

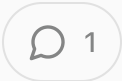
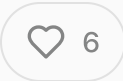
Rethinking the ML Pipeline: Why "Train Wide, Filter Smart" is a Game-Changer for AI Factor

What if this pre-filtering is limiting our AI Factor models' potential?



ANDREAS HIMMELREICH

OCT 30, 2025





Systematic AI Investing Portfolios is a reader-supported publication. To receive new posts and support my work, consider becoming a free or paid subscriber.

[Subscribe](#)

Train Wide, Filter Smart

For years, the standard approach in quantitative investing has been to pre-filter a universe (e.g., “large caps,” “profitable firms”) before training a machine learning model. The goal: reduce noise and focus the signal.

But what if this pre-filtering is limiting our models’ potential?

The powerful alternative: Train Wide, then Filter Smart!

The Architecture Shift:

1. Train Wide: Let your (nonlinear!) model (e.g., LightGBM, ExtraTrees) learn from a broad and noisy universe. Expose it to the full market ecosystem—the quality stocks, the speculative ones, everything.
2. Filter Smart: Apply simple, robust quality filters (like $\text{CurFYEPSMean} > 0$) only at the final portfolio construction phase.

Why This Delivers Superior Results:

- Deeper Market Intelligence: A model trained on a wide universe understands regime changes, factor interactions, and market dynamics that are invisible in a pre-filtered sandbox. **It learns the context of what makes a signal valuable.**
- Transforms the Signal Distribution: This method doesn’t just filter out “bad stocks.” It fundamentally reshapes the output distribution of AI-driven factors, especially under Z-Score normalization. The result? A dramatic reduction in performance “spikiness” and significantly smoother equity curves.
- Enables Concentrated, Low-Volatility Portfolios: By cleaning the signal after the prediction, you can run highly concentrated portfolios based on ML rankings without inheriting the volatility of the broader, noisier universe. My LightGBM strategies n

exhibit the stability I previously only associated with ensemble methods like ExtraTrees.

This “Train Wide, Filter Smart” paradigm separates the model’s job — pattern recognition — from the portfolio manager’s job—risk and quality control. It leverages the full power of ML while ensuring the final output is institutional-grade, defensible and robust.

Buy Rules: The Devil’s in the Detail

Building on the “Train Wide, Filter Smart” framework, a crucial technical insight emerges: Not all buy filters are created equal. The effectiveness of a filter is deeply dependent on the normalization method of your underlying AI factors.

Through rigorous backtesting, a clear pattern has emerged:

For “Rank & Date” Normalized Models: Complex percentile filters like $\text{FRank}(\text{EPS_Revision}) > 80$ work well. The rank-based system is inherently robust to outliers, making it a suitable partner for other relative, cross-sectional ranking rules.

For “Z-Score & Date” Normalized Models: Simple, absolute quality filters like $\text{CurFYEPSMean} > 0$ are dramatically superior. The Z-Score method is highly sensitive to distribution tails. Using a FRank filter here often injects the very spiky, extreme values that Z-Scores amplify, leading to volatile performance.

The Underlying Principle:

Your portfolio construction layer shouldn’t fight your feature engineering layer. A FRank filter on a Z-Score model often tries to clean a noisy signal with another noisy relative signal. In contrast, a simple quality gate ($\text{EPS} > 0$) creates the stable, well-behaved distribution that Z-Score normalization requires to shine.

This explains a key nuance: why a Zscore & Date system with a FRank buy rule can work well on the S&P 500 but fails on small caps? The S&P 500 universe is inherent quality filter; the data is much less noisy to begin with, so the complex rule doesn't fight the normalization.

Why are AI Factor Models robust?

Building a robust AI-driven strategy hinges on a deep understanding of its internal architecture. The true strength isn't in a single "magic" formula, but in a decentralized, multi-layered system designed to withstand market shifts.

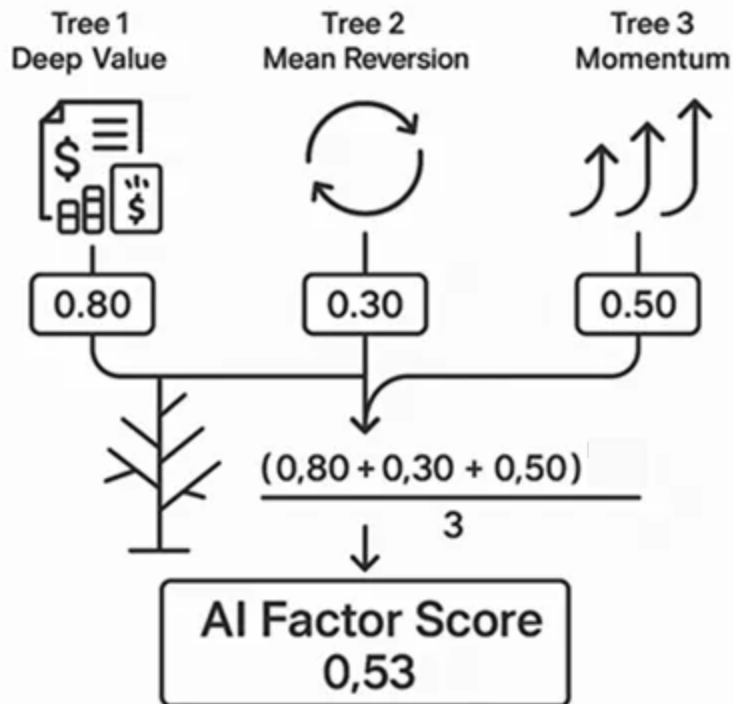
Let's break down the engine: an ensemble of decision trees. With a typical setup of features, a depth of 5, and 500 trees, we are not building one model. We are building a vast network of micro-rules.

A single tree of depth 5 can have up to 32 leaf nodes — each the end of a unique decision path. The total system capacity is profound:

$500 \text{ trees} \times 32 \text{ paths/tree} = 16,000 \text{ potential decision paths.}$

Factoring in the combinations across 179 features, the number of effective, unique interacting conditions easily surpasses 25,000. These are not complex, hand-crafted rules, but simple, stochastic micro-logic statements.

How Tree-Based Models Rank Stocks



The core robustness comes from the ensemble method (s above!). Each of the 500 trees is an independent expert, casting a single “vote.” The final prediction is a averaged consensus of this entire committee. For a catastrophic failure, a regime shift must invalidate a majority of these 25,000+ paths across hundreds of independent trees simultaneously — a statistical improbability.

AI Factor stays robust, even with buy rules on the portfolio strategy level!

This leads to a critical question: If the model is so complex, why doesn't a strict b rule like “CurFYEPSMean > 0” destroy its subtle intelligence?

The answer is that the buy rule filters the portfolio, not the model's knowledge. The model's 25,000+ rules represent a pre-trained understanding of the entire market landscape. Applying a quality gate doesn't delete this knowledge; it focuses its application on a cleaner, more stable segment.

From the 25,000+ available rules, the model activates a large subset — numbering in the thousands — to rank every stock that passes the filter. The diversity and number of rules that remain in play are more than sufficient to generate nuanced, intelligent predictions.

We are not reducing a sophisticated brain to a few simple ideas. We are giving it a cleaner dataset to process. The buy rule ensures input stability; the ensemble of trees ensures the output alpha is robust and intelligent.

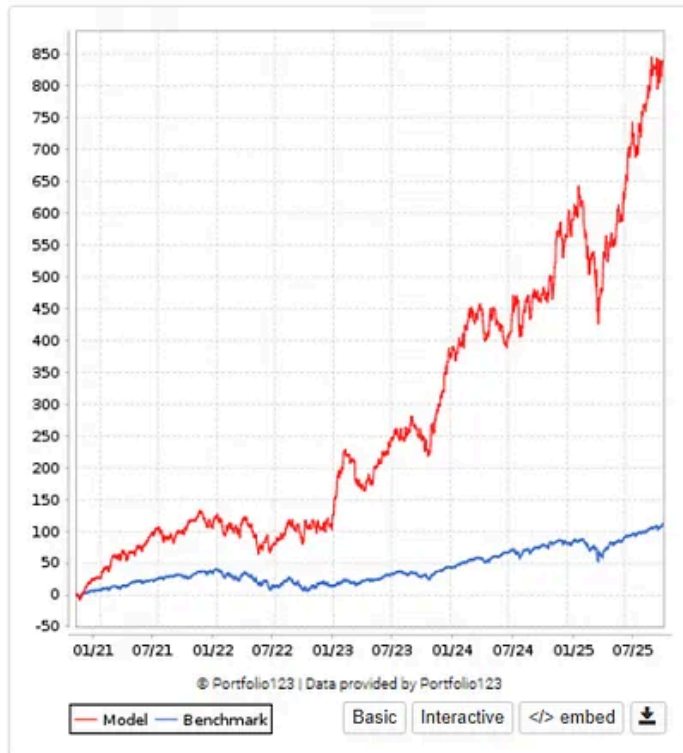
I have built > 200 Portfolio Strategies with extensive Buy rules, none of them failed OOS Live so far (OOS spanning from weeks to a year).

Example 1

Portfolio Strategy based on an AI Factor System with a Universe of $\text{AvgDailyTot}(20) < 2000000000$ (so leaning to the mid and small caps).

Not bad:

On 02/28/2022, the model's market exposure deviated 5.04% above the target, an excess of 0.04%.



General Info

[PDF Report](#)

Total Market Value (inc. Cash)	931,61
Cash	-12,61
Number of Positions	
Last Trades (5)	10/2
Period	10/17/20 - 10/2
Sizing Method	Dynamic We
Next Reconstitution (Every Week)	11/03/25 In 5
Next Rebalance (Every Week)	11/03/25 In 5
Mode	Auton
PIT Method - Prelim	
Benchmark	S&P 500 (SPY:U
Universe	No OTC Exchange + min 200 mi Fin:
Ranking System	200

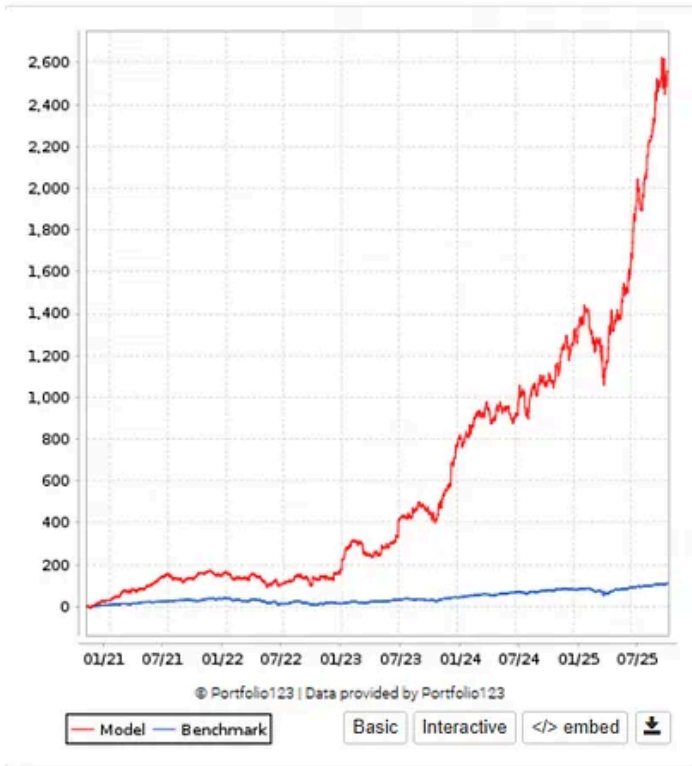
Quick Stats as of 10/28/2025

Total Return	831.1
Benchmark Return	112.1
Active Return	719.1
Annualized Return	55.1
Annual Turnover	278.1
Max Drawdown	-29.1
Benchmark Max Drawdown	-24.1
Overall Winners	(398/712) 55.1
Sharpe Ratio	
Correlation with S&P 500 (SPY:USA)	

Now the same strategy with a liquidity filter below 5 Million
 $\text{AvgDailyTot}(20) < 5000000$

Used in: 2 (Live Book) / 16 (Simulated Book)

On 05/02/2022, the model's market exposure deviated 6.20% above the target, an excess of 1.20%.



General Info

[PDF](#) [Report](#)

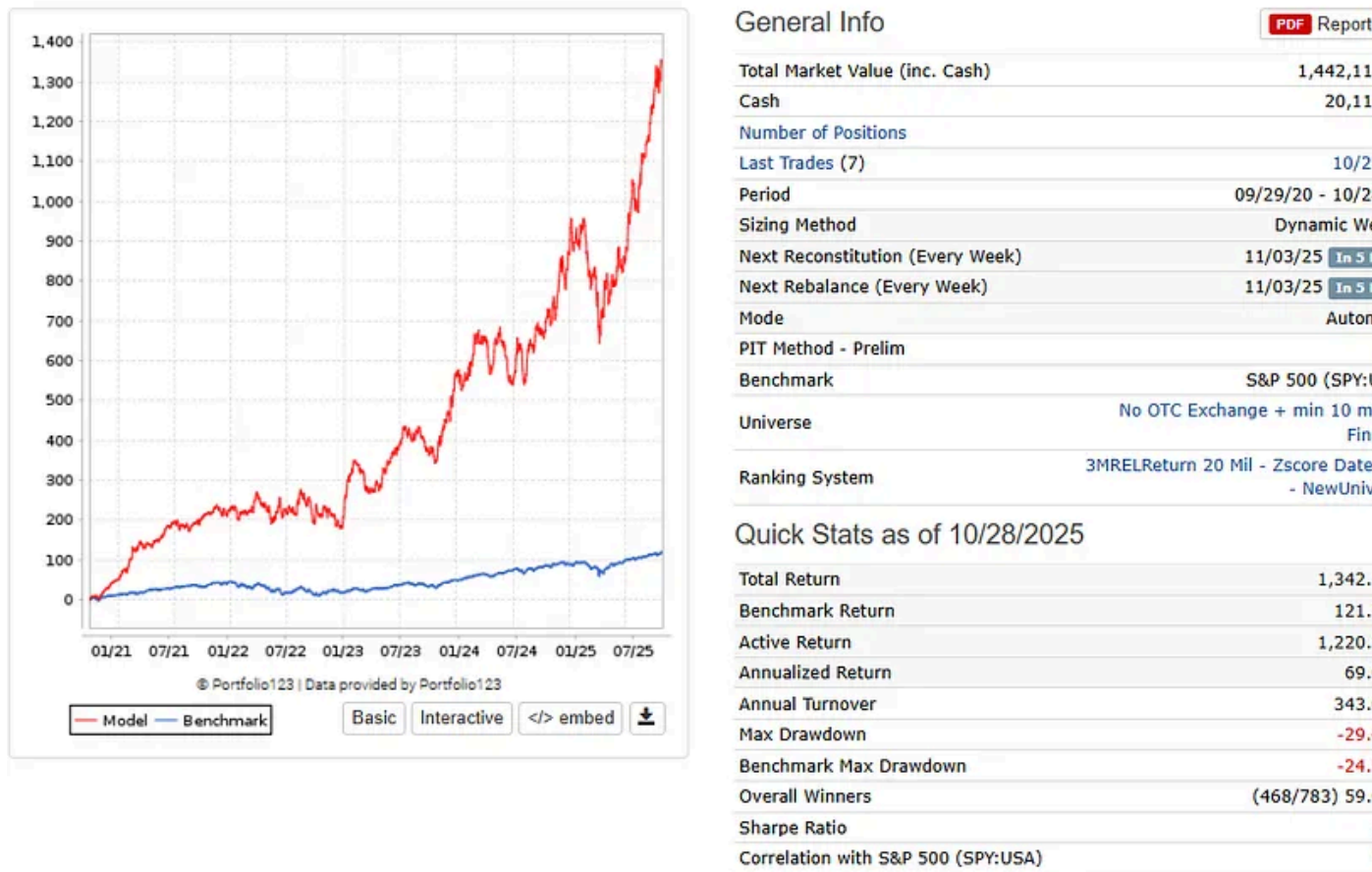
Total Market Value (inc. Cash)	2,653,687
Cash	-70,500
Number of Positions	
Last Trades (10)	10/27
Period	10/17/20 - 10/28/25
Sizing Method	Dynamic We
Next Reconstitution (Every Week)	11/03/25 In 5 D
Next Rebalance (Every Week)	11/03/25 In 5 D
Mode	Autom
PIT Method - Prelim	
Benchmark	S&P 500 (SPY:U
Universe	No OTC Exchange + min 200 mi Fina
Ranking System	200

Quick Stats as of 10/28/2025

Total Return	2,553.6
Benchmark Return	112.2
Active Return	2,441.4
Annualized Return	91.7
Annual Turnover	402.7
Max Drawdown	-28.9
Benchmark Max Drawdown	-24.5
Overall Winners	(339/619) 54.0
Sharpe Ratio	1
Correlation with S&P 500 (SPY:USA)	0

Example 2

Nothing against it, nice total return strategy:



And here comes the kicker, same strategy just with the following rules follows —>

Buy Rules (Implicit AND)

copy to screen

Buy1

MedianDailyTot(20) > 50000

Buy3

Rank > 97

Buy3

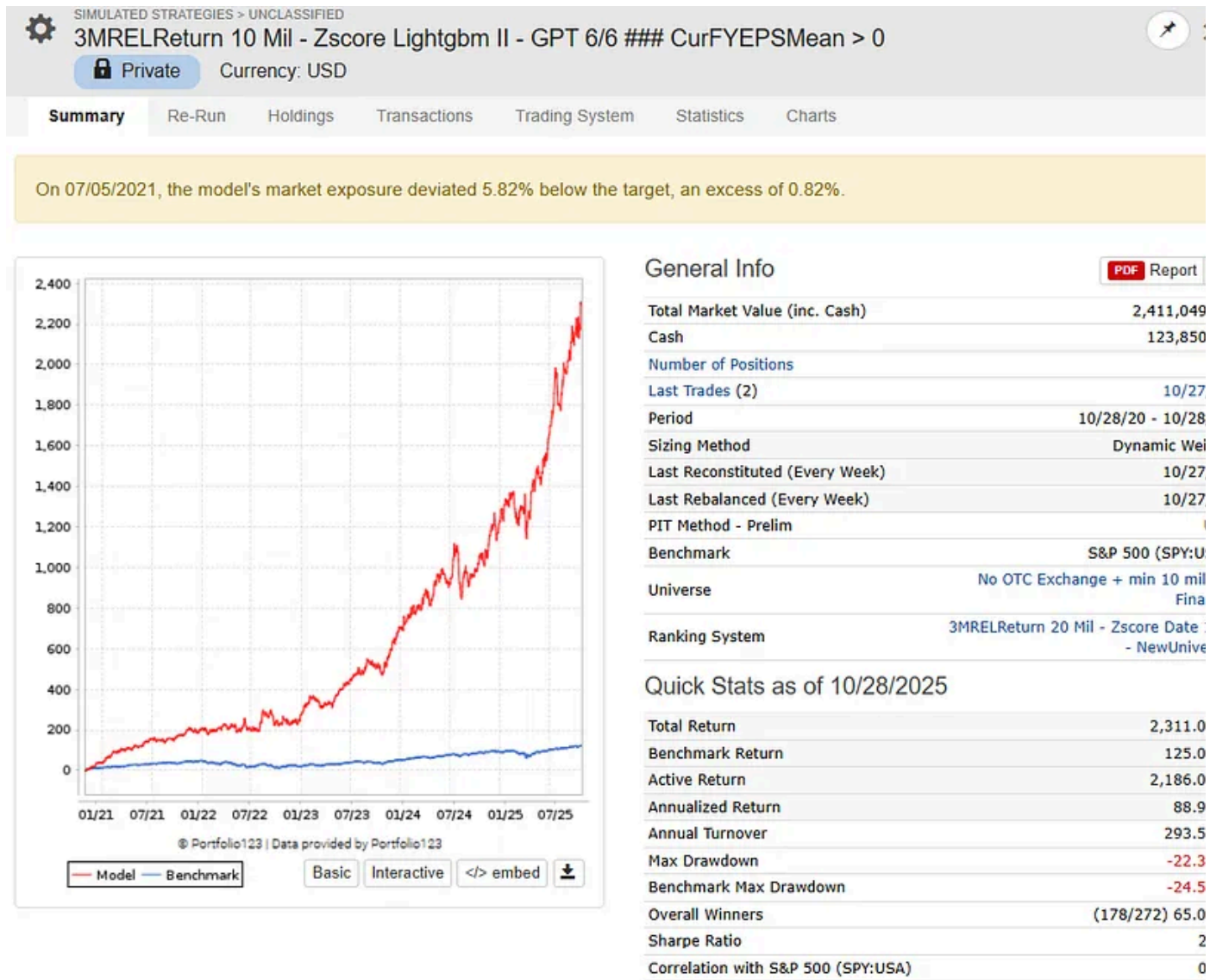
CurFYEPSMean > 0 or AltmanZOrig > 0

Sell Rules (Implicit OR)

copy to screen

Rank

Rank < 80



While I wouldn't run this as a standalone portfolio, it makes for a powerful satellite allocation designed to boost returns in a larger, core strategy!

Some examples are here: https://x.com/GfI_Himmelreich/status/19813538452602311:

The Junk Problem? Yes and No!

It's a common surprise: a sophisticated ML model, trained on hundreds of factors, v often load up on speculative stocks. We expect it to learn wisdom, but it only learns patterns.

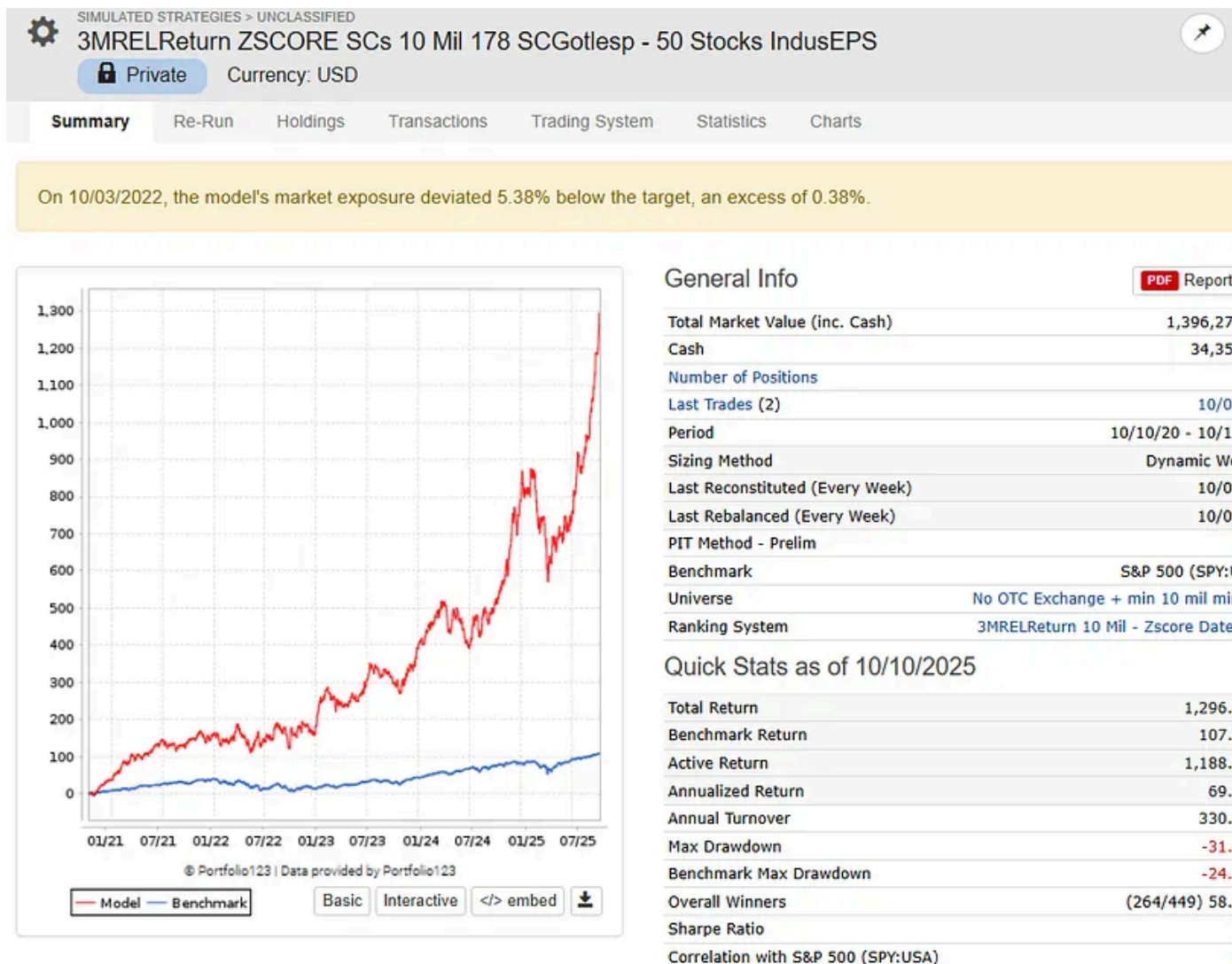
The reason is fundamental. ML optimizes for statistical prediction. Its sole objective is to minimize error, with no innate concept of risk or drawdowns (if the target is price based!).

This leads to the “Junk Factor” problem. And here’s the crucial part: The model is often statistically correct. Junk stocks — with their extreme volatility and binary outcomes — can be highly predictable and immensely profitable: but only during (later) bull markets.

The ML brilliantly identifies these explosive, short-term patterns. It doesn’t understand that they are regime-dependent and come with tail risk.

This is why a pure ML portfolio can be spiky and volatile. It’s chasing patterns that work until they don’t — often spectacularly.

Stuff like this (in 2008 my best guess is, it would have a 70% DD and be up to new highs in 12 Months):



But: I am not routing against those systems, they can be great in a system book!

Why?

Strategies like this, despite their high volatility, can provide valuable uncorrelated alpha and can be extremely effective as a long leg in a market-neutral book or as a satellite allocation to boost overall returns.

Uncorrelated Alpha Source: Its 0.61 correlation to the SPY, while positive, suggests it's not just a leveraged beta bet. It's capturing a different type of risk premia (likely high-octane, small-cap momentum), which is a diversifying return stream.

1. **Ideal for a Long/Short Book:** This is the key. You could pair this high-volatility long portfolio with a different, more stable short portfolio (e.g., shorting low-quality value traps or stable low-momentum stocks). The combined book would aim to be market-neutral, harvesting the pure “junk momentum” alpha while hedging out the brutal market-directional drawdowns.
2. **Satellite Allocation:** In a core-satellite framework, this strategy could be a small high-conviction “satellite” intended to boost the total return of a much larger, stable “core” portfolio. The core ensures survival; the satellite provides the kick.

Some of the best Traders trade Junk ;-)

By the way, it's no coincidence that the best discretionary traders I know — the ones who consistently print > 100% years trading breakouts — are masters at knowing exactly when to trade this so-called “junk.”

They don't avoid it; they exploit it with impeccable timing and market feel. Their success is the ultimate proof of concept, showing what is possible in these volatile corners of the market.

Train Wide, Filter Smart and weed out the junk (if you want!)

The ML's purpose is raw, unbiased pattern recognition across the entire market landscape — a task it performs with immense scale.

Our role is to apply the economic judgment and risk management it inherently lacks. By training the model wide and then applying a decisive quality filter at the portfolio level, we forge a powerful synergy.

Alternatively, we can consciously choose to harness those “junky” but historically profitable signals, fully aware of their regime-dependent nature. **The key is that it remains our deliberate, strategic decision, not the model’s blind directive.**

Best Regards

Andreas

Subscribe to Systematic AI Investing Portfolios

By Andreas Himmelreich · Launched a year ago

Systematic Investing Portfolios — Built with AI & Machine Learning General Posts are free! ♦ Include Live (AI) Strategy Portfolios (Small, Mid/Large, Large Cap) for paying subscribers

Subscribe

By subscribing, you agree Substack's [Terms of Use](#), and acknowledge its [Information Collection Notice](#) and [Privacy Policy](#).



6 Likes · 8 Restacks

← Previous

Next

Discussion about this post

Comments

Restacks



Write a comment...



Brent Harless  19 Dec

 Liked by [Andreas Himmelreich](#)

So good.

 LIKE (1)  REPLY

