

Algorithmic Trading

Arthur Li

December 13, 2025

Contents

1	Preface and Prerequisites	1
1.1	Brief Overview	1
1.2	Reading Roadmap	1
1.3	Overview of Systematic Investments	2
1.3.1	Alpha Models Overview	3
1.3.2	Risk Models	4
1.3.3	Transaction Cost Models	5
1.3.4	Portfolio Construction Models	5
1.3.5	Execution Model	6
1.3.6	Research	7
1.3.7	Risk Assessment	8
1.4	Exploratory Data Analysis	10
1.4.1	Data Taxonomy	10
1.4.2	Financial Data Structures	14
1.4.3	Data Labelling Techniques	18
1.4.4	Data Sample Weights	20
1.4.5	Fractionally Differentiated Features	21
2	Research	24
2.1	Research in AI	24
2.1.1	Reading AI Papers	24
2.1.2	Writing AI Papers	25

1 Preface and Prerequisites

1.1 Brief Overview

Lorem Ipsum

To be completed once most of the book is done

1.2 Reading Roadmap

Lorem Ipsum. To be completed once most of the book is done

This content builds upon the foundational works of Rishi K. [Narang \(2013\)](#), Raja [Velu \(2020\)](#), and Marcos Lopez [Prado \(2018\)](#), among others, whose insights form the backbone of our discussion.

1.3 Overview of Systematic Investments

A schematic of a live 'production' trading strategy is shown below, but does not include everything else necessary to create the strategy (i.e., research tools).

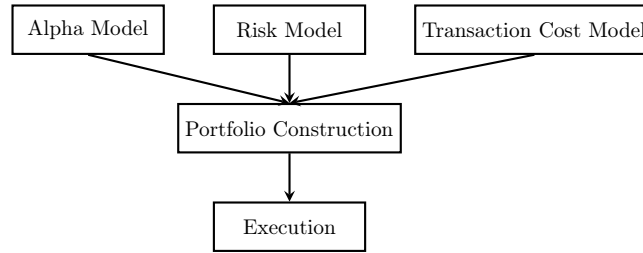


Figure 1: Live Production Trading Strategy Overview

At its core, the trading system is organised into three primary modules:

The trading system has three modules:

- i. Alpha model: predicts the future of the instruments considered for trading, i.e. directional alpha
- ii. Risk model: limits amount of exposure to factors that are unlikely to generate returns but could drive losses, i.e. directional exposure limit on an asset class
- iii. Transaction cost model: determine if the cost of the trades needed to migrate from current portfolio to new portfolio is desirable to the portfolio construction model.

These models feed into a portfolio construction model that balances the tradeoffs of profit and risk to determine the best portfolio to hold. The model finds the differences in trades that need to be executed.

The execution model then takes the required trades, and using inputs such as urgency in which the trades need to be executed and dynamics of liquidity in the markets, executes the trades in an efficient and low cost manner.

Method 1.3.1. *Chains of Production for Alpha Signals*

- i. Data Curation: for collecting, cleaning, indexing, storing, adjusting, and delivering all data to production chain. Requires experts in market microstructure and data protocols such as FIX.
- ii. Feature Analysis: transform raw data into informative signals. Requires experts in information theory, signal extraction and processing, visualisation, labelling, weighting, classifiers, feature importance techniques. Feature analysts collect and catalogue libraries of findings.
- iii. Strategists: informative features are transformed into actual investment algorithms. Strategists will parse libraries of features for ideas to develop an investment strategy. Require data scientists with deep knowledge of financial markets and economy. Features may be discovered by black box, but strategy is developed in a white box.
- iv. Back-testers: assess profitability of investment strategy under various scenarios. Requires data scientists with deep understanding of empirical and experimental techniques. Good back-tested incorporates in analysis meta-information on how strategy was created.
- v. Deployment Team: integrate strategy code into production line. Requires algorithm specialists and mathematical programmers. To ensure deployed solution is logically identical to prototype, and to optimise implementation sufficiently such that production latency is minimised.
- vi. Portfolio Oversight: once strategy is deployed, follows lifecycle.
 1. Embargo: initially, strategy is run on data observed after end date of backtest. If embargoed performance is consistent with backtest, strategy is promoted to next stage.
 2. Paper Trading: strategy run on live, real-time feed. Performance accounts for data parsing latencies, calculation latencies, execution delays, and other time lapses between observation and positioning.
 3. Graduation: strategy manages real position, whether in isolation of as part of ensemble. Performance evaluated precisely, including attributed risk, returns, and costs.
 4. Re-allocation: based on production performance, allocation is re-assessed frequently and automatically. Strategy allocation follows a concave function, Initial allocation is small. As time passes and strategy performs as expected, allocation is increased. Over time, performance decays and allocations become gradually smaller.
 5. Decommission: if strategy perform below expectations for sufficiently extended period of time, strategy is discontinued.

1.3.1 Alpha Models Overview

Theory-driven models tests theories of why markets behave in a manner, and see if they can be used to predict the future. Strategies utilising price-related data are trend and mean reversion; strategies utilising fundamental data are value/yield, growth and quality. Usually more than one model is used in combination.

Definition 1.3.2. *Theory Driven Models*

- i. Trend Following: markets move in given direction long enough that the trend can be identified. As more data support the bull/bear thesis in an uncertain market, more market participants will adopt the same thesis and hence move the asset price to a new equilibrium.
Moving average crossover indicator strategy has less than one point of return for every point of downside risk taken, as market behaviour are unstable and episodic.
- ii. Mean Reversion: markets move in opposite direction to the prevailing trend. Short-term imbalances between buyers and sellers due to liquidity forces prices to move abruptly in one direction, which increases probability of trend reversion as liquidity issue is resolved.
Statistical arbitrage bets on convergence of prices of similar stocks whose prices have diverged.
Longer-term trends can occur despite smaller oscillations around these trends occurring in the shorter term, hence both strategies may be used in conjunction.
- iii. Value/Yield: value strategies uses ratios of fundamental factor against the price of the instrument, inverted to keep the ratio consistent. The higher the yield, the cheaper the instrument.
Buying undervalued security and selling overvalued security is a *carry trade*. The difference between yield received and yield paid is the *carry*.
Quant Long Short (QLS) ranks stocks by attractiveness based on various factors such as value, then buy the higher-ranked stocks while shorting the lower-ranked stocks.
- iv. Growth: make predictions based on asset's expected or historically observed level of economic growth. Forward-looking growth expectations are typically used as a metric.
Growth is trending, and strongest growers are becoming more dominant relative to competitors. Macro growth factors may be used on foreign exchange, while micro growth factors may be used on companies.
- v. Quality: All else being equal, it is better to long high quality and short low quality. Capital safety is important. Factors include earnings quality, equity-to-debt ratios etc.

Data-driven models are more difficult to understand, with more complicated mathematics. Relies on data mining, more technically challenging and far less widely practiced. Typically more used in high-frequency space, as they can discern how market behaves without caring about the economic theory or rational.

Method 1.3.3. *Strategy Parameters*

An implementation approach requires a forecast target, time horizon, bet structure, investment universe, model specification, and run frequency.

- i. Forecast Target: models may forecast direction, magnitude, duration of move, and may include probability into the forecast. Signal strength is of importance, defined by a larger expected return and/or higher likelihood of return. A higher level of signal strength results in a bigger bet taken on the position.
- ii. Time Horizon: models may have forecast horizons ranging from microseconds to years. There are more variability between short-term and long-term strategies, as short-term strategies are making very large number of trades compared to long-term strategies.
- iii. Bet Structure: models can be made to forecast an instrument relative in itself or to others. For relative forecasts, smaller clusters (pairs) or larger clusters (sectors) may be used. For pairs, few assets can be compared precisely and directly. Large cluster grouping may eliminate impact of general movement of the sector and hence focus on the relative movement of stocks within the sector, allowing for clearer distinction between group behaviour and idiosyncratic behaviour. Clusters may be created either via statistical methods or using heuristics (i.e., fundamentally defined industry groups).
Statistical methods may be fooled by data, leading to bad grouping. Heuristic grouping may be imprecise for conglomerates, and may be too rigid. Relative alpha strategies tend to exhibit smoother returns during normal times than intrinsic alpha strategies, but may face incorrect groupings during stressful periods. This may be mitigated by utilising several grouping techniques in concert.
- iv. Investment Universe: choices made on geography, asset class, instrument class, and exclusions. Liquidity is preferred so estimations of transaction costs are reliable. Large quantities of high quality data is required, which is found in highly liquid and developed markets. Instruments with consistent behaviour is preferred, hence biotech stocks are excluded due to sudden, violent price changes. Hence, the most common asset classes and instruments modelled are common stocks, futures (on bonds and equity indices) and forex.

- v. Model Specification: focuses on definition of the strategy mathematically, and may be the source of alpha. Specification details in terms of machine learning or data mining techniques are also defined, to assist in fitting models to the data and setting parameter values. Refitting frequency is also defined to refresh the model and make it adapt to current market conditions; may lead to greater risk of overfitting.
- vi. Run Frequency: defined from monthly to real time frequency. Increasing frequency of runs lead to greater number of transactions and hence higher transaction costs, and risk of moving portfolio based on noisy data. Less frequency of runs lead to smaller number of larger-sized trades, hence may move the market with block trades; may also miss opportunities to trade at more favourable prices.

Method 1.3.4. *Blending of Models*

Most common approaches are linear models, nonlinear models, and machine learning models. If models are not combined, then several portfolios are constructed based on output from each model, then combined using portfolio construction techniques. The best method depends on the model.

- i. Linear Models: require independence of factors, and each factor to be additive. To determine the weight of each alpha factor, multiple regression techniques may be used.
- ii. Nonlinear Models: used when factors are not independent, or the relationship changes over time. Conditional models base the weight of one factor on the reading of another factor. Rotational models assign weights of factors that fluctuate over time based on updated calculations of the various signal's weights, giving higher weights to factors with better performance recently.
- iii. Machine Learning Models: developing machine learning strategies takes as much effort to produce one true investment strategy as to produce a hundred. The complexities include data curation and processing, HPC infrastructure, software development, feature analysis, execution simulations, backtesting etc. Decades ago, macroscopic alpha based on simple tools like econometrics are common, but this is quickly diminishing. Microscopic alpha however, becomes more abundant, but requires heavy ML tools.

1.3.2 Risk Models

Risk model concerns the intentional selection and sizing of exposures to improve the quality and consistency of returns. By pursuing an alpha, we want to be invested in the movement of the exposure to profit in the long run.

Method 1.3.5. *Factor-Based Models*

Factor-based models decompose asset returns into contributions from systematic factors and idiosyncratic components. The most common factors include:

- i. Market Factor: Captures the overall movement of the market.
- ii. Size Factor: Reflects the differential risk associated with companies of varying market capitalizations.
- iii. Value Factor: Accounts for risk due to discrepancies between market prices and fundamental valuations.
- iv. Momentum Factor: Measures the tendency of asset prices to continue in their current trajectory.

This allows traders to understand which elements drive portfolio risk and adjust exposures accordingly.

Method 1.3.6. *Statistical Models*

Statistical risk models leverage historical data and probabilistic techniques to quantify risk parameters.

- i. Historical Simulation: Directly computing risk metrics from past return distributions.
- ii. Monte Carlo Simulation: Generating a large number of potential future return scenarios to estimate risk under diverse conditions.
- iii. Parametric Methods: Employing analytical formulas based on assumed return distributions to calculate key risk measures.

These are useful for dynamically updating risk assessments as new market data become available.

Method 1.3.7. *Limiting Size of Risk*

The quantitative risk models that limit the size of risk varies by the manner in which size is limited, how risk is measured, and what is having its size limited.

Size limits can be limited by hard constraints and penalties. A hard limit may be arbitrary, hence penalty functions may be built to allow a position to increase beyond the limit level, only if the alpha model expects a significantly larger return. The levels of limits and penalties may be determined from either theory or data.

To measure risk, there are two methodologies. The first is longitudinal, and measures risk through the volatility of an instrument. The second is to measure the correlation or covariance between assets (dispersion).

Size limiting may be applied to single positions and groups of positions (sectors, asset classes). It may also be applied to various types of risks and the amount of portfolio leverage.

Method 1.3.8. *Limiting the Types of Risk*

To eliminate unintentional exposure as there is no expectation of being compensated sufficiently for accepting them. This can be achieved through theoretical or empirical risk models.

- i. Theory-Driven Risk Models: focuses on systematic risk factors, derived from economic theory. Systematic risks cannot be diversified away. Equity may have market risk, sector risk, market capitalisation risk etc. Fixed income may have interest rate risk.
- ii. Empirical Risk Models: uses historical data to determine the unnamed systematic risks that should be measured and mitigated. Uses principal component analysis (PCA) to discern unnamed systematic risks that may correspond to named risk factors. Used by statistical arbitrage traders who are betting on exactly the component of an asset's return not explained by systematic risks.

1.3.3 Transaction Cost Models

Trade is made only if it increases the odds or magnitude of return (from alpha model), or if it decreases the odds or magnitudes of loss (from risk model). However, this improvement should be higher than cost of trading. The transaction cost model is not designed to minimise cost of trading, only to inform portfolio construction engine the cost of making any given trade.

Remark 1.3.9. *Transaction Cost Components*

- i. Commissions and Fees: paid to brokerages (access to other market participants), exchanges (improved transaction security) and regulators (operational infrastructure) for the services provided. The bank's infrastructure is used by quants, where the brokerage commissions are rather small on a per-trade basis. Brokers also collect clearing and settlement fees. Clearing is the activity involving regulatory reporting and monitoring, tax handling, and handling failure, taken place in advance of settlement. Settlement is the delivery of securities in exchange for payment in full.
- ii. Slippage: the change in price between the time the quant system decides to transact and the time when the order is at the exchange for execution. Trend-following strategies suffer most from slippage as assets are already moving in desired direction; mean-reverting strategies suffer the least from slippage. The lower the latency to market, the smaller the slippage. The more volatile an asset, the bigger the slippage.
- iii. Market Impact: measures how much an order moves the market by its demand for liquidity. The impact of the trade on the market is unknown until the trade has already been completed. There may also be interaction between slippage and market impact (i.e., selling when a stock is trending upwards).

Definition 1.3.10. *Types of Transaction Cost Models*

- i. Flat Model: cost of trading is the same, regardless of size of order. Model is reasonable if size traded is nearly always about the same, and liquidity remains sufficiently constant.
- ii. Linear Model: cost of trading increases at a constant rate relative to size of order. Better estimate than flat transaction cost model.
- iii. Piece-Wise Linear Model: using piece-wise linear functions to model costs. Balance between simplicity and accuracy; better accuracy than flat or linear models.
- iv. Quadratic Model: most computationally intensive, but also most accurate.

1.3.4 Portfolio Construction Models

Comes in two major forms: rule-based, optimisers. Rule-based models are based on heuristics, can be exceedingly simple or rather complex, and derived from human experience (trial and error). Optimisers comprises of an objective function and uses algorithms to reach the end goal.

Definition 1.3.11. *Rule-Based Models*

- i. Equal Position Weighting: used if portfolio manager believes that if a position is good enough to own, no other information is needed in determining its size. Strength of signal is not used as input in weighting. Model assumes that there is sufficient statistical strength and power to predict not only direction but also magnitude relative to other forecasts in the portfolio. Portfolio takes few large bets on 'best' forecast, many smaller bets on less dramatic forecasts; may take excess risk in an idiosyncratic event on a seemingly attractive position, resulting in adverse selection bias.
- ii. Equal Risk Weighting: adjust position sizes inversely to volatilities or a measure of risk. More volatile positions given smaller allocations, less volatile positions given larger allocations. When unit of risk is equalised, it is almost always a backward-looking measurement such as volatility. If volatility changes with time, then model will be misled.

- iii. Alpha-Driven Weighting: position size based primarily on alpha model. Alpha signal determines size of position, but usually with size limits. Constraints used also includes limits on size of total bet on a group. May also have a function that relates the magnitude of forecast to size of position. If model used in futures trend following, might suffer sharp drawdowns. Reliance on accuracy of alpha.
- iv. Decision-Tree Weighting: decision path to arrive at the allocation for given instrument, depending on type of alpha model and type of instrument. Constraints may include percentage limits for allocation. Model size grows dramatically if more alpha models or more types of positions are included.

Remark 1.3.12. *Optimisers Models Parameters*

Harry Markowitz's mean variance optimisation (MVO) as the pioneer model. Models are based on principles of modern portfolio theory (MPT). Inputs include asset expected return (mean), asset variance, expected correlation matrix. Other inputs include size of portfolio in currency terms, desired risk level (volatility or expected drawdown), and other constraints such as liquidity, universe limits.

Model uses an objective function and an algorithm to seek the goal, usually maximising return of portfolio relative to volatility of portfolio returns.

- i. Expected Return: alpha models as basis of expected return, which also includes expected direction.
- ii. Expected Volatility: stochastic volatility forecasting methods is commonly used, as volatility may have high and low periods, with occasional jumps. GARCH model is most used.
- iii. Expected Correlation: as instrument correlations are not stable over time, it is more appropriate to group assets together before computing correlation within the group.

Method 1.3.13. *Optimisation Techniques*

- i. Unconstrained Optimisation: most basic form with no constraints. Might provide a single-instrument portfolio, where all money will be invested in instrument with highest risk-adjusted return.
- ii. Constrained optimisation: constraints include position limits, limits on various groupings of instruments. Might result in constraints driving the portfolio construction more than the optimiser.
- iii. Black-Litterman Optimisation: blends investor expectations with a degree of confidence about those expectations, and these with historical precedent evident in the data. Adjusts historically observed correlation levels by utilising investor's forecast of return for the various instruments.
- iv. Grinold and Kahn's Approach: builds a portfolio of signals, instead of sizing positions. To build factor portfolios, each of which are usually rule-based portfolios based on a single type of alpha forecast. Each portfolio backtested, then series of returns are then treated as instruments of a portfolio by the optimiser. Number of factor portfolios is more manageable, usually not more than 20. What is optimised is then a handful of factor portfolios. The model allows for inclusion of risk model, transaction cost model, portfolio size, and risk targets as inputs.
- v. Resampled Efficiency: to improve the inputs to optimisation by addressing oversensitivity to estimation error. To resample data using Monte Carlo simulation to reduce estimation error in inputs to the optimiser.
- vi. Data-Mining Approaches: machine learning techniques such as supervised learning or genetic algorithms used, as MVO involves searching many possible portfolios to find the best.

1.3.5 Execution Model

Two basic ways to execute trade: through electronic, or through human intermediary. For electronic execution, achieved through direct market access (DMA), which allows traders to utilise the infrastructure and exchange connectivity of brokerage firms to trade directly on electronic markets.

Execution algorithms can be acquired through building, using broker's, or a third-party software vendors.

Brokerages offer portfolio bidding, where the 'blind' portfolio for transaction is described by characteristics such as valuation ratios of longs and shorts, sector breakdown, market capitalisation etc. Broker then quote a fee in basis points in terms of the gross market value of portfolio traded. Hence, certainty is provided by the broker to the trader. Once agreement reached, broker receives fee and assumes risk of trading out the portfolio at future market prices, which may be better or worse than prices guaranteed.

Remark 1.3.14. *Order Execution Algorithm Parameters*

- i. Aggressive vs Passive: algorithm make decision of passive vs aggressive order, depending on how immediately the trader wants to do the trade. Market orders are considered aggressive. Limit order at current best order is fairly aggressive, while limit order below current bid is passive. Many exchanges pay providers of liquidity for placing passive orders, charging traders for using liquidity provided. Orders that cross the spread are using liquidity by using a passive order placed by another

trader, reducing liquidity available. Paying for liquidity sweetens deal for passive order, only if order is actually executed; passive trader gets better transaction price and a commission rebate from the exchange. Momentum strategies uses more aggressive orders; mean reversion uses more passive orders. A stronger, more certain signal will be executed with greater aggressiveness than a weaker or less certain signal. A middle ground will be to put limit orders between best current bid and offer.

- ii. Large vs Small Order: a large order may be broken into many smaller orders over a window of time, but risk price moving in adverse direction. Size of chunk depends on transaction cost model estimate, and analysis of correct level of aggressiveness.
- iii. Hidden vs Visible Order: a queue as a visible order gives away a bit of information. Hidden order will provide no information to the market, staving off imbalances, but reduces priority of trade in the queue. Algorithmic trading utilising hidden order is 'iceberging', which is taking a single larger order and chopping it into many smaller chunks, most posted to order book as hidden orders.
- iv. Order Routing: if there are several pools of liquidity for the same instrument, smart order routing will be used, which determines which pool of liquidity is most suitable for sending a given order. Depth of liquidity on various ECNs and connectivity speeds are also considered in smart order routing.
- v. Cancelling and Replacing Orders: traders may place larger number of orders with no intention of execution, then rapidly cancelling them and replacing them with other orders. This allows gaining of information on how market responds to the changing depth of the book, providing information on how to profit from the pattern of reaction. If trader wants to buy a large number of shares, he may enter a large number of small orders to sell the shares further away from market and cancel, improving market perception.

Definition 1.3.15. *High Frequency Trading*

Alpha driving strategies on extremely near-term bets (seconds or less) are *microstructure alphas*, focusing on liquidity patterns in order book. Larger quants may also use this to guide execution models, improving costs of entering trades. Small differences over a single trade add up significantly in the long run. To trade microstructure alpha as independent high frequency strategies, large investments in infrastructure and research must be done. Machine learning techniques may also be used to discern patterns in execution of other player orders. The more inferior the execution models, the easier it is to discern the pattern, allowing the ML strategy to profit from these patterns in the future. Patterns in the shorter timescale are somewhat stable.

Definition 1.3.16. *HFT Shark Strategy*

Designed to detect large orders that are iceberged, by sending series of very small trades; if each of these small orders get filled quickly, this may be a sign of a large and iceberged order. The shark simply front-run this large, hidden order by placing visible trades in front of the iceberged order. The iceberg strategy must then push prices up to execute trades. When the iceberged order is complete, prices will be pushed up favourably for the shark, which can then exit the position with a quick and relatively riskless profit.

Remark 1.3.17. *HFT Trading Infrastructure*

Using a broker that act as trading agent allows the infrastructure requirements to be handled by the broker, instead of dealing with the regulatory and other constraints.

High frequency strategies may use colocation or sponsored access. Colocation setup is where trader attempts to place trading servers as physically close to the exchange as possible.

Financial Information eXchange (FIX) protocol is the choice of real-time electronic communication among users. The software that implements the FIX protocol is free and open source (FIX engine). High frequency traders will likely build their own FIX engines to ensure optimal speeds.

1.3.6 Research

Definition 1.3.18. *Scientific Method*

1. Researcher observe a phenomenon in the market and construct a theory.
2. Researcher seeks out information to test the theory.
3. Researcher tests the theory, and with enough confidence, risk some capital on the validity of the theory.

Remark 1.3.19. *Sources of Alpha Idea Generation*

1. Observing the market, using the scientific method to test the theory
2. Academic literature, requiring significant time to read academic journals, working papers, and conference presentations for ideas. Literature from other fields such as astronomy, physics, or psychology, may provide ideas relevant to quant finance problems.
3. Migration of a researcher or portfolio manager from one quant shop to another.

4. Lessons from activities of discretionary traders

Remark 1.3.20. *Model Quality Assessment*

- i. Cumulative profit graph: if profit profile is not smooth, with long periods of inactivity, sharp losses and gains, then the model may have issues
- ii. Average annual rate of return: indicates how well the strategy made on historical data
- iii. Variability of returns: the less variable the level of returns, the better the strategy. May look at lumpiness of returns, which is the portion of strategy's total returns that comes from periods that are significantly above average (measures consistency of returns).
- iv. Worse Peak-to-Valley Drawdowns: measures maximum decline from any cumulative peak in profit curve. The lower the drawdown the better the strategy. Also, to measure recovery period after drawdowns; the shorter the recovery period the better the strategy.
- v. Predictive Power: R-squared statistic may be used, which shows how much of the variability of the predicted asset have been accounted for. A exceedingly high R^2 in would be 0.05 out of sample. Instrument returns may be bucketed by deciles; a model with reliable predictive power is one that appropriately buckets the instruments correctly.
- vi. Percentage Winning Trades, Winning Time Periods: whether the strategy tends to make profits from a small portion of trades that do very well, or from a large number of trades.
- vii. Ratios of Returns vs Risk: Statistics such as risk-adjusted return, Sharpe ratio, information ratio, Sterling ratio, Calmer ratio, Omega ratio.
- viii. Relationship with Other Strategies: value-add of new strategy compared with results of existing strategy with and without the new idea.
- ix. Time decay: understand strategy returns if trades are initiated on lagged basis after receiving a trading signal. Determine strategy sensitivity to timeliness with information received, and crowdedness of strategy.
- x. Sensitivity to specific parameters: high quality strategy has small changes in outcomes from slight changes in parameters. Or else this may be a sign that model may be overfitted.
- xi. Overfitting: plot a graph of parameter value vs function outcome; a good model has a flatter curve with no jumps. Models that are parsimonious (less parameters) uses less assumptions, hence less overfitting.

Remark 1.3.21. *Other Considerations in Model Testing*

Overestimation of trading costs may cause portfolio to hold positions for longer than optimal, and underestimation may result in high portfolio turnover and bleed from trading costs. Assumptions on availability of short positions must also be made; hard-to-borrow lists must be taken into consideration.

1.3.7 Risk Assessment**Definition 1.3.22.** *Model Risks*

Quant models has model risk, the risk that the model does not accurately describe, match, or predict the real-world phenomenon. Each component of the quant model may all have model risk.

- i. Inapplicability of Modelling: occurs when quant model is mistakenly applied to a problem. May also occur with misapplication of a technique to a given problem.
- ii. Model Misspecification: occurs when the model doesn't fit the real world. Model may work fine most of the time, but fail when an extreme event occurs.
- iii. Implementation Errors: errors in programming or architecting systems. Architectural error may also occur when models are loaded in a wrong sequence.

Definition 1.3.23. *Regime Change Risk*

Quant models are based on relationships prevalent in historical data. If there is a regime change, the historical relationships and behaviour may be altered, hence the model may lose effectiveness.

Definition 1.3.24. *Exogenous Shock Risk*

Risks driven by information that is not internal to the market, i.e., terrorist attacks, start of wars, bank bailouts, change in regulation such as in shorting rules. May require discretionary overrides.

Definition 1.3.25. *Contagion Risk*

Happens when other investors hold the same strategies. First part of risk factor relates to how crowded the quant strategy is. Second part relates to what else is held by other investors that could force them to exit the quant strategy in a panic (ATM effect).

Quant liquidation crisis may be driven by size and popularity of quantitative strategies, subpar returns from operators leading up to the crisis, the practice of funds cross-collateralising many strategies against each other, and risk targeting (risk managers target a specific level of volatility for their funds or strategies).

Method 1.3.26. *Risk Monitoring Methods*

- i. Exposure Monitoring Tools: with current positions held, the positions are grouped for the various exposures (i.e., valuation, momentum level, volatility) to monitor gross and net exposure to various sectors and industries, buckets of market capitalisation, various style factors.
- ii. Profit and Loss Monitors: with current portfolio, compare that with previous day closing price. Intraday performance charts are used. May also look at source of profit, hit rate (percentage of time strategy makes money on a given position).
- iii. Execution Monitors: shows progress of executions, i.e., which orders are currently being worked on, which ones are completed, with transaction size and prices. Fill rates for limit orders are used for more passive execution strategies. Slippage and market impact are also monitored.
- iv. System Performance Monitors: checks for software and infrastructure errors. Checks performance of CPUs, speeds of various stages of automated processes, latency in communication of messages.

1.4 Exploratory Data Analysis

1.4.1 Data Taxonomy

A brief overview of the types of data used in systematic trading.

Four essential types of financial data

Fundamental Data	Market Data	Analytics	Alternative Data
Assets	Price/Yield/IV	Analyst Recommendation	Satellite/CCTV
Liabilities	Volume	Credit Ratings	Google Searches
Sales	Dividend/Coupons	Earnings Expectations	Twitter/Chats
Costs/Earnings	Open Interest	News Sentiment	Metadata
Macro Variables	Quotes/Cancellations
...	Aggressor Side		
	...		

Remark 1.4.1. *Fundamental Data Characteristics*

- Data published is indexed by last date included in report, which precedes date of release.
- Data is often backfilled or re-instated, and data vendor may overwrite initial values with corrections.
- Data is extremely regularised and low frequency.

Remark 1.4.2. *Market Data Characteristics*

- Raw feed contains unstructured information, such as FIX messages (allow full construction of trading book), or full collection of BWIC (bids wanted in competition) responses.
- FIX data is not trivial to process, $\sim 10\text{TB}$ generated on daily basis

Remark 1.4.3. *Analytics Data Characteristics*

- Derivative data as processed based on original source. Signal already extracted from the original source.
- Costly, methodology used in production may be biased or opaque.

Remark 1.4.4. *Alternative Data Characteristics*

- Produced by individuals, business processes, and sensors.
- Primary information that has not made it to other sources.
- Cost and privacy concerns. May be useful if it annoys data infrastructure team.

Definition 1.4.5. *Reference Data*

- Trading Universe: evolving daily to incorporate new listings, de-listings etc. Knowing when a particular stock no longer trades is important to avoid survivor bias.
- Symbology Mapping: ISIN, SEDOL, RIC, Bloomberg Tickers etc. Data is not static, symbols may change, complicating historical data merges. Mapping needs to persist as point-in-time data and allow for historical 'as-of-date' usage, require implementation of bi-temporal data structure.
- Ticker Changes: for reasons described in symbology mapping. To maintain historical table of ticker changes to seamlessly go up and down time series data.
- Corporate Actions Calendars: contain stock and cash dividends (announcement, execution date), stock splits, reverse splits, rights offer, mergers and acquisitions, spin off, free float or shares outstanding adjustments, quotation suspensions etc.
For dividends, announcements may coincide with more volatility, jumps in price time series. Allow building of strategies that look to benefit from the added volatility.
For stock splits, reverse splits, rights offers, all historical data need to be adjusted backward to reflect the split (both volume and price).
For M&A, spin-offs, to account for changes in valuation, hence used in Merger Arbitrage strategies.
Suspensions result in gaps in data, may impact backtesting.
- Static Data: country, sector, primary exchange, currency, and quote factor. May be used to group instruments based on fundamental similarities (hence for pairs trading). Maintaining a table of quotation currency per instrument necessary to aggregate positions at portfolio level.
- Exchange Specific Data: individual exchanges have variety of differences to be accounted for when designing trading strategies. First group of information concerns the hours and dates of operations:

1. Holiday Calendar: Strategies trading simultaneously in several markets and leveraging correlation may not perform as well if one market is closed and another is open.
2. Exchange Sessions Hours: Different sessions (Pre-Market, Continuous Core, After-Hour etc.); auction times and respective cutoff times for order submission; lunch break restricting intraday trading and auctions before/after lunch; settlement times for futures market. Daylight Saving Time (DST) adjustments; length of trading hours during course of the year; different trading hours by venue.
3. Disrupted Days: Exchange outages or trading disruptions, market data issues. To be recorded so they can be filtered out when building or testing strategies.

Second group of information governing the mechanics of trading:

1. Tick Size: Minimum eligible price increment; may vary by instrument and as a function of price.
 2. Trade and Quote Lots: Minimum size increment for quotes or trades.
 3. Limit-Up and Limit-Down Constraints: Maximum daily fluctuations of securities, and whether trading is paused or can only be traded at better prices than the threshold.
 4. Short Sell Restrictions: Restrict short sells not to trade at price worse than last price, or not to create a new quote that will be lower than the lowest prevailing quote. Impact ability to source liquidity.
- vii. Market Data Condition Codes: vary per exchange and asset class, and each market event may be attributed to several codes at once. To build mapping table of condition codes and what they mean (i.e., auction trade, lit or dark trade, cancelled or corrected trade, regular trade, off-exchange trade reporting, block-size trade, trade originating from multi-leg order such as option spread trade etc.). To access liquidity for trading algorithm, trades published for reporting purposes must be excluded and not be used to update some of the aggregated daily data used in construction of trading strategies.
- viii. Special Day Calendars: days with distinct liquidity characteristics to be accounted for in both execution strategies and in alpha generation process. These (non-exhaustive) irregular events may be:
1. Half trading days preceding Christmas and following thanksgiving in US
 2. Ramadan even in Turkey
 3. Taiwan market opening on weekend to make up for lost trading days during holiday periods
 4. Korean market changing trading hours on day of nationwide university entrance exam
 5. Brazilian market opening late on day following the Carnival
 6. Last trading days of months and quarters (investors rebalance portfolios)
 7. Index rebalancing dates, where intraday volume distribution is significantly skewed toward EODs
 8. Options and futures expiry dates (quarterly/monthly expiry, Triple Witching in US, Special Quotations in Japan) where excess trading volume and different intraday patterns result from hedging activity and portfolio adjustments.

Model normal days first. Special days are modelled either independently, or using normal days as baseline.

- ix. Futures-Specific Reference Data: to know which contract was live at any point of time by using expiry calendar, and the most liquid contract. Equity index futures are most liquid for first contract available (front month), energy futures such as oil are more liquid for second contract. Hence to know which contract carry the most significant price formation characteristics, and what is true liquidity available. Note there is no real standardised expiry frequency that applies across markets. When computing rolling-window metrics, to account for potential roll dates (due to investors rolling forward positions) that may have happened during the time span. May blend volume time series prior to roll date and after roll date. Futures market also have different market phases during the day with significantly different liquidity characteristics. Various market data metrics (volume profile, average spread, average bid-ask sizes etc) should be computed separately for each market phase by maintaining a table of start and end times of each session for each contract.
- x. Options-Specific Reference Data (Options Chain): expiry date and strike price combination (option chain). Map of equity tickers to option tickers with strike and expiry dates allow for design for more complex investment and hedging strategies (i.e., distance to strike, change in open interest of puts and calls).
- xi. Market-Moving News Releases: macroeconomic announcements. To maintain calendar of dates and times of their occurrences to assess their impact on strategies. Central bank announcements or meeting minutes releases about major economies (FED/FOMC, ECB, BOE, BOJ, SNB), Non-Farm Payrolls, Purchasing Managers' Index, Manufacturing Index, Crude Oil Inventories etc. Stock-specific releases such as earnings calendars, specialised sector events (for healthcare, biotech etc).

- xii. Related Tickers: tickers that are related to each other as they fundamentally represent the same underlying asset. Allows efficient opportunity exploitation. Primary tickers to composite tickers mapping (for markets with fragmented liquidity), dual listed/fungible securities in US and Canada, ADR or GDR, local and foreign boards in Thailand etc.
- xiii. Composite Assets: ETFs, Indexes, Mutual Funds etc. May be used to achieve desired exposures, or as cheap hedging instruments, and provide arbitrage opportunities when they deviate from NAV. To maintain information such as time series of their constituents and value of any cash component, divisor used to translate NAV into quoted price, constituent weights.
- xiv. Latency tables: for higher frequency trading strategies. Contains distribution of latency between different data centres for more efficient order routing, and reordering data that are recorded in different locations.

Definition 1.4.6. *Market Data Feed*

- i. Level I Data (Trade and BBO Quotes): trades and top of book quotes. Enough to reconstruct Best Bid and Offer (BBO). Also contains information in form of trade status (cancelled, reported late etc), trade and quote qualifiers (odd lot, normal trade, auction trade, Intermarket Sweep, average price reporting, on which exchange etc). May be used to analyse sequence of events and decide if a given print should be used to update the last price and total volume traded at a point in time.
- ii. Level II Data (Market Depth): addition of quote depth data, displays all lit limit order book updates (price changes, addition or removal of shares quoted) at any level in the book, for all of the lit venues in fragmented markets.
- iii. Level III Data (Full Order View): message data. Each order arriving is attributed a unique ID for tracking over time, and is precisely identified when it is executed, cancelled, or amended. Possible to build a full (with national depth) book at any moment intraday. Example from US market:
 1. Timestamp: number of milliseconds after midnight
 2. Ticker: equity symbol (up to 8 char)
 3. Order: Unique order ID
 4. T: message type. 'B' is add buy order; 'S' is add sell order; 'E' is execute outstanding order in part; 'C' is cancel outstanding order in part; 'F' is execute outstanding order in full; 'D' is delete outstanding order in full; 'X' is bulk volume for cross event; 'T' is execute non-displayed order
 5. Shares: order quantity for 'B', 'S', 'E', 'X', 'C', 'T' messages. Zero for 'F', 'D' messages
 6. Price: order price for 'B', 'S', 'X', 'T' messages. Zero for cancellation and executions. Last 4 digits are decimal, padded to right with zeroes. Divide by 1000 to convert to currency value.
 7. MPID: Market Participant ID associated with transaction (4 char)
 8. MCID: Market Centre Code for originating exchange (1 char)

A few special types of orders worth mentioning are:

1. Order subject to price sliding: execution price may be one cent worse than display price at NASDAQ; ranked at locking price as hidden order, displayed at the price one minimum price variation inferior to locking price. New order ID will be used if order is replaced as a display order.
2. Pegged order: based on NBBO, not routable, new timestamp given upon repricing; display rule vary over exchanges
3. Mid-point peg order: non-displayed, may result in half-penny execution
4. Reserve order: displayed size is ranked as displayed limit order; reserve size is behind non-displayed orders and pegged orders in priority.
Minimum display quantity is 100, amount replenished from reserve size when it falls below 100 shares; New timestamp created, displayed size re-ranked upon replenishment.
5. Discretionary order: displayed at one price while passively trading a more aggressive discretionary price. Order becomes active when shares are available within discretionary price range. Order ranked last in priority. Execution price may be worse than display price.
6. Intermarket sweep order: can be executed without need for checking prevailing NBBO.

Using these data, we may model: the pattern of inter-arrival times of various events; arrival and cancellation rates as a function of distance from nearest touch price; arrival and cancellation rates as a function of other information, such as in the queue on either side of the book, order book imbalance etc.

Once modelled, we may analyse: the impact of market order on limit order book; chances for limit order to move up the queue from given entry position; probability of earning the spread; expected direction of price movement over a short horizon.

Definition 1.4.7. *Binned Data*

- i. Open, High, Low, Close (OHLC) and Previous Close Price: indication on trading activity and intraday volatility. Distance traveled between lowest and highest points is indication of market sentiment. Previous close has to be adjusted for corporate actions and dividends.
- ii. Last Trade before Close (Price/Size/Time): how much the close price may have jumped in final moments of trading; how stable it is as a reference value for next day.
- iii. Volume: trading activity indicator, especially when level jumps from long term average. Collect volume breakdown between lit and dark venues for execution strategies.
- iv. Auctions Volume and Price: price discovery event when significant volume prints occur.
- v. VWAP: indication of trading activity on the day. Easier to algorithmically execute large orders with VWAP than a single print.
- vi. Short Interest/Days-to-Cover/Utilisation: good proxy for investor position. Short pressure an indication of upcoming short term moves: large short interest is bearish view from institutional investors. Utilisation level of available securities to borrow gives indication of how much room is left for future shorting. Days-to-Cover to assess magnitude of potential short squeeze (if sellers unwind position, fraction of available daily liquidity needed); larger value indicates larger potential of sudden upswing on heavily shorted securities.
- vii. Futures Data: insight into activity or large investors through open interest data. Offer arbitrage opportunities if their basis exhibits mis-pricing compared to one's dividend estimates.
- viii. Index-Level Data: source of relative measures for instrument specific features (index OHLC, volatility). Normalised features identify individual instruments deviating from their benchmarks.
- ix. Options Data: information on position of traders through open interest and Greeks.
- x. Asset Class Specific: yield/benchmark rates (repo, 2y, 10y, 30y), CDS spreads, US Dollar Index

Definition 1.4.8. *Granular Intraday Microstructure Activity*

- i. Number and Frequency of Trades: proxy for activity level, and how continuous it is. Low number of trades mean harder execution, and may be more volatile
- ii. Number and Frequency of Quote Updates: similar proxy for activity level
- iii. Top of Book Size; proxy for liquidity of instruments (larger top of book size makes it possible to trade larger order quasi immediately)
- iv. Depth of Book (price and size): similar proxy for liquidity
- v. Spread Size (average, median, time weighted average): proxy for cost of trading. Parametrised distribution used to identify opportunities if they are cheap or expensive
- vi. Trade size (average, median): to identify intraday liquidity opportunities when examining volume available in the order book.
- vii. Ticking time (average, median): representation of how often one should expect changes in the order book first level. For execution algorithms for which the frequency of updates (adding/cancelling child orders, re-evaluating decisions etc.) should commensurate with characteristics of the traded instrument.

Daily distribution can be used as start of day estimates and updated intraday with online Bayesian updates. Last group of daily data is derived from previous two groups but stored pre-computed to save time during research phase, or to be used as normalising values:

- i. X-day Average Daily Volume (ADV) / Average Auction Volume
- ii. X-day Volatility (close-to-close, open-to-close etc)
- iii. Beta with respect to index or sector (plain beta, or asymmetric up-days/down-days beta)
- iv. Correlation matrix

When binning data, this may be grouped into bins ranging from a few seconds to 30 minutes. Minute bar data are used for volume and spread profiles to prevent introducing excess noise due to market friction.

Definition 1.4.9. *Fundamental Data and Other Data*

- i. Key Ratios: Earnings Per Share (EPS), Price-to-Earning (P/E), Price-to-Book (P/B), etc.
- ii. Analyst Recommendations: aggregated values given consensus valuation
- iii. Earnings data: estimations by research analysts provide quarterly earning estimates which can be used as indication of performance of stock before actual value is published

- iv. Holders: sudden changes in ownership indicate changes in sentiment by sophisticated investor
- v. Insiders Purchase/Sale: indicator of future stock price moves from group of people who have access to best possible information about the company
- vi. Credit Ratings: credit downgrades resulting in higher funding costs have negative impact on equity prices

1.4.2 Financial Data Structures

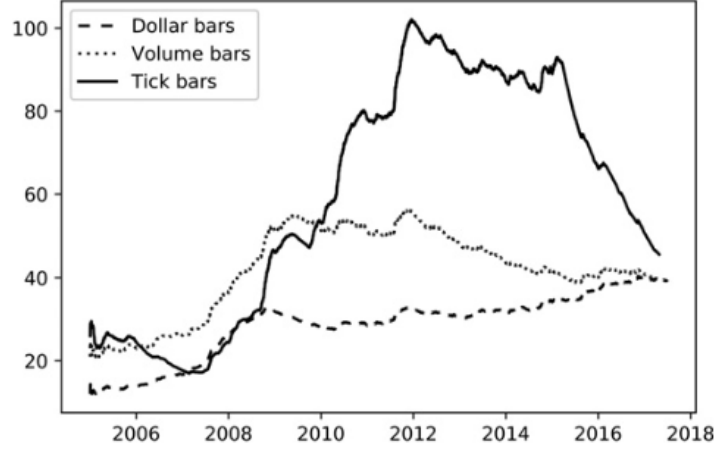


Figure 2: Average daily frequency of tick, volume, and dollar bars

Remark 1.4.10. *Standard BARS*

Method to transform a series of observations arriving at irregular frequency into a homogeneous series derived from regular sampling.

- i. Time Bars: obtained by sampling information at fixed intervals. Information collected includes timestamp, volume-weighted average price (VWAP), open price, close price, high price, low price, volume etc. To be avoided as markets do not process information at constant time interval. Time bars oversample information in low-activity periods and under-sample information in high-activity periods. Time bars exhibit poor statistical properties, i.e., serial correlation, heteroscedasticity, non-normality of returns.
- ii. Tick Bars: sample variables extracted each time a pre-defined number of transactions take place. Allows synchronisation of sampling with a proxy of information arrival. Sampling as a function of trading activity creates returns closer to IID Normal (Thierry and Helyette (2000)). When constructing tick bars, to be aware of outliers, as many exchanges carry out auction at open and at close; order book accumulates bids and offers without matching. Order fragmentation introduces some arbitrariness in number of ticks. Matching engine protocols may split one fill into multiple artificial partial fills as a matter of operational convenience.
- iii. Volume Bars: samples every time a pre-defined amount of security's units that have been exchanged. Achieves better statistical properties than sampling tick bars. Convenient artefact for studying market microstructure theories.
- iv. Dollar Bars: samples an observation every time a pre-defined market value is exchanged. Used when the analysis involves significant price fluctuations. Robust against corporate actions such as splits, reverse splits, issuance of new shares, buying back existing shares. Bar size could be dynamically adjusted as a function of free-floating market cap of a company or outstanding amount of issued debt.

Remark 1.4.11. *Information-Driven Bars*

Method to sample more frequently when new (micro-structural) information arrives to the market.

- i. Tick Imbalance Bars: sample bars whenever tick imbalance exceeds expectations. To determine tick index T such that accumulation of signed ticks exceeds a given threshold. Let $\{(p_t, v_t)\}_{t=1, \dots, T}$ be sequence of ticks where p_t and v_t is the price and volume associated with tick t . Let tick rule define a sequence $\{b_t\}_{t=1, \dots, T}$ where

$$b_t = \begin{cases} b_{t-1} & \text{if } \Delta p_t = 0 \\ \frac{|\Delta p_t|}{\Delta p_t} & \text{if } \Delta p_t \neq 0 \end{cases}$$

The tick imbalance at time T is defined as

$$\theta_T = \sum_{t=1}^T b_t$$

Compute expected value of θ_T at beginning of the bar,

$$E_0[\theta_T] = E_0[T](P[b_t = 1] - P[b_t = -1]) = E_0[T](2P[b_t = 1] - 1)$$

where $E_0[T]$ is expected size of tick bar, $P[b_t = 1]$ and $P[b_t = -1]$ is unconditional probability that a tick is classified as a buy and sell. In practice, $E_0[T]$ and $(2P[b_t = 1] - 1)$ may be estimated as an exponentially weighted moving average of T and b_t values from prior bars.

Define the tick imbalance bar (TIB) as a T^* contiguous subset of ticks such that

$$T^* = \arg \min_T \{|\theta_T| \geq E_0[T] | 2P[b_t = 1] - 1\}$$

where the size of expected imbalance is implied by $|2P[b_t = 1] - 1|$.

When θ_T is more imbalanced than expected, a low T will satisfy the conditions.

TIBs are produced more frequently under presence of informed trading (asymmetric information that triggers one-side trading). TIBs are buckets of trades containing equal amounts of information.

- ii. Volume/Dollar Imbalance Bars: sample bars when volume or dollar imbalances diverge from expectations. First, define imbalance at time T as

$$\theta_T = \sum_{t=1}^T b_t v_t$$

where v_t may represent ether number of securities traded (VIB) or dollar amount traded (DIB).

The expected value of θ_T at the beginning of the bar is then computed as

$$\begin{aligned} E_0[\theta_T] &= E_0 \left[\sum_{t|b_t=1}^T v_t \right] - E_0 \left[\sum_{t|b_t=-1}^T v_t \right] \\ &= E_0[T](P[b_t = 1]E_0[v_t|b_t = 1] - P[b_t = -1]E_0[v_t|b_t = -1]) \\ &= E_0[T](v^+ - v^-) \end{aligned}$$

where the initial expectation of v_t is decomposed into component contributed by buys and sells. Then

$$E_0[\theta_T] = E_0[T](2v^+ - E_0[v_t])$$

In practice, $E_0[T]$ and $(2v^+ - E_0[v_t])$ may be estimated as exponentially weighted moving average of T and $b_t v_t$ values from prior bars. Next, define VIB or DIB as a T^* -contiguous subset of ticks such that

$$T^* = \arg \min_T \{|\theta_T| \geq E_0[T] | 2v^+ - E_0[v_t]\}$$

where the size of expected imbalance is implied by $|2v^+ - E_0[v_t]|$.

When θ_T is more imbalanced then expected, a low T will satisfy the conditions.

VIB and DIB addresses concerns on tick fragmentation and outliers, and also addresses the issues of corporate actions, as the bar size is adjusted dynamically.

- iii. Tick Runs Bars: sample bars when the sequence of buys in overall volume diverges from expectations. For the case when large traders sweep order book, use iceberg orders, or slice parent orders into multiple children, all leaving a trace of runs in the $\{b_t\}_{t=1,\dots,T}$ sequence. Define length of current run as

$$\theta_T = \max \left\{ \sum_{t|b_t=1}^T b_t - \sum_{t|b_t=-1}^T b_t \right\}$$

The expected value of θ_T at beginning of bar is computed as

$$E_0[\theta_T] = E_0[T] \max\{P[b_t = 1], 1 - P[b_t = 1]\}$$

In practice, $E_0[T]$ and $P[b_t = 1]$ may be estimated as exponentially weighted moving average of T and

proportion of buy ticks from prior bars. Next, define TRB as T^* -contiguous subset of ticks such that

$$T^* = \arg \min_T \{ \theta_T \geq E_0[T] \max\{P[b_t = 1], 1 - P[b_t = 1]\} \}$$

where the expected count of ticks from runs is implied by $\max\{P[b_t = 1], 1 - P[b_t = 1]\}$.

When θ_T exhibits more runs than expected, a low T will satisfy these conditions.

Instead of measuring length of longest sequence, count number of ticks of each side without offsetting.

- iv. Volume/Dollar Runs Bars: sample bars when volume or dollars traded by one side exceed expectation for a bar. First, define volume or dollars associated with a run as

$$\theta_T = \max \left\{ \sum_{t|b_t=1}^T b_t v_t - \sum_{t|b_t=-1}^T b_t v_t \right\}$$

where v_t may either represent volume (VRB) or dollar amount exchanged (DRB). The expected value of θ_T at beginning of the bar is then

$$E_0[\theta_T] = E_0[T] \max\{P[b_t = 1]E_0[v_t|b_t = 1], (1 - P[b_t = 1])E_0[v_t|b_t = -1]\}$$

In practice, $E_0[T]$, $P[b_t = 1]$, $E_0[v_t|b_t = 1]$, $E_0[v_t|b_t = -1]$ may be estimated as exponentially weighted moving average of T , proportion of buy ticks, buy volumes, and sell volumes from prior bars. Next, define a volume runs bar (VR) as T^* -contiguous subset of ticks such that

$$T^* = \arg \min_T \{ \theta_T \geq E_0[T] \max\{P[b_t = 1]E_0[v_t|b_t = 1], (1 - P[b_t = 1])E_0[v_t|b_t = -1]\} \}$$

expected volume from runs is implied by $\max\{P[b_t = 1]E_0[v_t|b_t = 1], (1 - P[b_t = 1])E_0[v_t|b_t = -1]\}$.

When θ_T exhibits more runs than expected, volume from runs is greater than expected, a low T will satisfy these conditions.

Definition 1.4.12. Multi-Product Series: ETF Trick

To model a basket of securities as if it was a single cash product. To transform any complex multi-product dataset into a single dataset that resembles a total return ETF.

Method 1.4.13. ETF Trick

Produce a time series that reflects the value of \$1 invested. Changes in the series will reflect changes in PnL, series will be strictly positive, and implementation shortfall will be taken into account. The bars contain:

- i. Raw open price of instrument $i = 1, \dots, I$ at bar $t = 1, \dots, T$: $o_{i,t}$
- ii. Raw close price of instrument $i = 1, \dots, I$ at bar $t = 1, \dots, T$: $p_{i,t}$
- iii. USD value of one point of instrument $i = 1, \dots, I$ at bar $t = 1, \dots, T$: $\varphi_{i,t}$. This includes forex rate.
- iv. Volume of instrument $i = 1, \dots, I$ at bar $t = 1, \dots, T$: $v_{i,t}$
- v. Carry, dividend, or coupon paid by instrument i at bar t : $d_{i,t}$. Variable can also be used to charge margin costs or costs of funding.

All instruments $i = 1, \dots, I$ were tradable at bar $t = 1, \dots, T$. Even if some instruments were not tradable over entirety of time interval $[t-1, t]$, at least they were tradable at times $t-1$ and t .

For basket of securities with allocation vector ω_t rebalanced (or rolled) on bars $B \subseteq \{1, \dots, T\}$, the \$1 investment value $\{K_t\}$ is derived as

$$h_{i,t} = \begin{cases} \frac{\omega_{i,t} K_t}{o_{i,t-1} \varphi_{i,t} \sum_{i=1}^I |\omega_{i,t}|} & \text{if } t \in B \\ h_{i,t-1} & \text{otherwise} \end{cases}$$

$$\delta_{i,t} = \begin{cases} p_{i,t} - o_{i,t} & \text{if } (t-1) \in B \\ \Delta p_{i,t} & \text{otherwise} \end{cases}$$

$$K_t = K_{t-1} + \sum_{i=1}^I h_{i,t-1} \varphi_{i,t} (\delta_{i,t} + d_{i,t})$$

where $K_0 = 1$ is the initial AUM. Variable $h_{i,t}$ is the holdings of instrument i at time t , $\delta_{i,t}$ is change of market value between $t-1$ and t for instrument i . Note profits or losses are being reinvested whenever $t \in B$, hence preventing negative prices. Dividends $d_{i,t}$ are already embedded in K_t .

The purpose of $\omega_{i,t} \left(\sum_{i=1}^I |\omega_{i,t}| \right)^{-1}$ is to de-lever the allocations.

Let τ_i be transaction cost associated with trading \$1 of the instrument. Three additional variables that the strategy needs to know for every observed bar t are:

- i. Rebalance Costs: variable cost $\{c_t\}$ associated with allocation rebalance is

$$c_t = \sum_{i=1}^I (|h_{i,t-1}| p_{i,t} + |h_{i,t}| o_{i,t+1}) \varphi_{i,t} \tau_i \quad \forall t \in B$$

Note c_t is not embedded in K_t , as shorting will generate fictitious proceeds when allocation is rebalanced. In code, $\{c_t\}$ is treated as a (negative) dividend.

- ii. Bid-Ask Spread: the cost $\{\tilde{c}_t\}$ of buying or selling one unit of this ETF,

$$ildec_t = \sum_{i=1}^I |h_{i,t-1}| p_{i,t} \varphi_{i,t} \tau_i$$

When a unit is bought or sold, strategy must charge this cost \tilde{c}_t .

- iii. Volume: volume traded $\{v_t\}$ is determined by least active member in the basket. Let $v_{i,t}$ be volume traded by instrument i over bar t . The number of tradable basket units is

$$v_t = \min_i \left\{ \frac{v_{i,t}}{|h_{i,t-1}|} \right\}$$

Transaction costs functions may not be linear, and can be simulated by the strategy.

Method 1.4.14. ETF Trick: Computation of Allocation Vector with PCA

Consider an IID multivariate Gaussian process with means vector μ of size $N \times 1$, and covariance matrix V of size $N \times N$. First, perform spectral decomposition $VW = W\Lambda$, where columns in W are reordered so that elements of Λ diagonal are sorted in descending order. Second, given allocations vector ω , portfolio risk is

$$\sigma^2 = \omega' V \omega = \omega' W \Lambda W' \omega = \beta' \Lambda \beta = (\Lambda^{1/2} \beta)' (\Lambda^{1/2} \beta)'$$

where β is projection of ω on orthogonal basis. Third, Λ is a diagonal matrix, thus

$$\sigma^2 = \sum_{n=1}^N \beta_n^2 \Lambda_{n,n}$$

The risk attributed to the n th component is

$$R_n = \beta_n^2 \Lambda_{n,n} \sigma^{-2} = [W' n]_n^2 \Lambda_{n,n} \sigma^{-2}$$

with $R' 1_N = 1$, and 1_N is a vector of N ones.

Note $\{R_n\}_{n=1,\dots,N}$ is distribution of risks across orthogonal components.

Next, compute vector ω which delivers user-defined risk distribution R . Note from earlier,

$$\beta = \left\{ \sigma \sqrt{\frac{R_n}{\Lambda_{n,n}}} \right\}_{n=1,\dots,N}$$

which represents allocation in new (orthogonal basis).

The allocation in old basis is $\omega = W\beta$. Rescaling ω re-scales σ , hence keeping risk distribution constant.

Method 1.4.15. ETF Trick: Single Futures Roll

To work with non-negative rolled series, derive price series of \$1 investment as follows:

- i. Compute time series of rolled futures prices
- ii. Compute return r as rolled price change divided by previous roll price
- iii. Form a price series using these returns

These methods allow us to produce a continuous, homogeneous, and structured dataset from collection of unstructured financial data. Note however, that several ML algorithms do not scale well with sample size. ML algorithms achieve higher accuracy when they attempt to learn from relevant examples.

Method 1.4.16. *Sampling for Reduction*

To reduce the amount of data used to fit ML algorithm, downsampling could be used.

- i. Sequential sampling at constant step size (linspace sampling)
- ii. Sampling randomly using uniform distribution (uniform sampling)

Note both samples do not necessarily contain subset of most relevant observations.

Method 1.4.17. *Event-Based Sampling: CUSUM Filter*

Bets are often placed after some event takes place, hence to let ML algorithm learn whether there is an accurate prediction function under these circumstances, CUSUM filter could be used.

This is a quality-control method, to detect shift in mean value of measured quantity away from a target value. Let $\{y_t\}_{t=1,\dots,T}$ be IID observations arising from a locally stationary process. The cumulative sums are

$$S_t = \max\{0, S_{t-1} + y_t - E_{t-1}[y_t]\}, \quad S_0 = 0$$

An action will be recommended at the first t satisfying $S_t \geq h$ for some threshold h (filter size).

Note $S_t = 0$ whenever $y_t = E_{t-1}[y_t] - S_{t-1}$, The zero floor means some downward deviations will be skipped.

The filter is set up to identify a sequence of upside divergences from any reset level zero.

The threshold is activated when

$$S_t \geq h \Leftrightarrow \exists \tau \in [1, t] \mid \sum_{i=\tau}^t (y_i - E_{i-1}[y_t]) \geq h$$

This concept of run-ups can be extended to include run-downs, giving symmetric CUSUM filter.

$$\begin{aligned} S_t^+ &= \max\{0, S_{t-1}^+ + y_t - E_{t-1}[y_t]\}, \quad S_0^+ = 0 \\ S_t^- &= \min\{0, S_{t-1}^- + y_t - E_{t-1}[y_t]\}, \quad S_0^- = 0 \\ S_t &= \max\{S_t^+, -S_t^-\} \end{aligned}$$

1.4.3 Data Labelling Techniques**Method 1.4.18.** *Labelling with Fixed-Time Horizon Method*

Given features matrix X with I rows, $\{X_i\}_{i=1,\dots,I}$ drawn from some bars with index $t = 1, \dots, T$, where $I \leq T$, let an observation X_i be assigned a label $y_i \in \{-1, 0, 1\}$,

$$y_i = \begin{cases} -1 & \text{if } r_{t_{i,0}, t_{i,0}+h} < -\tau \\ 0 & \text{if } |r_{t_{i,0}, t_{i,0}+h}| \leq \tau \\ 1 & \text{if } r_{t_{i,0}, t_{i,0}+h} > \tau \end{cases}$$

$$r_{t_{i,0}, t_{i,0}+h} = \frac{p_{t_{i,0}+h}}{p_{t_{i,0}}} - 1$$

where τ is a pre-defined constant threshold, $t_{i,0}$ is index of bar immediately after X_i takes place, $t_{i,0} + h$ is index of h -th bar after $t_{i,0}$, and $r_{t_{i,0}, t_{i,0}+h}$ is price return over bar horizon h .

Remark 1.4.19. *Limitations of Fixed-Time Horizon Method*

- i. Time bars do not exhibit good statistical properties (as seen earlier)
- ii. The same threshold τ is applied regardless of observed volatility.
Compute daily volatility at intraday estimation points, applying span of n days to an exponentially weighted moving standard deviation.

Method 1.4.20. *Labelling with Triple-Barrier Method*

Labels an observation according to first barrier touched out of three barriers.

- i. Set two horizontal barriers and one vertical barrier. Horizontal barriers are defined by profit-taking and stop-loss limits, which are a dynamic function of estimated volatility (realised or implied). Third barrier is the number of bars elapsed since the position was taken (expiration limit).
- ii. If upper barrier is touched first, label observation as 1. If lower barrier is touched first, label observation as -1 . If vertical barrier is touched first, either label by sign of the return or with 0.

Note that the method is path-dependent. To label an observation, need to account for entire path spanning $[t_{i,0}, t_{i,0} + h]$ where h defines the vertical barrier (expiration limit). Let $t_{i,1}$ be the time of first barrier touch with return as $r_{t_{i,0}, t_{i,1}}$. The horizontal barriers may not be symmetric.

Remark 1.4.21. *Triple-Barrier Method Configurations*

Denote a barrier configuration by triplet $[pt, sl, t1]$ which are the upper barrier, lower barrier, physical barrier. Set value as 0 if barrier is inactive, and 1 if barrier is active.

The three useful configurations are:

- i. $[1, 1, 1]$: to realise profit, but have set a maximum tolerance for losses and a holding period.
- ii. $[0, 1, 1]$: to exit after a number of bars, unless stopped-out.
- iii. $[1, 1, 0]$: take profit as long as not stopped-out.

The three less realistic configurations are:

- i. $[0, 0, 1]$: equivalent to fixed-time horizon method.
- ii. $[1, 0, 1]$: position held until a profit is made or maximum holding period is exceeded, without regard for immediate unrealised losses
- iii. $[1, 0, 0]$: position is held until a profit is made. Could lock in loose position for years.

The two illogical configurations are:

- i. $[0, 1, 0]$: aimless. Hold position until stopped-out.
- ii. $[0, 0, 0]$: no barriers. Position locked forever, no label generated.

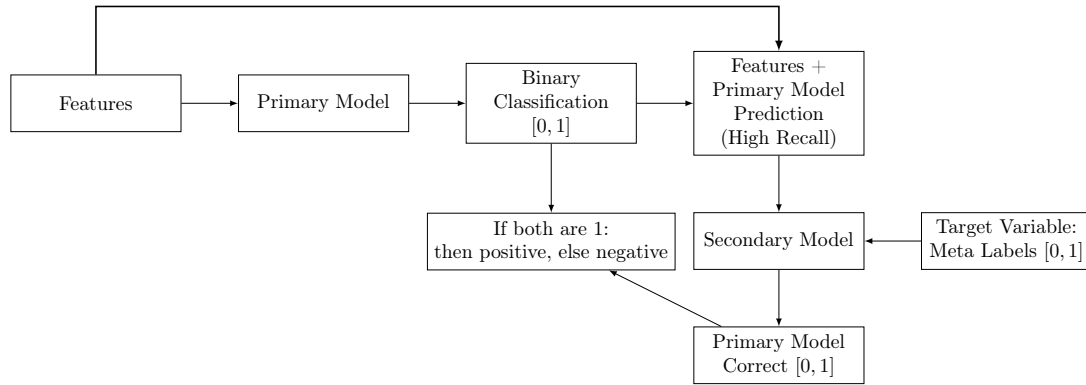


Figure 3: Meta-Labeling Process

Method 1.4.22. *Meta-Labeling*

The technique is particularly helpful to achieve higher F1-scores.

First, build a model that achieves high recall, even if precision is not particularly high. Second, correct for low precision by applying meta-labelling to positives predicted by primary model.

Meta-labelling will filter out false positives, where majority of positives have been identified by primary model.

The second model's purpose is to determine if the positive from primary model is true or false.

- i. Train a primary model (binary classification)
- ii. A threshold level is determined at which the primary model has a high recall, ROC curves could be used to help determine a good level.
- iii. Typical features of second model are as follows:
 - i. Primary model features concatenated with predictions from first model.
 - ii. Market state
 - iii. Features indicative of false positives
 - iv. Distribution related
 - v. Recent model performance

Meta Labels are used as target variable in second model. Fit the second model

- iv. Prediction from the secondary model is combined with the prediction from the primary model and only where both are true, is your final prediction true.

Remark 1.4.23. *Limitations of Meta-Labeling*

- i. If model has overfit the data, meta-labelling will not add much value

- ii. If every trade is not treated as an independent observation, the meta-model is forced to determine day-to-day exposures, which is the wrong way to apply the technique
- iii. Technique trades recall for precision. Require a large number of trades to train on, while being happy with reduction in trade frequency

1.4.4 Data Sample Weights

Note that most of ML literature is based on IID assumption, and ML applications usually fail in finance as these assumptions are unrealistic in the case of financial time series.

Remark 1.4.24. *Overlapping Outcomes*

Let label y_i be assigned to an observed feature X_i , where $y_i = f([t_{i,0}, t_{i,1}])$ is a function over the interval. When $t_{i,1} > t_{j,0}$ and $i < j$, then y_j will depend on common return $r_{t_{j,0}, \min\{t_{i,1}, t_{j,1}\}}$ (over interval $[t_{j,0}, \min\{t_{i,1}, t_{j,1}\})$). The series of labels $\{y_i\}_{i=1, \dots, J}$ are not IID whenever there is overlap between any two consecutive outcomes, i.e., $\exists i \mid t_{i,1} > t_{i+1,0}$. If this is resolved by restricting bet horizon to $t_{i,1} \leq t_{i+1,0}$, there is no overlap, but this will lead to coarse models where features sampling frequency is limited by horizon used to determine outcome. To investigate outcomes that lasted a different duration, samples have to be resampled with different frequency. In addition, if path-dependent labelling technique is to be applied, the sampling frequency will be subordinated to first barrier's touch. Hence, to use $t_{i,1} > t_{i+1,0}$, leading to overlapping outcomes.

Method 1.4.25. *Estimating Uniqueness of Label*

Let two labels y_i and y_j be concurrent at time t , both a function of at least one common return $r_{t-1,t} = \frac{p_t}{p_{t-1}} = 1$. To compute the number of labels that are a function of given return $r_{t-1,t}$:

- i. For each $t = 1, \dots, T$, form a binary array $\{1_{t,i}\}_{i=1, \dots, I}$ where $1_{t,i} \in \{0, 1\}$.
Variable $1_{t,i} = 1$ if and only if $[t_{i,0}, t_{i,1}]$ overlaps with $[t-1, t]$ and $1_{t,i} = 0$ otherwise.
- ii. Compute the number of labels concurrent at t , $c_t = \sum_{i=1}^I 1_{t,i}$

Method 1.4.26. *Average Uniqueness of Label*

To estimate label's uniqueness (non-overlap) across its lifespan.

- i. Uniqueness of label i at time t is $u_{t,i} = 1_{t,i} c_t^{-1}$.
- ii. Average uniqueness of label i is average $u_{t,i}$ over label's lifespan, $\bar{u}_i = (\sum_{t=1}^T u_{t,i})(\sum_{t=1}^T 1_{t,i})^{-1}$.

Note that $\{\bar{u}_i\}_{i=1, \dots, I}$ are not used for forecasting the label, hence there is no information leakage.

Remark 1.4.27. *IID and Oversampling*

Probability of not selecting item i after I draws with replacement on set of I items is $(1 - I^{-1})^I$. As $I \rightarrow \infty$, note that $(1 - I^{-1})^I \rightarrow e^{-1}$. Number of unique observations drawn to be expected is $(1 - e^{-1}) \approx \frac{2}{3}$. If maximum number of overlapping outcomes is $K \leq I$, probability of not selecting a particular item i after I draws with replacement on set of I items is $(1 - K^{-1})^I$. As sample size increase, probability can be approximated as $(1 - I^{-1})^{I^K} \approx e^{-\frac{K}{I}}$. Implication is that incorrectly assuming IID draws lead to oversampling.

Method 1.4.28. *Sampling with Bootstrap, Redundancy*

Sampling with bootstrapping on observations where $I^{-1} \sum_{i=1}^I \bar{u}_i \ll 1$, in-bag observations will increasingly be redundant to each other, and very similar to out-of-bag observations. Two solutions may be:

- i. Drop overlapping outcomes before performing bootstrap.
As overlaps are not perfect, dropping an observation due to overlap will lead to extreme loss in information.
- ii. Utilise the average uniqueness $I^{-1} \sum_{i=1}^I \bar{u}_i$ to reduce undue influence of outcomes that contain redundant information. Ensure in-bag observations are not sampled at frequency much higher than uniqueness.

Method 1.4.29. *Sequential Bootstrap*

Draws made according to changing probability that controls for redundancy.

- i. Observation X_i is drawn from uniform distribution, $i \sim U[1, I]$.
Probability of drawing any value i is $\delta_i^{(1)} = I^{-1}$.
- ii. Second draw, to reduce probability of drawing observation X_j with highly overlapping outcome.
Let φ be sequence of draws (may include repetitions), where $\{\varphi^{(1)}\} = \{i\}$.
Uniqueness of j at time t is $u_{t,j}^{(2)} = 1_{t,j}(1 + \sum_{k \in \varphi^{(1)}} 1_{t,k})^{-1}$, which is the uniqueness from adding alternative j 's to existing sequence of draws $\varphi^{(1)}$.

Average uniqueness of j is average $u_{t,j}^{(2)}$ over j 's lifespan, $\bar{u}_j^{(2)} = (\sum_{t=1}^T u_{t,j}) (\sum_{t=1}^T 1_{t,j})^{-1}$.

A second draw can be made based on updated probabilities $\{\delta_j^{(2)}\}_{j=1,\dots,I}$:

$$\delta_j^{(2)} = \bar{u}_j^{(2)} \left(\sum_{k=1}^I \bar{u}_k^{(2)} \right)^{-1}$$

where $\sum_{j=1}^I \delta_j^{(2)} = 1$. Do a second draw, update $\varphi^{(2)}$, and re-evaluate $\{\delta_j^{(3)}\}_{j=1,\dots,I}$.

iii. Process is repeated until I draws have taken place.

Process draws samples much close to IID, verified by increase in $I^{-1} \sum_{i=1}^I \bar{u}_i$.

Method 1.4.30. *Weighting Observations by Uniqueness and Absolute Return*

Let labels be a function for return sign ($\{-1, 1\}$ for standard label, $\{0, 1\}$ for meta-label). The sample weights can be defined in terms of sum of attributed returns over event's life-span, $[t_{i,0}, t_{i,1}]$,

$$\tilde{w}_i = \left| \sum_{t=t_{i,0}}^{t_{i,1}} \frac{r_{t-1,t}}{c_t} \right|, \quad w_t = \tilde{w}_i \left(\sum_{j=1}^I \tilde{w}_j \right)^{-1}$$

where $\sum_{i=1}^I w_i = I$. The method weigh an observation as a function of absolute log returns that can be attributed uniquely to it. Lower returns should be assigned higher weights.

Method 1.4.31. *Time Decay Weighting*

To let sample weights decay as new observations arrive.

Let $d[x] \geq 0 \forall x \in [0, \sum_{i=0}^I \bar{u}_i]$ be time-decay factors multiplying sample weights from earlier.

The final weight has no decay, $d[\sum_{i=1}^I \bar{u}_i] = 1$, and all other weights will adjust relative to that.

Let $c \in (-1, 1]$ be user-defined parameters that determines decay function as follows:

- i. If $c \in [0, 1]$, then $d[1] = c$ with linear decay
- ii. If $c \in (-1, 0)$, then $d[-c \sum_{i=1}^I \bar{u}_i] = 0$, with linear decay between $[-c \sum_{i=1}^I \bar{u}_i, \sum_{i=1}^I \bar{u}_i]$, and $d[x] \forall x \leq -c \sum_{i=1}^I \bar{u}_i$.

If given linear piecewise function $d = \max\{0, a + bx\}$, requirements are met by following boundary conditions:

- i. $d = a + b \sum_{i=1}^I \bar{u}_i = 1 \Rightarrow a = 1 - b \sum_{i=1}^I \bar{u}_i$
- ii. Contingent on c :
 1. $d = a + b \cdot 0 = c \Rightarrow b = (1 - c) (\sum_{i=1}^I \bar{u}_i)^{-1} \quad \forall c \in [0, 1]$
 2. $d = a - bc \sum_{i=1}^I \bar{u}_i = 0 \Rightarrow b = [(c + 1) \sum_{i=1}^I \bar{u}_i]^{-1} \quad \forall c \in (-1, 0)$

In the implementation, decay takes place according to cumulative uniqueness. Note that

- i. $c = 1$ means there is no time decay
- ii. $0 < c < 1$ means weights decay linearly over time, but every observation still receives strictly positive weight, regardless of age
- iii. $c = 0$ means weights converge linearly to zero over time
- iv. $c < 0$ means oldest portion cT of observations receive zero weight (erased from memory)

Method 1.4.32. *Class Weighting*

Weights for underrepresented labels. Critical in classification problems where the most important classes have rare occurrences. To assign higher weights to samples associated with those rare labels.

1.4.5 Fractionally Differentiated Features

Standard stationarity transformations (i.e. integer differentiation) reduce signal by removing memory. Although stationarity is necessary for inferential purposes, it is rarely the case that we want all memory to be erased.

Fractionally differentiated processes exhibit long-term persistence and anti-persistence, hence enhancing the forecasting power compared to standard ARIMA approach.

Definition 1.4.33. *BackShift Operator*

Let B be the backshift operator applied to a matrix of real-valued features $\{X_t\}$, where $B^k X_t = X_{t-k}$ for any integer $k \geq 0$. By binomial expansion, we then have

$$\begin{aligned} (1 - B)^d &= \sum_{k=0}^{\infty} \binom{d}{k} (-B)^k = \sum_{k=0}^{\infty} \prod_{i=0}^{k-1} (d-i) \frac{(-B)^k}{k!} = \sum_{k=0}^{\infty} (-B)^k \prod_{i=0}^{k-1} \frac{d-i}{k-i} \\ &= 1 - dB + \frac{d(d-1)}{2!} B^2 - \frac{d(d-1)(d-2)}{3!} B^3 + \dots \end{aligned}$$

Remark 1.4.34. *Properties of Fractionally Differentiated Features*

Let d be a real (non-integer) positive number. The arithmetic series consists of dot product

$$\begin{aligned} \tilde{X}_t &= \sum_{k=0}^{\infty} \omega_k X_{t-k} \\ \omega &= \left\{ 1, -d, \frac{d(d-1)}{2!}, -\frac{d(d-1)(d-2)}{3!}, \dots, (-1)^k \prod_{i=0}^{k-1} \frac{d-i}{k-i}, \dots \right\} \\ X &= \{X_t, X_{t-1}, \dots, X_{t-k}, \dots\} \end{aligned}$$

where ω are the weights, X are the values. Properties of these features are:

- i. Long memory: if d is a positive integer number, then

$$\prod_{i=0}^{k-1} \frac{d-i}{k-i} = 0 \quad \forall k > d$$

and memory beyond that point is cancelled.

- ii. Iterative weight generation: given sequence of weights ω , for $k = 0, \dots, \infty$, the weights are

$$\omega_k = -\omega_{k-1} \frac{d-k+1}{k}, \quad \omega_0 = 1$$

- iii. Convergence: For $k > d$, if $\omega_{k-1} \neq 0$, then

$$\left| \frac{\omega_k}{\omega_{k-1}} \right| = \left| \frac{d-k+1}{k} \right| < 1$$

and $\omega_k = 0$ otherwise. Hence weights converge asymptotically to zero.

For positive d and $k < d+1$, then $\frac{d-k+1}{k} \geq 0$, which makes initial weights alternate in sign.

For non-integer d , once $k \geq d+1$, ω_k will be negative if $\text{int}[d]$ is even, and positive otherwise.

In summary, $\lim_{k \rightarrow \infty} \omega_k = 0^-$ when $\text{int}[d]$ is even, and $\lim_{k \rightarrow \infty} \omega_k = 0^+$ when $\text{int}[d]$ is odd.

In special case $d \in (0, 1)$, that $-1 < \omega_k < 0 \forall k > 0$. Alternate weight signs makes $\{\tilde{X}_t\}_{t=1, \dots, T}$ stationary, as memory wanes or is offset over the long run.

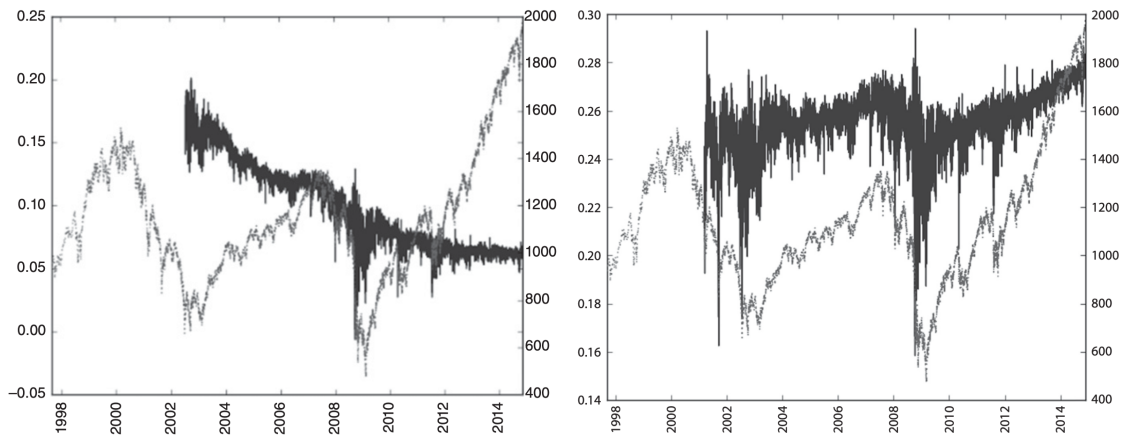


Figure 4: Fractional differentiation controlling for weight loss with expanding and fixed-width window

Method 1.4.35. *Expanding Window*

Given time series T with real observations $\{X_t\}_{t=1,\dots,T}$, for each l , the relative weight loss is defined as

$$\lambda_l = \sum_{j=T-l}^T |\omega_j| \bigg/ \sum_{i=0}^{T-1} |\omega_i|$$

Given tolerance level $\tau \in [0, 1]$, determine value l^* such that $\lambda_{l^*} \leq \tau$ and $\lambda_{l^*+1} > \tau$. This value l^* corresponds to the first results $\{\tilde{X}_t\}_{t=1,\dots,l^*}$, where weight-loss is beyond acceptable threshold $\lambda_t > \tau$.

From Remark 1.4.34, it is clear λ_{l^*} depends on convergence speed of $\{\omega_k\}$, which in turn depends on $d \in [0, 1]$. For $d = 1$, $\omega_k = 0 \ \forall k > 1$, and $\lambda_l = 0 \ \forall l > 1$, hence it suffices to drop \tilde{X}_1 .

As $d \rightarrow 0^+$, l^* increases, and larger portion of initial $\{\tilde{X}_t\}_{t=1,\dots,l^*}$ needs to be dropped to keep the weight loss $\lambda_{l^*} < \tau$. Note that there will be negative drift caused by negative weights added to initial observations as window is expanded. By controlling for weight loss, negative drift is still substantial as $\{\tilde{X}_t\}_{t=l^*+1,\dots,T}$ are computed on an expanding window.

Method 1.4.36. *Fixed-Width Window*

Drop weights after their modulus $|\omega_k|$ decreases below a given threshold τ . This is equivalent to finding the first l^* such that $|\omega_{l^*}| \geq \tau$ and $|\omega_{l^*+1}| \leq \tau$, setting a new variable $\tilde{\omega}_k$:

$$\tilde{\omega}_k = \begin{cases} \omega_k & \text{if } k \leq l^* \\ 0 & \text{if } k > l^* \end{cases}, \quad \tilde{X}_t = \sum_{k=0}^{l^*} \tilde{\omega}_k X_{t-k} \quad \text{for } t = T - l^* + 1, \dots, T$$

Note that the same vector of weights is used across all estimates of $\{\tilde{X}_t\}_{t=l^*,\dots,T}$, hence avoiding negative drift caused by expanding window's added weights.

Distribution has skewness and excess kurtosis from memory, but it is stationary.

2 Research

2.1 Research in AI

2.1.1 Reading AI Papers

Remark 2.1.1. *Reading Wide*

Navigate through literature reading small amounts of individual research papers, build and improve mental model of research topic.

On *Papers with Code*, check out top benchmark models, read abstract and take note of key information only. Next step for each model, look into the datasets used by the papers and make notes.

On *Google Scholar*, check if there is any survey papers (reviews and describe state of problem space with challenges and opportunities) to get up to speed, read each paper, and make notes.

Next find related works (recently published) to understand how researchers in the field traditionally approached the problems, and what the emerging trends are. Related works will populate reading list.

Remark 2.1.2. *Reading Deep*

First pass will not understand more than 10% of research paper, and may require reading of another more fundamental paper. Then read subsequent passes until 70% ~ 80% understanding.

In introduction, highlight problems and challenges, solutions to challenges, main contributions of the work. Should be able to extract a problem-solution chain from the introduction and proposed solution.

In methods, maintain list of concepts not yet understood. If there is link to paper references, keep track.

In experiments, highlight data setup, main table of results. Keep track of unfamiliar evaluation metrics.

Remark 2.1.3. *Practical AI Research Tools*

Experiment tracking tools include 'Weights & Biases', 'Tensorboard', 'Neptune'.

If using 'Weights & Biases', model artefacts can be stored directly on a system of choice.

To train deep learning pipelines by using a config which can modify depending on which dataset, model, or configuration is used, Python's Hydra package may be used.

Remark 2.1.4. *Identifying Gaps in Research Paper*

- i. Identify gaps in research question by comparing with the research hypotheses of compiled papers
- ii. Identify gaps in experimental gaps, such as shortcomings in evaluation of methods, the way the comparisons were chosen or implemented, and whether the experimental setup tests the research hypothesis decisively.
- iii. Identify gaps through expressed limitations, implicit and explicit.

2.1.2 Writing AI Papers

Remark 2.1.5. *Generating Ideas for Building on Research Paper*

- i. Change task of interest. Can the main ideas be applied to a different modality, a different data type? Can the method or learned model be applied to a different task? Can the outcome of interest be changed?
- ii. Change the evaluation strategy. Can it be evaluated on a different dataset or different metric? Explore why something works well/breaks. Make different comparisons.
- iii. Change the proposed method. Can training dataset or data elements be changed? Can the pre-training/training strategy be changed? Can the deep learning architecture or problem formulation be changed?

Remark 2.1.6. *Iterating on Research Ideas*

- i. Search for whether idea has been tried. Construct titles for new paper ideas, see if there is Google result.
- ii. Read important related works and follow up works.
- iii. Get feedback from domain experts on drafted ideas in written form. Email to authors of the work being built on, share idea and plan, ask their opinion.

Remark 2.1.7. *Global Structure of ML Papers*

ML Papers follow the following pattern (in 6 to 7 sections):

- i. Abstract (answers 5 to 6 canonical questions in 100 ~ 250 words):
 1. What is the background and gap?
 2. What is the key desideratum?
 3. What is the proposed solution?
 4. What are its main components?
 5. What are its strengths?
 6. What are the notable results (with tasks, numbers)?
- ii. Introduction (start and end position are nearly identical across papers):
 1. Context, success of prior approaches
 2. Weaknesses or gaps of prior methods
 3. Desiderata for an improved solution
 4. Proposed method: key components, contributions
 5. High-level experimental overview, positive results
- iii. Related Work/Background (typically 2 to 3 subsections), each paragraph conveys:
 1. High-level mapping of approach categories
 2. Evolution of methods over time
 3. How the proposed method compares to each category
 4. What gaps persist
- iv. Methods (content always includes):
 1. Overall approach description
 2. Architecture, input/output flow
 3. Loss functions, training objectives
 4. Implementation details
 5. (Sometimes) dataset descriptions or task-specific usage
- v. Experiments (how results are conveyed. Always structured as such):
 1. Overall evaluation setup
 2. Dataset description (if not earlier)
 3. Implementation details
 4. Results per task type (tables/figures)
 5. Ablations at the end
 6. References to earlier figures and comparisons to prior models

- vi. Conclusion, Broader Impacts (mirrors abstract but expands):
 - 1. Solution, components
 - 2. Strengths, notable results
 - 3. Future directions
 - 4. Limitations and societal considerations
 - 5. Motivations for follow-up work

Remark 2.1.8. *Figure Progression of ML Papers*

ML Papers follow the following pattern for figures:

- i. Method overview diagram (always first)
- ii. Lower-level architecture or objective illustration
- iii. Comparisons vs. prior models across multiple tasks
- iv. Ablation studies
- v. (Optional) qualitative examples, predictions, dataset samples

References

- Narang, R. K. (2013). *Inside the Black Box: The Simple Truth About Quantitative Trading*. Wiley Finance Series. Wiley.
- Prado, M. L. D. (2018). *Advances in Financial Machine Learning*. Wiley Finance Series. Wiley.
- Thierry, A. and G. Helyette (2000, October). Order flow, transaction clock, and normality of returns. *The Journal of Finance* 55(5), 2259–2284.
- Velu, R. (2020). *Algorithmic Trading and Quantitative Strategies*. CRC Press.