

# The Prior-Informed Modeling Framework *versus* Data-Driven, Machine-Learning Approaches

## Motivation

In modern finance, analysts often employ machine-learning (ML) approaches to verify the accuracy and consistency of empirical findings using train–test splits, cross-validation, and large-scale model comparison. These practices are valuable for prediction and can reduce overfitting at the model-selection stage. However, purely data-driven pipelines typically do not formalize the inferential questions that are central to anomaly testing; they do not quantify the probability that an apparent effect is a false discovery after controlling for the benchmark, they do not encode the ex-ante plausibility or base rate of true effects, and they do not explicitly adjust for the size and dependence of the search space. The finance literature documents that extensive search across many, often dependent, strategies inflates spurious findings (data-snooping and multiple-testing risk), that in-sample success frequently fails out of sample, and that seemingly profitable rules can be overfit and non-implementable once realistic frictions are included (Bailey et al., 2014; De Prado, 2018; Hansen, 2005; Sullivan, Timmermann, and White, 1999; Welch and Goyal, 2008; White, 2000).

The prior weight on null hypothesis,  $\pi_{\text{null}}$ , is designed to address these gaps. It operationalizes prior skepticism that comes from economic theory, historical replication rates, and domain knowledge about implementability, turning that information into explicit prior odds against a true effect. The posterior odds then combine these prior odds with the evidence from the replication sample, yielding a posterior probability of a false discovery that is

interpretable, reproducible, and comparable across studies. In short, while ML verification is useful for prediction,  $\pi_{\text{null}}$ , provides the missing inferential layer: it penalizes wide and dependent searches in a transparent way, reflects realistic base rates of success, and disciplines signals that are likely fragile under trading frictions, thereby aligning anomaly assessment with the substantive concerns established in the finance literature.

## Data-snooping and Multiple Testing

When many candidate signals, specifications, and hyper-parameters are tried, nominal p-values overstate evidence. This is the data-snooping problem: repeated search increases the chance of finding seemingly significant results by luck. Bootstrap-based reality checks and related tests explicitly correct for searching over a dependent model set (Sullivan, Timmermann, and White, 1999; White, 2000), and the SPA test formalizes inference when many competing forecast rules are dependent (Hansen, 2005). Purely data-driven pipelines often rely on cross-validation to prevent overfitting, but cross-validation selects from the search space without quantifying the family-wise or false-discovery error that arises from the search itself. As the tested universe grows, so does the expected number of false positives unless skepticism is increased commensurately.

Our framework essentially encodes this skepticism as *prior odds* against a true positive. The posterior odds satisfy

$$\frac{P(H_{\text{null}} \mid \text{data})}{P(H_{\text{alt}} \mid \text{data})} = \underbrace{\frac{\pi_{\text{null}}}{1 - \pi_{\text{null}}}}_{\text{prior odds}} \times \underbrace{BF_{\text{alt-null}}(\text{data})}_{\text{evidence}}$$

so larger  $\pi_{\text{null}}$  raises the evidence bar required to report a discovery. This acts as a principled penalty for a wide search, converting pre-data, base-rate beliefs into an explicit threshold on the posterior null probability. In other words, more searching implies higher prior skepticism, which your decision rule translates into tighter discovery criteria.

## Dependence among Strategies

Anomaly signals and model outputs are often highly correlated because they share characteristics, construction rules, or latent risk exposures. Standard multiple-testing controls can be optimistic when tests are dependent, or very conservative when designed to be valid under any dependence, and many ML selection steps do not report any explicit error control for dependence. The finance literature documents substantial comovement and overlap among strategies, highlighting that dependence is a first-order feature rather than a nuisance (e.g., Sullivan, Timmermann, and White, 1999; White, 2000; Hansen, 2005).

A explicit prior-skepticism parameter absorbs part of this dependence risk by tilting posterior odds against marginal signals whose apparent strength may be inflated by correlation with already searched features. Because your posterior-null mapping is monotone in the signed statistic, dependence cannot manufacture confidence *ex nihilo*; it must overcome prior odds that reflect realistic base rates. In practice, setting a higher  $\pi_{\text{null}}$  for families known to be crowded or overlapping functions as a transparent and reproducible correction for dependence-driven optimism.

## Out-of-sample Fragility and Backtest Overfitting under Frictions

Data-driven selection often produces strategies that degrade out of sample, particularly once transaction costs, shorting constraints, and capacity are incorporated. Large-scale studies show predictive regressions that look promising in sample fail to beat simple benchmarks out of sample (Goyal and Welch, 2008). Backtest overfitting quantifies how searching many rules almost guarantees exceptional in-sample performance for at least one of them, even when no true signal exists, and how this interacts with implementability (Bailey et al., 2014; López de Prado, 2018).

A prior  $\pi_{\text{null}}$  value chosen with implementability in mind (for example, higher  $\pi_{\text{null}}$  when microcaps or high turnover are involved) pre-discounts fragile signals. It then requires stronger out-of-sample posterior evidence to overturn this skepticism. This pairing of theory- and evidence-based priors with replication-period data yields a decision metric that is both probabilistically interpretable and economically grounded.

## Integrating $\pi_{\text{null}}$ with the Decision Rule

The Bayesian decision rule reports a discovery only when the posterior probability of the null is at or below a target threshold  $p^c$

$$P(H_{\text{null}} \mid \text{data}) \leq p^c$$

where

$$P(H_{\text{null}} \mid \text{data}) = \frac{1}{1 + \frac{1-\pi_{\text{null}}}{\pi_{\text{null}}} BF_{\text{alt-null}}}.$$

So it combines prior skepticism with likelihood evidence in a single interpretable number. As the search space expands or dependence strengthens, users can increase  $\pi_{\text{null}}$  to reflect lower base rates of true effects or greater concern about overlap. As implementability worsens,  $\pi_{\text{null}}$  can be raised for the affected anomaly families. The rule then updates these prior odds with replication-period statistics, producing a posterior null probability that is transparent about assumptions and comparable across studies.

## Practical Calibration

A practical recipe is to map a family-level replication base rate  $b$  to  $\pi_{\text{null}} \approx 1 - b$ , then adjust by theory strength and implementability. Examples include that stronger economic rationale or robust risk-based stories warrant slightly lower  $\pi_{\text{null}}$ ; crowded or overlapping strategy families, high turnover, microcap reliance, or tight capacity constraints warrant higher  $\pi_{\text{null}}$ .

## Concluding Remark

Data-driven, machine-learning models are powerful for prediction, yet they do not by themselves formalize the inferential costs of extensive search, dependence across strategies, or realistic trading frictions. The prior null probability,  $\pi_{\text{null}}$ , operationalizes ex-ante skepticism about existence of an anomaly, and then lets the data determine when the prior skepticism has been overcome. The result is a decision-relevant posterior probability of a false discovery that is transparent, reproducible, and aligned with the substantive knowledge built by the finance literature on data-snooping, dependence, and implementability.

## Reference

Bailey, D.H., Borwein, J., Lopez de Prado, M. and Zhu, Q.J., 2017. The probability of backtest overfitting. *J. Comput. Finance* 20, 39–69.

De Prado, M.L., 2018. *Advances in financial machine learning*. John Wiley & Sons.

Hansen, P.R., 2005. A test for superior predictive ability. *J. Bus. Econ. Stat.* 23 (4), 365–380.

Sullivan, R., Timmermann, A., White, H., 1999. Data-snooping, technical trading rule performance, and the bootstrap. *J. Finance*, 54 (5), 1647–1691.

Welch, I., Goyal, A., 2008. A comprehensive look at the empirical performance of equity premium prediction. *Rev. Financ. Stud.*, 21 (4), 1455–1508.

White, H. 2000. A reality check for data snooping. *Econometrica*, 68 (5), 1097–1126.