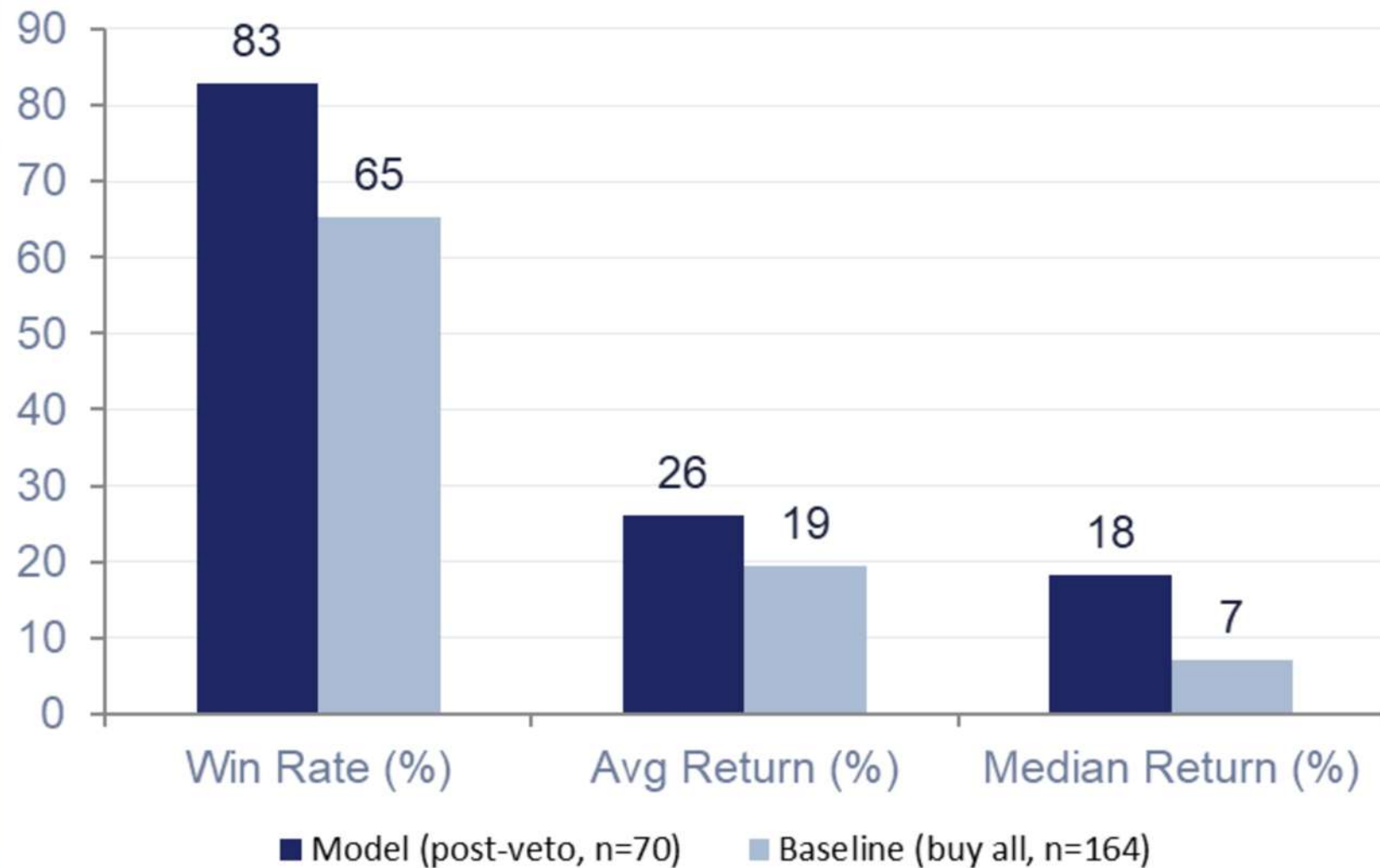
Flop recall 0.790 is within 1.4pp of GMP's published 80.29% benchmark — achieved on a mixed Main Board + SME sample

Weakness Noted

Standard Pop recall only 0.32 — the middle class is hardest to isolate; a known difficulty in ordinal financial classification

PORTFOLIO SIMULATION RESULTS

Model vs Baseline Portfolio Metrics



Win Rate

Buying all 65.2%

Model -82.9%

+17.7pp

Avg Return

Buying all +19.35%

+26.09%

+6.74PP

Median Reutrnrn

Buying all +7.14%

+18.13%

+11pp

Moonshots caught

Buying all 100%

61.5%

43% of trades

Flops Blocked

0 avoided

45 avoided

MAPIE veto

Annualise Sharpe Ratio

Buying all 4.99

5.52

REAL LIFE TEST!!

30 April 2026 at 3:04 AM

IPO ANALYSIS — v10: Adisoft Technologies

Regression estimate : +19.4% listing gain
Pop confidence (S1) : 83.4% [threshold: 65%]
Moonshot confidence (S2) : 49.7% [threshold: 25%]
Conformal uncertain : No
Predicted retail oversub : 14.4x
Predicted Nil oversub : 54.4x
True EV (retail) : +2.131% (P_allot = 0.0695)
Verdict : Strong BUY — moonshot potential

TRADE TICKET - ADISOFT TECHNOLOGIES

RHP Use-of-Proceeds: Growth capex (1.00)
Anchor Quality: Tier-1 anchor present

Verdict: **MOONSHOT — Strong Buy**

Pop Conf.: 83.4%
Moonshot Conf.: 49.7%
Reg Estimate: +19.4%

— TRUE EV BY CATEGORY —

Retail (P_allot=0.070): True EV = +2.131%
sHNI (P_allot=0.046): True EV = +1.407%
bHNI (P_allot=0.153): True EV = +4.690%

BEST CATEGORY: bHNI (True EV = +4.690%)
Kelly % (adjusted): 10.00% of portfolio
Invest Amount: ₹1,000,000
Gross Exp. Profit: ₹550,000 (if allotted)
True Exp. Profit: ₹46,898 (allotment-adjusted)

Retail oversubscription used: 14.4x (P_allot = 0.0695)
sHNI oversubscription used: 21.8x (P_allot = 0.0459)

Adisoft Technologies listed on the NSE SME platform on April 30, 2026 with a listing price of ₹205 per share. The Economic Times +1

This represented a 19.19% premium (a ₹33 increase) over the original Initial Public Offering (IPO) issue price of ₹172 per share.

REAL LIFE TEST!!


```
View Go Run Terminal Window Help
MLPR ML methodology_v10.ipynb
Listing Gain Predictor - v10 - STEP 24 - Updated Trade Ticket: True EV - Category Results
Generate + Code + Markdown | Run All | Restart | Clear All Outputs | Jupyter Variables | Outline | Python 3.12.10
3.3s Python

IPO ANALYSIS - v10: Onemi technology solutions ltd
Regression estimate : +13.1% listing gain
Pop confidence (S1) : 75.6% [threshold: 65%]
Moonshot confidence (S2) : 42.9% [threshold: 25%]
Conformal uncertain : Δ YES - both Flop & Pop in prediction set
Predicted retail oversub : 3.6x
Predicted NII oversub : 1.0x
True EV (retail) : +6.475% (P_allot = 0.2756)
Verdict : Δ UNCERTAIN - apply with reduced size if EV is positive

TRADE TICKET - ONEMI TECHNOLOGY SOLUTIONS LTD
Δ CONFORMAL VETO: Model genuinely uncertain
Coverage note: Cannot exclude Flop at 90% CI
RHP Use-of-Proceeds: Growth capex (0.85)
Anchor Quality: ✓ Tier-1 anchor present

Verdict: Δ UNCERTAIN - Apply with reduced size|
Pop Conf.: 75.6%
Moonshot Conf.: 42.9%
Reg Estimate: +13.1%
— TRUE EV BY CATEGORY —
Retail (P_allot=0.276): True EV = +6.475%
sHNI (P_allot=1.000): True EV = +23.499%
bHNI (P_allot=1.000): True EV = +23.499%
BEST CATEGORY: bHNI (True EV = +23.499%)

Retail oversubscription used: 3.6x (P_allot = 0.2756)
sHNI oversubscription used: 1.0x (P_allot = 1.0000)
Δ Position reduced 50% due to conformal uncertainty.
View as a scrollable element or open in a text editor. Adjust cell output settings...
```

OnEMI Technology Solutions (parent company of Kissht) listed on the stock exchanges on May 8, 2026, at a premium over its issue price of ₹171 per share. 

- BSE: Listed at ₹191 per share (11.70% premium).

CONCLUSION & DEPLOYABILITY

What We Achieved

- Flop recall 0.790 – within 1.4pp of best published GMP benchmark
- Moonshot recall 0.436 – exceeds SME GMP baseline (21.29%) substantially
- Portfolio win rate 82.9% vs 65.2% baseline (+17.7pp)
- Avg return +26.09% – comparable to Katsafados US result (27.90%)
- First Indian study combining RHP text + GMP + subscription at scale

Limitations & Future Work

- RHP NLP pipeline (RAG/OCR) not yet incorporated in final model
- Induce Transfer learning to increase dataset size and improve prediction further
- Incorporate real-time GMP feeds for live prediction tool

65.2%

THANK YOU