

# Portfolio Optimization: Theory, Methods, and LLM Integration

## A Systematic Framework for Quantitative Portfolio Construction

2025

### Abstract

This guide presents a systematic framework for constructing optimized equity portfolios by combining modern quantitative methods with large language model (LLM) augmentation. The theoretical foundations: moment estimation, regime-switching models, Bayesian view integration, risk budgeting, hierarchical allocation, and robust optimization are developed as a self-contained treatment of each method's mathematical basis and practical considerations. Regime-switching estimation receives dedicated treatment through Hidden Markov Models and Deep Markov Models, which provide time-varying moment estimates that adapt to the prevailing market environment and inform dynamic risk measure selection. Each chapter maps mathematical formulations to the algorithmic procedures that realize them, establishing a clear path from theory to implementation. LLM integration is treated as a first-class component throughout: language models contribute to universe screening, regime classification, view generation, risk measure selection, constraint specification, and backtest interpretation. The result is a unified pipeline architecture where prior construction, optimization, and validation compose as interchangeable estimation stages amenable to cross-validation and hyperparameter tuning.

### Contents

<b>1</b>	<b>Data Preparation and Universe Construction</b>	<b>3</b>
1.1	Return Computation and Data Representation . . . . .	3
1.2	Outlier Treatment and Data Cleaning . . . . .	4
1.3	Pre-Selection Pipeline . . . . .	6
<b>2</b>	<b>Moment Estimation and Prior Construction</b>	<b>9</b>
2.1	Expected Return Estimation . . . . .	9
2.2	Covariance Matrix Estimation . . . . .	12
2.3	Factor Model Construction . . . . .	17
2.4	Regime-Switching Models . . . . .	19
2.5	The Empirical Prior . . . . .	23
2.6	LLM-Augmented Moment Estimation . . . . .	24
<b>3</b>	<b>View Integration and Bayesian Updating</b>	<b>26</b>
3.1	The Black-Litterman Framework . . . . .	26
3.2	Entropy Pooling: Non-Linear View Integration . . . . .	30
3.3	Opinion Pooling: Combining Multiple Expert Views . . . . .	32
3.4	LLM-Driven View Generation and Integration . . . . .	33

<b>4</b>	<b>Risk Measures, Diversification, and Hierarchical Methods</b>	<b>37</b>
4.1	Risk Measures: From Variance to Tail Risk . . . . .	37
4.2	Risk Budgeting and Equal Risk Contribution . . . . .	42
4.3	Maximum Diversification . . . . .	44
4.4	Distance Measures and Codependence . . . . .	44
4.5	Hierarchical Clustering . . . . .	46
4.6	Hierarchical Risk Parity . . . . .	46
4.7	Hierarchical Equal Risk Contribution . . . . .	47
4.8	Nested Clusters Optimization . . . . .	48
4.9	Regime-Driven Risk Adaptation . . . . .	48
4.10	LLM-Augmented Risk and Diversification . . . . .	50
<b>5</b>	<b>Portfolio Optimization and Robust Methods</b>	<b>52</b>
5.1	Objective Functions . . . . .	52
5.2	Constraints . . . . .	53
5.3	Robust Optimization Under Parameter Uncertainty . . . . .	55
5.4	Distributionally Robust CVaR . . . . .	56
5.5	Synthetic Data and Stress Testing . . . . .	56
5.6	Benchmark Tracking . . . . .	57
5.7	Naive Allocation Methods . . . . .	58
5.8	Ensemble Optimization . . . . .	59
<b>6</b>	<b>Validation, Model Selection, and Production Pipeline</b>	<b>60</b>
6.1	Walk-Forward Backtesting . . . . .	60
6.2	Combinatorial Purged Cross-Validation . . . . .	61
6.3	Multiple Randomized Cross-Validation . . . . .	62
6.4	Performance Scoring . . . . .	62
6.5	Hyperparameter Tuning . . . . .	64
6.6	Pipeline Architecture . . . . .	65
6.7	Rebalancing Frameworks . . . . .	66

# 1 Data Preparation and Universe Construction

Robust portfolio optimization depends on the quality of its inputs. This chapter establishes the data preparation pipeline that transforms raw price data into a well-conditioned universe suitable for optimization.

## 1.1 Return Computation and Data Representation

The fundamental input to any portfolio optimization is the matrix of asset returns. Two conventions exist, arithmetic (linear) returns and logarithmic returns, and the choice between them has material consequences for portfolio construction.

**Arithmetic returns** express the fractional change in price over a single period:

$$r_{i,t} = \frac{P_{i,t}}{P_{i,t-1}} - 1$$

where  $P_{i,t}$  denotes the adjusted closing price of asset  $i$  at time  $t$ . The critical property of arithmetic returns for portfolio optimization is **cross-sectional additivity**: the portfolio return at time  $t$  decomposes exactly as a weighted sum of constituent returns,

$$r_{P,t} = \sum_{i=1}^N w_i \cdot r_{i,t}$$

where  $w_i$  is the weight assigned to asset  $i$ . This identity holds exactly and enables the familiar representation of portfolio expected return as  $\mathbb{E}[r_P] = \mathbf{w}^\top \boldsymbol{\mu}$  and portfolio variance as  $\sigma_P^2 = \mathbf{w}^\top \boldsymbol{\Sigma} \mathbf{w}$ , on which the entire mean-variance apparatus depends.

**Logarithmic returns** are defined as:

$$r_{i,t}^{\log} = \ln\left(\frac{P_{i,t}}{P_{i,t-1}}\right)$$

Log returns aggregate across time by simple addition, a consequence of the telescoping property of logarithms:

$$\ln\left(\frac{P_T}{P_0}\right) = \sum_{t=1}^T \ln\left(\frac{P_t}{P_{t-1}}\right)$$

This temporal additivity makes log returns convenient for multi-period compounding analysis and for statistical modeling under the assumption of normally distributed innovations. However, log returns do **not** aggregate across assets:  $\ln(1+r_P) \neq \sum_i w_i \ln(1+r_i)$ . This violation of cross-sectional additivity renders log returns unsuitable as direct inputs to portfolio optimization. Using log returns in a mean-variance optimizer produces incorrect portfolio returns and misleading risk estimates.

**Multi-period scaling** of return moments introduces a further subtlety. The naive approach of scaling mean returns linearly with the horizon ( $\mu_T = T\mu$ ) and volatility by the square root of time ( $\sigma_T = \sigma\sqrt{T}$ ) relies on the assumption of independent and identically distributed returns. While this approximation is serviceable for short horizons, it introduces systematic error for horizons exceeding one year, where compounding effects, mean reversion, and volatility clustering become material.

The correct approach assumes **log-normal price dynamics**, under which the multi-period gross return follows from the continuously compounded return distribution. Specifically, if single-period log returns are normally distributed with mean  $\mu$  and variance  $\sigma^2$ , then the expected cumulative arithmetic return over horizon  $T$  is:

$$\mathbb{E}[R_T] = \exp\left(\mu T + \frac{1}{2}\sigma^2 T\right) - 1$$

and the variance of the cumulative arithmetic return is:

$$\text{Var}[R_T] = \exp(2\mu T + \sigma^2 T) \left( \exp(\sigma^2 T) - 1 \right)$$

These expressions capture the convexity adjustment arising from compounding and correctly propagate uncertainty over extended horizons.

In practice, adjusted price series are first converted to arithmetic returns for use in downstream estimation. Multi-period adjustment is handled by applying the log-normal compounding correction to both expected returns and covariance estimates before they enter the optimizer. Setting the investment horizon to the target number of periods (for example, 252 for annualization with daily data) ensures that the moments fed into the optimizer reflect the compounding dynamics of wealth accumulation rather than the naive linear scaling assumption.

## 1.2 Outlier Treatment and Data Cleaning

Financial return data frequently contain anomalous observations arising from data vendor errors, corporate actions improperly adjusted, exchange glitches, or genuine but extreme market events. Failing to address outliers corrupts moment estimates: a single erroneous return of +500% can dominate a covariance matrix estimated from hundreds of observations. Conversely, indiscriminate removal of genuine extreme returns understates tail risk and biases the portfolio toward insufficient hedging.

A principled **three-group methodology** balances these concerns by applying graduated treatment based on the magnitude of the deviation from the cross-sectional or time-series mean:

1. **Values exceeding 10 standard deviations:** These observations are removed entirely from the estimation sample. Returns of this magnitude almost certainly reflect data errors: erroneous price entries, missing adjustment factors, or transmission failures. Their inclusion would dominate moment estimates regardless of the estimation method employed.

2. **Values between 3 and 10 standard deviations:** These observations are **winsorized** to the  $\pm 3\sigma$  boundary:

$$X_{\text{winsorized}} = \mu + \text{sign}(X - \mu) \times 3\sigma$$

Winsorization preserves the direction and approximate timing of extreme returns while capping their influence on moment estimates. Unlike truncation (which discards the observation entirely), winsorization retains the data point and its ordinal rank, reducing distortion in correlation estimates.

3. **Values within 3 standard deviations:** These observations are retained without adjustment. They represent the core distribution of returns from which moments are estimated.

The thresholds of 3 and 10 standard deviations reflect a balance between robustness and information loss. Under a Gaussian distribution, returns beyond  $3\sigma$  occur with probability 0.27%, roughly once per year for daily data. Returns beyond  $10\sigma$  have vanishing probability under any reasonable model and warrant removal.

**Missing data imputation** arises when assets have sporadic gaps in their return histories (holidays in one market but not another, trading halts, or data vendor outages). Simple approaches such as zero-fill or forward-fill introduce biases: zero-fill deflates volatility estimates, while forward-fill creates artificial serial correlation. More principled methods include imputation from sector or industry averages (preserving cross-sectional relationships) and regression-based prediction from correlated factors (preserving the covariance structure of the broader universe).

**Survivorship bias** constitutes one of the most insidious data quality issues in historical analysis. If the estimation sample includes only assets that survived to the present, historical return and risk estimates are upwardly biased, because the worst-performing assets (those that delisted, defaulted, or were acquired at distressed valuations) are systematically excluded. Proper treatment requires including delisted stocks in the historical sample with correct handling of delisting returns, the final return realized by investors when an asset ceases trading.

**Look-ahead bias** arises when the estimation procedure uses information that was not available to market participants at the time of the investment decision. Accounting data is particularly susceptible: quarterly earnings are reported with lags of several weeks, and point-in-time databases must be used to ensure that only data available as of the estimation date enters the model. Reporting lags must be explicitly accounted for to ensure that all information used in the optimization was genuinely available to market participants at the relevant decision point.

**Data validation** encompasses systematic checks for impossible or inconsistent values: negative prices or volumes, sudden jumps exceeding plausible limits, and inconsistencies across related fields (e.g., a closing price outside the day's high-low range). These checks serve as a first line of defense against data corruption and should be applied before any statistical treatment.

### 1.3 Pre-Selection Pipeline

Before optimization can proceed, the raw universe of assets must be refined through a series of systematic filters that enforce data quality requirements, remove redundancy, and focus the optimization on a well-conditioned set of investable instruments. Each filter addresses a specific pathology that, if left untreated, would degrade optimization quality through numerical instability, spurious diversification, or wasted degrees of freedom on uninvestable assets.

#### 1.3.1 Complete History Selection

The covariance matrix lies at the heart of virtually every portfolio optimization method. Its estimation requires a complete panel of contemporaneous returns; assets with gaps, late-starting histories, or early terminations introduce missing-value artifacts that corrupt the covariance estimate. Pairwise deletion (estimating each covariance element from whatever overlapping observations exist) can produce a covariance matrix that is not positive semi-definite, violating the fundamental requirement for a valid optimization input.

The most conservative and numerically reliable approach retains only assets with complete return histories across the full estimation window. This filter ensures that the covariance matrix is estimated from a balanced panel, guaranteeing positive semi-definiteness when combined with a sufficient number of observations relative to the number of assets.

#### 1.3.2 Zero-Variance Filtering

Assets with negligible return dispersion (suspended stocks, instruments with stale prices that have not been updated, or data errors that produce constant price series) contribute no information to portfolio diversification. Worse, their near-zero variance creates **numerical instability** in covariance matrix inversion: the inverse of a matrix with near-zero eigenvalues amplifies estimation noise, producing extreme and unstable portfolio weights.

Removing zero-variance or near-zero-variance assets before covariance estimation eliminates this source of numerical pathology at negligible cost to portfolio opportunity.

#### 1.3.3 Correlation Filtering

Highly correlated assets provide **redundant diversification**: a portfolio holding multiple assets with pairwise correlations exceeding 0.95 concentrates risk in a single underlying factor despite appearing diversified by asset count. This redundancy inflates the effective dimensionality of the optimization problem without contributing genuine risk reduction, and it exacerbates estimation error in the covariance matrix by introducing near-collinearity.

Systematic correlation filtering identifies all asset pairs whose sample correlation exceeds a specified threshold (typically  $\rho > 0.90$  to 0.95) and drops one asset from each pair, retaining the asset with superior risk-adjusted characteristics (e.g., higher Sharpe ratio, greater liquidity, or lower estimation uncertainty). Typical threshold values range from 0.90 to 0.95, balancing redundancy removal against excessive universe shrinkage.

### 1.3.4 Extreme Performance Selection

When the investable universe is large and factor-scoring or signal-generation has been applied, optimization benefits from focusing on a subset of assets with the strongest factor exposures. Rather than optimizing across hundreds or thousands of assets, many of which contribute negligible expected alpha, the universe is reduced to the  $k$  highest-ranked (or lowest-ranked, depending on the signal direction) assets.

This concentration serves two purposes. First, it reduces the dimensionality of the optimization problem, improving the condition number of the covariance matrix and the stability of the resulting portfolio. Second, it ensures that optimization effort is allocated to assets where the investment thesis is strongest, rather than diluted across marginal positions.

### 1.3.5 Pareto-Optimal Selection

A complementary approach to universe reduction identifies assets that are **Pareto-dominated** in the mean-variance sense. Asset  $i$  is Pareto-dominated if there exists another asset  $j$  such that:

$$\mu_j \geq \mu_i \quad \text{and} \quad \sigma_j \leq \sigma_i$$

with at least one inequality strict. In words, asset  $j$  offers equal or higher expected return with equal or lower risk. A dominated asset  $i$  is strictly inferior: no rational mean-variance investor would prefer it to  $j$  regardless of risk aversion.

Eliminating Pareto-dominated assets removes strictly inferior alternatives from the feasible set without discarding any asset that could contribute to the efficient frontier. The remaining **non-dominated** assets form the Pareto front in the mean-variance plane, and only these need enter the optimizer.

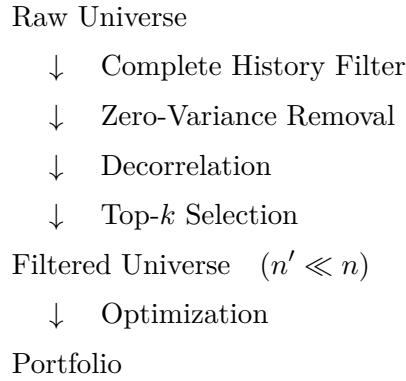
### 1.3.6 Expiring Asset Filtering

In mixed portfolios that include derivatives, fixed-income instruments, or futures contracts alongside equities, some assets may approach their expiration or maturity date within the investment horizon. Allocating to such assets creates operational complications: the optimizer may assign weight to an instrument that will cease trading before the next rebalancing date, forcing unplanned liquidation and incurring transaction costs.

Filtering out assets whose remaining life is shorter than the investment horizon prevents this mismatch, ensuring that all instruments in the optimization universe remain tradeable throughout the intended holding period.

### 1.3.7 Sequential Pipeline Composition

The filters described above compose sequentially, each operating on the output of the previous stage. The complete pre-selection pipeline follows the chain:



The ordering matters: complete history selection must precede covariance-dependent filters (since the latter require a well-formed return matrix), and correlation filtering should precede dimensionality reduction (since redundant assets inflate the apparent universe size). Extreme performance selection or Pareto-optimal selection typically comes last among the filters, operating on the cleaned and de-duplicated universe to produce the final set of candidates for optimization.

The **pipeline architecture** treats this entire chain, from raw price data through filtering through optimization, as a single composite estimation procedure. This design brings three important benefits:

1. **Hyperparameter tuning:** The threshold parameters of each filter (correlation bound, number of extremes to retain, minimum history length) become hyperparameters of the composite estimator, searchable via cross-validation alongside the optimizer's own parameters.
2. **Cross-validation consistency:** When the pipeline is evaluated under walk-forward or combinatorial cross-validation, the filters are re-fitted on each training fold. This prevents information leakage that would occur if filter decisions were made once on the full sample and then held fixed during validation.
3. **Reproducibility and modularity:** Each filter is an independent transformation that can be added, removed, or reordered without modifying the rest of the pipeline. New filters (e.g., liquidity screens, ESG exclusions) integrate by insertion at the appropriate position.

Composing these transformers with a final optimization estimator enables the entire chain to be fitted, evaluated, and scored through a uniform interface.



## 2 Moment Estimation and Prior Construction

Portfolio optimization requires estimates of the joint return distribution: at minimum, expected returns and the covariance matrix. The quality of these estimates determines whether the optimizer discovers genuine risk-return tradeoffs or merely exploits estimation noise. Classical mean-variance optimization is notoriously sensitive to its inputs: small perturbations in expected returns can produce wildly different portfolio weights, and covariance matrices estimated from finite samples inherit substantial sampling error that propagates directly into allocation decisions. This chapter presents the principal moment estimation methods, organized by the statistical quantity they target, and shows how they compose into prior distributions that feed the optimization pipeline. The progression moves from expected returns through covariance matrices to factor models, culminating in the assembly of these components into prior objects that serve as the interface between estimation and optimization.

### 2.1 Expected Return Estimation

Expected returns are simultaneously the most influential and the most difficult inputs to estimate. A portfolio's optimal weights depend linearly on expected returns in the unconstrained mean-variance problem, so estimation errors in the mean vector translate directly, and often dramatically, into misallocation. Empirically, the signal-to-noise ratio of expected return estimates is far lower than that of covariance estimates: with monthly data, estimating means to useful precision requires decades of observations, whereas the covariance matrix stabilizes over much shorter horizons. This asymmetry motivates the range of estimators discussed below, which trade off between fidelity to historical data and stability through various forms of regularization.

#### 2.1.1 Empirical Mean

The simplest and most direct estimator of expected returns is the sample average:

$$\hat{\mu}_i = \frac{1}{T} \sum_{t=1}^T r_{i,t}$$

where  $r_{i,t}$  denotes the return of asset  $i$  at time  $t$  and  $T$  is the total number of observations. The sample mean is an unbiased estimator of the true expected return under stationarity, and its variance decreases as  $\sigma_i^2/T$ . However, this convergence is slow relative to the precision demanded by portfolio optimization. Merton (1980) demonstrated that estimation error in expected returns dominates portfolio construction for samples shorter than approximately 25 years of monthly data. The core difficulty is that equity return distributions have high variance relative to their means: annualized Sharpe ratios rarely exceed unity, implying that the signal (the mean) is small relative to the noise (the standard deviation). As a result, the sample mean is an unbiased but extremely noisy estimator, and portfolios constructed from raw sample means tend to exhibit extreme and unstable weights.

### 2.1.2 Shrinkage Toward the Grand Mean

Shrinkage estimation addresses the instability of sample means by pulling individual estimates toward a common central value. The shrinkage estimator takes the form:

$$\hat{\mu}_i^{\text{shrunk}} = (1 - \alpha) \bar{\mu}_i + \alpha \cdot \bar{\mu}_{\text{grand}}$$

where  $\bar{\mu}_i$  is the sample mean of asset  $i$ ,  $\bar{\mu}_{\text{grand}} = \frac{1}{N} \sum_{i=1}^N \bar{\mu}_i$  is the grand mean across all  $N$  assets, and  $\alpha \in [0, 1]$  is the shrinkage intensity. When  $\alpha = 0$ , the estimator reduces to the sample mean; when  $\alpha = 1$ , all assets receive the same expected return, and any mean-variance optimizer will produce a minimum-variance portfolio.

The theoretical justification derives from the James-Stein estimator, which demonstrates that for  $N \geq 3$  Gaussian random variables, the sample mean is inadmissible: there exists a shrinkage estimator that uniformly dominates it in terms of total squared error. The James-Stein result provides the theoretically optimal  $\alpha$  for Gaussian returns, balancing the bias introduced by shrinkage against the variance reduction it achieves. In practice, shrinkage reduces the dispersion of expected return estimates, compressing extreme values toward the cross-sectional average. This compression produces more diversified portfolios that are less sensitive to estimation noise in any individual asset's mean return.

### 2.1.3 Exponentially Weighted Mean

When the return-generating process is non-stationary (as during regime changes, structural breaks, or evolving market microstructure) equal weighting of all historical observations may be inappropriate. The exponentially weighted mean assigns geometrically decaying weights to past observations:

$$\hat{\mu}_i = \sum_{t=1}^T w_t \cdot r_{i,t}, \quad w_t \propto (1 - \alpha)^{T-t}$$

where  $\alpha \in (0, 1)$  is the decay parameter. Higher values of  $\alpha$  concentrate weight on more recent observations, producing estimates that adapt quickly to changes in the return distribution. The effective sample size of the estimator is approximately  $1/\alpha$ , so  $\alpha = 0.06$  corresponds to roughly 17 observations of effective history.

The trade-off is between responsiveness and stability. Aggressive decay (large  $\alpha$ ) enables the estimator to track genuine regime shifts, but it also increases susceptibility to recent anomalies or transient market dislocations. Conservative decay (small  $\alpha$ ) approaches the equal-weighted sample mean and inherits its stability at the cost of slower adaptation.

### 2.1.4 Equilibrium (CAPM) Returns

Rather than estimating expected returns from historical data, the equilibrium approach derives implied returns from the assumption that observed market capitalization weights represent an

optimal portfolio. Starting from the first-order condition of a mean-variance investor holding the market portfolio, the implied excess returns are:

$$\boldsymbol{\Pi} = \delta \boldsymbol{\Sigma} \mathbf{w}_{\text{mkt}}$$

where  $\boldsymbol{\Sigma}$  is the covariance matrix,  $\mathbf{w}_{\text{mkt}}$  is the vector of market capitalization weights, and  $\delta$  is the risk aversion coefficient. The risk aversion parameter is typically estimated from the market portfolio's risk-return characteristics:

$$\delta = \frac{\mathbb{E}[R_{\text{mkt}}] - R_f}{\sigma_{\text{mkt}}^2}$$

where  $\mathbb{E}[R_{\text{mkt}}] - R_f$  is the expected market excess return and  $\sigma_{\text{mkt}}^2$  is the variance of market returns.

This approach has a crucial advantage: it requires no historical return estimation for individual assets. The only inputs are the covariance matrix (which is estimated with greater precision than means) and the risk aversion parameter. The resulting implied returns are internally consistent with market equilibrium and produce well-diversified portfolios when used as inputs to mean-variance optimization. For this reason, equilibrium returns serve as the most stable and theoretically grounded baseline for the Black-Litterman framework, where they function as the prior distribution that is subsequently updated with investor views.

### 2.1.5 Estimator Selection Guidance

The choice of expected return estimator should reflect the available data, the investment horizon, and the role the estimates play in the broader portfolio construction pipeline. The following table summarizes recommended choices under common conditions:

Condition	Recommended Estimator	Rationale
Long history ( $T > 500$ ), stable markets	Sample mean	Sufficient data for reliable estimation
Short history or regime uncertainty	Shrinkage toward grand mean	Reduces extreme estimates toward consensus
Recent regime change suspected	Exponentially weighted mean	Weights recent observations more heavily
Black-Litterman or view-based framework	Equilibrium returns	Stable prior for Bayesian updating
Regime-switching dynamics	HMM-blended expected returns	State-conditional estimates adapt to prevailing regime

In practice, the estimator choice interacts with downstream optimization. Mean-variance optimization amplifies estimation error in expected returns, so estimators with lower variance (shrinkage,

equilibrium) tend to produce superior out-of-sample performance even if they introduce modest bias. Risk-parity and minimum-variance approaches, which do not use expected returns at all, sidestep this problem entirely, a design choice that itself constitutes an implicit statement about the difficulty of mean estimation.

## 2.2 Covariance Matrix Estimation

The covariance matrix governs the risk structure of the portfolio. Although covariance estimates are generally more reliable than mean estimates for a given sample size, high-dimensional portfolios introduce challenges of their own: when the number of assets  $N$  approaches or exceeds the number of observations  $T$ , the sample covariance matrix becomes ill-conditioned or singular, and its eigenvalues exhibit systematic biases. The estimators presented below address these challenges through shrinkage, spectral cleaning, sparsity assumptions, and adaptive weighting.

### 2.2.1 Sample Covariance

The unbiased sample covariance matrix is:

$$\hat{\Sigma} = \frac{1}{T-1} \sum_{t=1}^T (\mathbf{r}_t - \hat{\boldsymbol{\mu}})(\mathbf{r}_t - \hat{\boldsymbol{\mu}})^\top$$

where  $\mathbf{r}_t \in \mathbb{R}^N$  is the vector of asset returns at time  $t$  and  $\hat{\boldsymbol{\mu}}$  is the sample mean vector. This estimator is unbiased and maximum likelihood under Gaussian returns, but its quality degrades rapidly as the ratio  $N/T$  increases. For  $N/T > 0.1$ , the eigenvalue spectrum of  $\hat{\Sigma}$  is distorted: the largest eigenvalues are biased upward and the smallest are biased downward. When  $N > T$ , the matrix is singular and cannot be inverted, making it unusable for standard mean-variance optimization.

Even when  $T > N$ , the effective condition number of  $\hat{\Sigma}$  may be large enough to produce numerically unstable portfolio weights. The inverse covariance matrix, which appears in the optimal portfolio formula  $\mathbf{w}^* = \Sigma^{-1}\boldsymbol{\mu}$ , amplifies small eigenvalues, causing extreme sensitivity to estimation noise in the least-variance eigenvectors.

### 2.2.2 Ledoit-Wolf Shrinkage

The Ledoit-Wolf estimator addresses the eigenvalue distortion of the sample covariance by shrinking it toward a structured target matrix:

$$\hat{\Sigma}_{\text{LW}} = \delta^* \mathbf{F} + (1 - \delta^*) \mathbf{S}$$

where  $\mathbf{S}$  is the sample covariance matrix,  $\mathbf{F}$  is a structured target, and  $\delta^* \in [0, 1]$  is the analytically optimal shrinkage intensity. The optimal  $\delta^*$  minimizes the expected Frobenius norm of the estimation error  $\|\hat{\Sigma}_{\text{LW}} - \Sigma_{\text{true}}\|_F^2$  and admits a closed-form expression that depends only on the sample data.

A common target is the constant-correlation model:

$$\mathbf{F} = \bar{\rho} \mathbf{D} \mathbf{J} \mathbf{D} + (1 - \bar{\rho}) \mathbf{D}^2$$

where  $\mathbf{D} = \text{diag}(\sigma_1, \dots, \sigma_N)$  contains the sample standard deviations,  $\mathbf{J}$  is the matrix of ones, and  $\bar{\rho}$  is the average pairwise sample correlation. This target preserves individual asset volatilities while imposing a common correlation structure, providing a well-conditioned matrix that the sample covariance is shrunk toward.

Empirical studies consistently find that Ledoit-Wolf shrinkage reduces portfolio volatility forecast errors by 30–50% compared to the raw sample covariance, with the largest improvements occurring in high-dimensional settings where  $N/T$  is substantial. The improvement arises because shrinkage corrects the systematic eigenvalue distortion: overestimated large eigenvalues are pulled down, underestimated small eigenvalues are pulled up, and the resulting matrix is better conditioned.

### 2.2.3 Oracle Approximating Shrinkage

The Oracle Approximating Shrinkage (OAS) estimator extends the Ledoit-Wolf framework by computing an analytically optimal shrinkage intensity without requiring the specification of a particular target structure. It approximates the oracle shrinkage (the intensity that would be chosen if the true covariance matrix were known) using only observable quantities from the sample. This estimator adapts automatically to the data characteristics, adjusting its shrinkage intensity based on the dimensionality ratio  $N/T$  and the structure of the sample covariance.

OAS typically performs comparably to Ledoit-Wolf while offering greater flexibility in the choice of structured target.

### 2.2.4 Random Matrix Theory Denoising

Random matrix theory provides a principled framework for separating signal from noise in the eigenvalue spectrum of the sample covariance matrix. Under the null hypothesis that returns are independent with common variance  $\sigma^2$ , the Marchenko-Pastur distribution describes the limiting eigenvalue density of the sample covariance. The upper bound of this distribution is:

$$\lambda_+ = \sigma^2 \left(1 + \sqrt{N/T}\right)^2$$

Eigenvalues of  $\hat{\Sigma}$  that fall below  $\lambda_+$  are indistinguishable from pure noise and carry no exploitable information about the true covariance structure. The denoising procedure replaces these noise eigenvalues with their average while preserving the eigenvalues and eigenvectors above the threshold, which are assumed to reflect genuine systematic risk factors.

Formally, the denoised covariance matrix is reconstructed as:

$$\hat{\Sigma}_{\text{denoised}} = \sum_{k:\lambda_k > \lambda_+} \lambda_k \mathbf{v}_k \mathbf{v}_k^\top + \bar{\lambda}_{\text{noise}} \sum_{k:\lambda_k \leq \lambda_+} \mathbf{v}_k \mathbf{v}_k^\top$$

where  $\bar{\lambda}_{\text{noise}}$  is the average of the noise eigenvalues, ensuring the trace (total variance) is preserved. This approach is particularly effective when  $N/T$  is large, as a greater fraction of eigenvalues falls within the noise band.

### 2.2.5 Detoning

Detoning removes the dominant market factor, corresponding to the largest eigenvalue and its associated eigenvector, from the covariance matrix. In most equity covariance matrices, the first principal component captures the market mode: the tendency of all stocks to move together. While this common factor is a genuine feature of the return distribution, it can obscure the underlying correlation structure that is relevant for diversification.

The detoned covariance matrix is:

$$\hat{\Sigma}_{\text{detoned}} = \hat{\Sigma} - \lambda_1 \mathbf{v}_1 \mathbf{v}_1^\top$$

where  $\lambda_1$  is the largest eigenvalue and  $\mathbf{v}_1$  is the corresponding eigenvector. The resulting matrix reveals the residual correlation structure net of the common market mode. This is useful when the objective is to construct portfolios that are diversified in a factor-neutral sense rather than merely market-directional. Detoning is often applied in combination with denoising: first the noise eigenvalues are cleaned, then the market factor is removed.

### 2.2.6 Gerber Statistic

The Gerber statistic provides a robust measure of co-movement that is less sensitive to the distributional assumptions underlying the Pearson correlation. Rather than measuring linear association across all return magnitudes, the Gerber statistic focuses exclusively on co-movements that exceed a significance threshold, filtering out small random fluctuations that may reflect noise rather than genuine co-dependence.

The Gerber covariance between assets  $i$  and  $j$  is defined through the concordance and discordance counts:

$$G_{ij}(\theta) = \frac{N_{ij}^{++} + N_{ij}^{--} - N_{ij}^{+-} - N_{ij}^{-+}}{N_{ij}^{++} + N_{ij}^{--} + N_{ij}^{+-} + N_{ij}^{-+}}$$

where  $N_{ij}^{++}$  counts the number of periods where both  $|r_{i,t}| > \theta_i$  and  $|r_{j,t}| > \theta_j$  with the same sign, and  $N_{ij}^{+-}$  counts periods where both exceed their thresholds but with opposite signs. The thresholds  $\theta_i$  are typically set as a fraction of each asset's standard deviation. Observations where either asset's return falls within the threshold band are excluded from the count entirely.

This construction makes the Gerber statistic robust to non-Gaussian return distributions, particularly those with heavy tails. By ignoring small returns, it avoids contamination from the many near-zero observations that dominate the sample but carry little information about tail dependence. The resulting covariance matrix tends to be better conditioned and more informative about the co-movement structure during periods of market stress.

### 2.2.7 Exponentially Weighted Covariance

Analogous to the exponentially weighted mean, the exponentially weighted covariance assigns geometrically decaying weights to past observations, producing an estimator that adapts to changing volatility and correlation regimes:

$$\hat{\Sigma}_t = (1 - \alpha) \hat{\Sigma}_{t-1} + \alpha \cdot \mathbf{r}_t \mathbf{r}_t^\top$$

where  $\alpha \in (0, 1)$  controls the decay rate. Small values of  $\alpha$  produce a slowly evolving estimate close to the equal-weighted sample covariance, while large values create a responsive estimate dominated by recent observations. The effective window length is approximately  $1/\alpha$  observations.

This estimator is particularly valuable when volatility and correlation are time-varying, a well-documented empirical regularity. During periods of market stress, correlations tend to increase, and an exponentially weighted estimator captures this dynamic more rapidly than the sample covariance. The trade-off is the same as for the exponentially weighted mean: responsiveness to genuine regime changes comes at the cost of increased susceptibility to transient fluctuations.

### 2.2.8 Graphical Lasso

The Graphical Lasso estimates a sparse precision matrix (inverse covariance matrix) by solving the  $\ell_1$ -penalized maximum likelihood problem:

$$\hat{\Theta} = \arg \max_{\Theta \succ 0} \{\log \det \Theta - \text{tr}(\mathbf{S}\Theta) - \lambda \|\Theta\|_1\}$$

where  $\mathbf{S}$  is the sample covariance,  $\Theta = \Sigma^{-1}$  is the precision matrix, and  $\lambda > 0$  is the regularization parameter. The  $\ell_1$  penalty forces many entries of  $\hat{\Theta}$  to exactly zero, producing a sparse graphical model of conditional dependencies. A zero entry  $\hat{\Theta}_{ij} = 0$  implies that assets  $i$  and  $j$  are conditionally independent given all other assets: their partial correlation is zero.

Sparsity in the precision matrix is a natural assumption for large equity universes: while many pairs of stocks exhibit marginal correlation (through shared exposure to common factors), far fewer pairs have significant conditional dependence after accounting for other assets. The sparse precision matrix directly yields the optimal portfolio weights in the minimum-variance problem, as  $\mathbf{w}_{\text{MV}} \propto \Theta \mathbf{1}$ , making this estimator particularly efficient for large-scale portfolio construction. The regularization parameter  $\lambda$  is typically selected by cross-validation.

### 2.2.9 Implied Covariance

The implied covariance estimator incorporates forward-looking information from options markets by blending implied volatilities with historical correlation estimates:

$$\hat{\Sigma}_{\text{implied}} = f(\sigma_{\text{IV}}, \rho_{\text{historical}})$$

where  $\sigma_{\text{IV}}$  is the vector of implied volatilities extracted from option prices and  $\rho_{\text{historical}}$  is the historical correlation matrix. The construction replaces the diagonal elements of the historical covariance (the variances) with squared implied volatilities while retaining the off-diagonal correlation structure from historical data.

This approach has a theoretical motivation: option-implied volatilities are forward-looking, reflecting the market's consensus expectation of future realized volatility over the option's life. They incorporate information beyond what is available in the historical return series, including anticipated events (earnings announcements, policy decisions) and current risk sentiment. Empirical evidence suggests that implied volatilities are generally superior forecasters of future realized volatility compared to historical estimates.

The practical requirement is access to implied volatility data, which must be routed through the estimation pipeline alongside return data via metadata routing.

### 2.2.10 Estimator Selection Guidance

The following table maps portfolio construction scenarios to recommended covariance estimators:

Condition	Recommended Estimator	Rationale
$N < 100, T > 500$	Sample covariance	Sufficient data for reliable estimation
$N > 100, T < 5N$	Ledoit-Wolf or OAS	Shrinkage corrects overfit eigenvalues
Fat-tailed returns	Gerber statistic	Robust to non-Gaussian co-movement
Known factor structure	Factor model covariance	Exploits structural dimensionality reduction
Large $N > 500$	Graphical Lasso (cross-validated)	Sparse precision matrix
Regime-sensitive strategies	Exponentially weighted covariance	Adapts to changing volatility
Discrete regime shifts	HMM regime-conditional covariance	State-specific estimates blended by filtered probabilities
High noise ratio $N/T > 0.5$	Random matrix theory denoising	Separates signal from noise



Condition	Recommended Estimator	Rationale
Options data available	Implied covariance	Forward-looking volatility

The covariance estimator choice should reflect both the dimensionality challenge and the downstream optimization objective. For minimum-variance portfolios, covariance accuracy is paramount because the optimizer depends entirely on the risk model. For mean-variance or Black-Litterman portfolios, covariance quality matters for risk estimation but interacts with expected return estimates; a well-conditioned covariance matrix prevents the optimizer from generating extreme positions along poorly estimated eigenvectors.

## 2.3 Factor Model Construction

Factor models decompose asset returns into systematic and idiosyncratic components. The canonical factors (value, momentum, quality, and growth) capture distinct sources of systematic return variation that have been extensively documented in the empirical asset pricing literature. Within the portfolio optimization pipeline, factor models serve primarily as a dimensionality reduction device: rather than estimating the full  $N \times N$  covariance matrix directly, the factor model estimates a  $K \times K$  factor covariance matrix plus  $N$  specific variances, drastically reducing the number of parameters and the associated estimation error.

### 2.3.1 Loading Matrix Estimation

The factor model relates asset returns to a set of  $K$  common factors through a linear regression:

$$r_{nt} = \alpha_n + \sum_{k=1}^K \beta_{nk} \cdot f_{kt} + u_{nt}$$

where  $r_{nt}$  is the return of asset  $n$  at time  $t$ ,  $f_{kt}$  is the return of factor  $k$  at time  $t$ ,  $\beta_{nk}$  is the loading of asset  $n$  on factor  $k$ ,  $\alpha_n$  is the intercept, and  $u_{nt}$  is the idiosyncratic residual with  $\mathbb{E}[u_{nt}] = 0$  and  $\text{Cov}(u_{nt}, f_{kt}) = 0$  for all  $k$ .

The loading matrix  $\mathbf{B} \in \mathbb{R}^{N \times K}$  is estimated by ordinary least squares regression:

$$\hat{\mathbf{B}} = (\mathbf{Y}^\top \mathbf{Y})^{-1} \mathbf{Y}^\top \mathbf{X}$$

where  $\mathbf{X} \in \mathbb{R}^{T \times N}$  contains asset returns and  $\mathbf{Y} \in \mathbb{R}^{T \times K}$  contains factor returns. Each column of  $\hat{\mathbf{B}}$  gives the factor exposures of one asset, and these exposures determine how systematic risk propagates from factor space to asset space.

### 2.3.2 Dimensionality Reduction

The central benefit of the factor model is the dramatic reduction in the number of parameters to estimate. The full covariance matrix requires  $\frac{N(N+1)}{2}$  parameters, while the factor model requires only  $\frac{K(K+1)}{2}$  parameters for the factor covariance plus  $N$  specific variances, for a total of  $\frac{K(K+1)}{2} + N$ . The savings are substantial in high-dimensional settings: for  $N = 500$  assets and  $K = 50$  factors, the factor model estimates  $\frac{50 \cdot 51}{2} + 500 = 1,775$  parameters compared to  $\frac{500 \cdot 501}{2} = 125,250$  for the full covariance, a reduction of 98.6%.

The implied covariance structure under the factor model is:

$$\mathbf{\Sigma} = \mathbf{BFB}^\top + \mathbf{D}$$

where  $\mathbf{B} \in \mathbb{R}^{N \times K}$  is the loading matrix,  $\mathbf{F} \in \mathbb{R}^{K \times K}$  is the factor covariance matrix, and  $\mathbf{D} \in \mathbb{R}^{N \times N}$  is a diagonal matrix of specific (idiosyncratic) variances. This decomposition ensures that the estimated covariance is positive semi-definite by construction (provided  $\mathbf{F}$  is positive semi-definite and  $\mathbf{D}$  has non-negative entries), and it is always invertible when the specific variances are strictly positive, regardless of the ratio  $N/T$ .

### 2.3.3 Factor Covariance and Specific Risk

The factor covariance matrix  $\mathbf{F}$  is estimated in the lower-dimensional factor space, where  $K \ll N$  ensures that even the sample covariance is well-conditioned. Any of the covariance estimators discussed in the previous section (Ledoit-Wolf, exponentially weighted covariance, random matrix theory denoising, and others) can be applied to factor returns rather than asset returns, providing additional regularization in factor space.

The specific risk matrix  $\mathbf{D}$  is estimated from the residuals of the factor regression. Specifically,  $D_{nn} = \hat{\sigma}_{u_n}^2 = \frac{1}{T-K-1} \sum_{t=1}^T \hat{u}_{nt}^2$ , where  $\hat{u}_{nt} = r_{nt} - \hat{\alpha}_n - \sum_k \hat{\beta}_{nk} f_{kt}$  are the regression residuals. The diagonality assumption (that idiosyncratic returns are uncorrelated across assets) is a strong but useful simplification that ensures  $\mathbf{D}$  is well-conditioned and that the factor model captures all systematic co-movement through the common factors.

A further advantage of the factor model structure is that it facilitates view specification in the Black-Litterman framework. Views can be expressed directly on factor returns rather than individual asset returns, and these factor-level views propagate to asset-level expectations through the loading matrix:  $\mathbb{E}[\mathbf{r}] = \mathbf{B}\mathbb{E}[\mathbf{f}]$ . This is conceptually natural when views derive from macroeconomic analysis (which naturally targets factor premia) rather than individual stock picking.

### 2.3.4 Factor Model Prior

The factor model assembles loading matrix, factor covariance, and specific risk into a complete prior distribution that can replace or complement the empirical prior in downstream optimization. The factor model prior accepts a specification for how factor-level moments are estimated. For

example, using the sample mean and covariance of factor returns as the factor-level distribution, these moments are then mapped to asset space through the loading matrix.

Factor returns are passed as auxiliary data alongside the asset return matrix, allowing factor models to participate in cross-validation and pipeline composition alongside other estimators.

## 2.4 Regime-Switching Models

Financial return distributions exhibit persistent shifts in their statistical properties across different market environments. Bull markets, bear markets, high-volatility crises, and calm recovery periods each produce distinct patterns of means, variances, and correlations. The static estimators discussed above treat the return-generating process as stationary, averaging over regime differences to produce a single set of moment estimates. Regime-switching models explicitly model the transitions between these states, producing time-varying moment estimates that adapt to the prevailing market environment.

### 2.4.1 Hidden Markov Models for Regime Detection

The Hidden Markov Model (HMM), introduced to financial econometrics by Hamilton (1989), posits that observed asset returns are generated by a process governed by a discrete latent state that evolves as a first-order Markov chain.

**Generative model.** Let  $z_t \in \{1, \dots, S\}$  denote the latent regime at time  $t$ . State transitions are governed by a transition matrix  $\mathbf{A} \in \mathbb{R}^{S \times S}$ :

$$p(z_t = j \mid z_{t-1} = i) = A_{ij}$$

with  $A_{ij} \geq 0$  and  $\sum_{j=1}^S A_{ij} = 1$  for each row  $i$ . The initial state distribution is  $\pi_{0,s} = p(z_1 = s)$ . Conditional on the latent state, the  $N$ -dimensional return vector is drawn from a state-specific Gaussian emission:

$$\mathbf{r}_t \mid z_t = s \sim \mathcal{N}(\boldsymbol{\mu}_s, \boldsymbol{\Sigma}_s)$$

where  $\boldsymbol{\mu}_s \in \mathbb{R}^N$  and  $\boldsymbol{\Sigma}_s \in \mathbb{R}^{N \times N}$  are the state-conditional mean vector and covariance matrix. The joint distribution of the complete observation and state sequences factorizes as:

$$p(\mathbf{r}_{1:T}, z_{1:T}) = p(z_1) \prod_{t=2}^T p(z_t \mid z_{t-1}) \prod_{t=1}^T p(\mathbf{r}_t \mid z_t)$$

Two-state models (bull/bear or low-volatility/high-volatility) capture the dominant regime structure in equity markets. Three-state models add an intermediate regime that accommodates transitional periods. The number of states  $S$  is a structural hyperparameter, selected by information criteria (AIC, BIC) or cross-validated predictive performance.

**Parameter estimation via Baum-Welch.** The model parameters  $\theta = \{\mathbf{A}, \boldsymbol{\pi}_0, \{\boldsymbol{\mu}_s, \boldsymbol{\Sigma}_s\}_{s=1}^S\}$  are estimated from observed returns using the Baum-Welch algorithm, a special case of expectation-maximization (EM). The algorithm alternates between computing expected sufficient statistics under the current parameter estimates and updating those parameters to maximize the expected complete-data log-likelihood.

The E-step employs the forward-backward algorithm. The forward variable  $\alpha_t(s) = p(\mathbf{r}_{1:t}, z_t = s)$  satisfies the recursion:

$$\alpha_t(s) = \left[ \sum_{s'=1}^S \alpha_{t-1}(s') A_{s',s} \right] p(\mathbf{r}_t \mid z_t = s)$$

initialized with  $\alpha_1(s) = \pi_{0,s} p(\mathbf{r}_1 \mid z_1 = s)$ . The backward variable  $\beta_t(s) = p(\mathbf{r}_{t+1:T} \mid z_t = s)$  satisfies:

$$\beta_t(s) = \sum_{s'=1}^S A_{s,s'} p(\mathbf{r}_{t+1} \mid z_{t+1} = s') \beta_{t+1}(s')$$

initialized with  $\beta_T(s) = 1$ . The smoothed state probabilities and pairwise transition probabilities follow as:

$$\begin{aligned} \gamma_t(s) &= p(z_t = s \mid \mathbf{r}_{1:T}) = \frac{\alpha_t(s) \beta_t(s)}{p(\mathbf{r}_{1:T})} \\ \xi_t(i, j) &= p(z_{t-1} = i, z_t = j \mid \mathbf{r}_{1:T}) = \frac{\alpha_{t-1}(i) A_{ij} p(\mathbf{r}_t \mid z_t = j) \beta_t(j)}{p(\mathbf{r}_{1:T})} \end{aligned}$$

The M-step updates the parameters using these sufficient statistics:

$$\begin{aligned} \hat{A}_{ij} &= \frac{\sum_{t=2}^T \xi_t(i, j)}{\sum_{t=2}^T \gamma_{t-1}(i)}, \quad \hat{\boldsymbol{\mu}}_s = \frac{\sum_{t=1}^T \gamma_t(s) \mathbf{r}_t}{\sum_{t=1}^T \gamma_t(s)} \\ \hat{\boldsymbol{\Sigma}}_s &= \frac{\sum_{t=1}^T \gamma_t(s) (\mathbf{r}_t - \hat{\boldsymbol{\mu}}_s)(\mathbf{r}_t - \hat{\boldsymbol{\mu}}_s)^\top}{\sum_{t=1}^T \gamma_t(s)} \end{aligned}$$

The E-step and M-step alternate until convergence, monotonically increasing the observed-data log-likelihood  $\log p(\mathbf{r}_{1:T} \mid \theta)$  at each iteration.

**Online filtering.** For production deployment, forward filtering provides the causal filtered state probability at the current time:

$$p(z_t = s \mid \mathbf{r}_{1:t}) \propto p(\mathbf{r}_t \mid z_t = s) \sum_{s'=1}^S A_{s',s} p(z_{t-1} = s' \mid \mathbf{r}_{1:t-1})$$

This recursive update uses only past and current observations, maintaining the temporal causality required for real-time portfolio construction. No future data contaminates the state estimate, a property that is critical for backtest validity.

#### 2.4.2 Regime-Conditional Moment Estimation

The HMM's state-conditional parameters yield time-varying moment estimates through probability-weighted blending. Given the filtered state probabilities at time  $t$ , the blended expected return vector is:

$$\hat{\boldsymbol{\mu}}_t = \sum_{s=1}^S p(z_t = s \mid \mathbf{r}_{1:t}) \boldsymbol{\mu}_s$$

The blended covariance matrix must account for both within-state and between-state variation in means:

$$\hat{\boldsymbol{\Sigma}}_t = \sum_{s=1}^S p(z_t = s \mid \mathbf{r}_{1:t}) \left[ \boldsymbol{\Sigma}_s + (\boldsymbol{\mu}_s - \hat{\boldsymbol{\mu}}_t)(\boldsymbol{\mu}_s - \hat{\boldsymbol{\mu}}_t)^\top \right]$$

The second term inflates the blended covariance by the cross-state dispersion of means. When the regime is ambiguous (filtered probabilities near uniform), this term is large, producing appropriately conservative risk estimates. When the regime is clearly identified (one state probability near unity), the blended moments approximate the conditional moments of the dominant state.

During a high-probability bear market regime, the blended estimates shift toward the bear-state parameters: lower expected returns, higher volatilities, and elevated correlations. During a bull regime, the reverse holds. Transitions between regimes are smooth, governed by the continuous evolution of filtered probabilities, which avoids the instability that hard regime switching would introduce.

These time-varying moments serve as direct inputs to the empirical prior and, through it, to the optimization pipeline. Any of the covariance estimators discussed in the preceding sections (shrinkage, denoising, Graphical Lasso) can be applied within each regime to regularize the state-conditional covariance estimates, a technique that is particularly valuable when the effective sample size per regime is small.

#### 2.4.3 Deep Markov Models

Hidden Markov Models are constrained by their discrete state space: the number of distinct regime configurations is fixed at  $S$ , and transitions between regimes are instantaneous. Deep Markov Models (DMMs) generalize this framework by replacing discrete latent states with continuous latent vectors and parameterizing the transition and emission distributions with neural networks. This generalization enables the model to capture complex, nonlinear dynamics in the return-generating process that discrete-state HMMs cannot represent.

**Generative model.** The DMM posits a continuous latent state  $\mathbf{z}_t \in \mathbb{R}^d$  that evolves through a neural-network-parameterized transition distribution:

$$p_\theta(\mathbf{z}_t \mid \mathbf{z}_{t-1}) = \mathcal{N}\left(\mathbf{z}_t \mid \boldsymbol{\mu}_\theta^{\text{trans}}(\mathbf{z}_{t-1}), \text{diag}\left(\boldsymbol{\sigma}_\theta^{\text{trans}}(\mathbf{z}_{t-1})^2\right)\right)$$

where  $\boldsymbol{\mu}_\theta^{\text{trans}}$  and  $\boldsymbol{\sigma}_\theta^{\text{trans}}$  are outputs of a neural network (the transition network) with parameters  $\theta$ . Observed returns are generated from a state-dependent emission distribution:

$$p_\theta(\mathbf{r}_t \mid \mathbf{z}_t) = \mathcal{N}\left(\mathbf{r}_t \mid \boldsymbol{\mu}_\theta^{\text{emit}}(\mathbf{z}_t), \text{diag}\left(\boldsymbol{\sigma}_\theta^{\text{emit}}(\mathbf{z}_t)^2\right)\right)$$

The joint distribution factorizes as:

$$p_\theta(\mathbf{r}_{1:T}, \mathbf{z}_{1:T}) = p(\mathbf{z}_1) \prod_{t=2}^T p_\theta(\mathbf{z}_t \mid \mathbf{z}_{t-1}) \prod_{t=1}^T p_\theta(\mathbf{r}_t \mid \mathbf{z}_t)$$

**Gated transitions.** A refinement introduces gating mechanisms that interpolate between linear and nonlinear dynamics:

$$\boldsymbol{\mu}_\theta^{\text{trans}}(\mathbf{z}_{t-1}) = (1 - \mathbf{g}_t) \odot (\mathbf{W}\mathbf{z}_{t-1} + \mathbf{b}) + \mathbf{g}_t \odot \mathbf{h}_\theta(\mathbf{z}_{t-1})$$

where  $\mathbf{g}_t = \sigma(\mathbf{W}_g \mathbf{z}_{t-1} + \mathbf{b}_g)$  is a sigmoid gate,  $\mathbf{W}\mathbf{z}_{t-1} + \mathbf{b}$  is the linear component,  $\mathbf{h}_\theta$  is a nonlinear neural network, and  $\odot$  denotes elementwise multiplication. When the gate activations are near zero, the transition is approximately linear, resembling a linear state-space model. When the gate activations approach one, the transition is fully nonlinear. This adaptive complexity allows the model to allocate nonlinearity only where the data requires it.

**Variational inference.** The posterior  $p_\theta(\mathbf{z}_{1:T} \mid \mathbf{r}_{1:T})$  is intractable due to the nonlinear dependencies introduced by the neural networks. Amortized variational inference addresses this by introducing a structured inference network  $q_\phi(\mathbf{z}_{1:T} \mid \mathbf{r}_{1:T})$  that factorizes as:

$$q_\phi(\mathbf{z}_{1:T} \mid \mathbf{r}_{1:T}) = \prod_{t=1}^T q_\phi(\mathbf{z}_t \mid \mathbf{z}_{t-1}, \mathbf{r}_{t:T})$$

Each factor conditions on both the previous latent state and future observations, which are processed by a backward-running recurrent neural network. The inference network and generative model are trained jointly by maximizing the evidence lower bound (ELBO):

$$\mathcal{L}(\theta, \phi) = \mathbb{E}_{q_\phi} \left[ \sum_{t=1}^T \log p_\theta(\mathbf{r}_t \mid \mathbf{z}_t) \right] - \sum_{t=1}^T \text{KL}[q_\phi(\mathbf{z}_t \mid \cdot) \parallel p_\theta(\mathbf{z}_t \mid \mathbf{z}_{t-1})]$$

The first term rewards accurate reconstruction of observed returns. The second term penalizes

deviation of the approximate posterior from the generative prior, regularizing the latent dynamics against overfitting.

**Predictive moment extraction.** Given observations  $\mathbf{r}_{1:t}$  up to the current time, the DMM produces a distribution over the next-period latent state and consequently a predictive distribution over returns. The predictive moments are:

$$\hat{\boldsymbol{\mu}}_{t+1} = \mathbb{E}_{p(\mathbf{z}_{t+1}|\mathbf{r}_{1:t})} \left[ \boldsymbol{\mu}_{\theta}^{\text{emit}}(\mathbf{z}_{t+1}) \right]$$

$$\hat{\boldsymbol{\Sigma}}_{t+1} = \mathbb{E}_{p(\mathbf{z}_{t+1}|\mathbf{r}_{1:t})} \left[ \text{diag} \left( \boldsymbol{\sigma}_{\theta}^{\text{emit}}(\mathbf{z}_{t+1})^2 \right) + \boldsymbol{\mu}_{\theta}^{\text{emit}}(\mathbf{z}_{t+1}) \boldsymbol{\mu}_{\theta}^{\text{emit}}(\mathbf{z}_{t+1})^{\top} \right] - \hat{\boldsymbol{\mu}}_{t+1} \hat{\boldsymbol{\mu}}_{t+1}^{\top}$$

These expectations are approximated by Monte Carlo sampling from the posterior predictive distribution. The resulting time-varying moments substitute for or complement the static estimators discussed earlier in this chapter.

#### 2.4.4 Model Selection Guidance

The regime-switching models occupy a distinct position in the estimator landscape, complementing rather than replacing the static estimators.

Condition	Recommended Approach	Rationale
Clear regime structure (bull/bear)	HMM with 2–3 states	Interpretable, parsimonious, well-understood
Smooth, continuous regime variation	Deep Markov Model	Captures nonlinear dynamics in continuous latent space
Limited data ( $T < 500$ )	Static estimators (shrinkage, factor model)	Insufficient data for reliable regime inference
High-dimensional universe ( $N > 100$ )	HMM on factor returns, mapped to assets	Reduces dimensionality of emission model
Regime transitions known a priori	Fixed-regime conditional estimation	Avoids estimation of transition dynamics

HMMs are preferred when interpretability matters and the regime structure is expected to be discrete and low-dimensional. DMMs are preferred when the underlying dynamics are complex, the data is abundant, and the practitioner is willing to trade interpretability for representational capacity. Static estimators remain appropriate when regime dynamics are not a primary concern or when the sample is too short for reliable regime inference.

### 2.5 The Empirical Prior

The empirical prior assembles an expected return estimator and a covariance estimator into a single object that provides the complete input specification for downstream optimization. This

assembly is the critical interface between the estimation stage and the optimization stage of the portfolio construction pipeline: the optimizer receives a prior distribution and extracts the moments it requires, without needing to know how those moments were estimated.

The empirical prior accepts an expected return estimator and a covariance estimator as components, allowing any combination of the estimators discussed above. A typical composition might pair shrinkage toward the grand mean for expected returns with Ledoit-Wolf shrinkage for the covariance matrix, producing a prior that is regularized in both the mean and covariance dimensions. The modularity of this design means that changing the covariance estimator (for instance, switching from Ledoit-Wolf to random matrix theory denoising) requires modifying a single component without altering any other part of the pipeline.

For multi-period optimization, the prior supports log-normal return assumptions. When enabled, the prior compounds single-period moments to produce multi-period expected returns and covariances consistent with the geometric growth properties of wealth accumulation. Setting the investment horizon (for instance, 252 periods for annualization when working with daily data) applies the appropriate compounding correction. This adjustment is important because mean-variance optimization at long horizons should account for the volatility drag that reduces compound returns below arithmetic returns:

$$\mathbb{E}[\log(1 + R_p)] \approx \mathbb{E}[R_p] - \frac{1}{2}\sigma_p^2$$

The prior serves as the modular interface that decouples estimation from optimization. Any estimator conforming to the prior protocol can be substituted without downstream changes, enabling systematic comparison of estimation approaches through cross-validation.

## 2.6 LLM-Augmented Moment Estimation

Large language models introduce a qualitative information channel into the otherwise quantitative moment estimation pipeline. While the estimators discussed above operate exclusively on numerical return data, LLMs can process unstructured text (economic reports, central bank communications, earnings transcripts, news articles) and translate qualitative assessments into quantitative adjustments to the estimation process. This augmentation does not replace statistical estimation but rather informs the configuration and parameterization of the estimators.

**Risk aversion calibration.** The equilibrium return estimator depends critically on the risk aversion parameter  $\delta$ , which determines the scale of implied returns. LLMs can classify the current macroeconomic regime (early-cycle, mid-cycle, late-cycle, or recession) by processing leading indicators such as purchasing managers' indices, yield curve slopes, and credit spreads. This regime classification maps to appropriate risk aversion parameters: expansionary regimes warrant lower risk aversion (higher implied returns, more aggressive positioning), while late-cycle or recessionary regimes warrant higher risk aversion (lower implied returns, more defensive positioning). The mapping from regime to  $\delta$  translates qualitative macroeconomic assessment into a precise numerical



input for the equilibrium return estimator.

**Factor weight adaptation.** Business cycle phase influences which systematic factors are likely to deliver premiums over the subsequent investment horizon. LLMs can recommend shifting factor emphasis based on current economic conditions; for instance, tilting toward quality and value factors in late-cycle environments where earnings resilience and margin of safety become more important, or toward momentum and growth factors in early-cycle recoveries where economic acceleration favors high-beta, growth-oriented exposures. These recommendations adjust the loading matrix or factor covariance inputs in the factor model prior, shaping the risk model to reflect the anticipated regime.

**Covariance regime selection.** The choice of covariance estimator itself can be informed by qualitative analysis. News sentiment analysis can detect whether the current environment favors responsive estimators (such as exponentially weighted covariance during regime transitions, when correlations and volatilities are shifting rapidly) or stable estimators (such as Ledoit-Wolf during calm periods when the return distribution is approximately stationary). An LLM monitoring news flow and policy communications can flag transitions between these regimes, triggering automated switches in the covariance estimation methodology.

**Confidence calibration.** Perhaps most subtly, LLMs can assess the reliability of moment estimates by cross-referencing quantitative signals with qualitative information. When historical factor premia are consistent with the narrative embedded in economic data and corporate earnings, the moment estimates can be trusted with higher confidence. When quantitative signals diverge from the qualitative picture (for instance, when momentum signals remain positive but leading indicators suggest deterioration) the LLM can flag the discrepancy, prompting wider uncertainty bands in the prior or greater shrinkage toward neutral estimates. This meta-level assessment of estimate reliability adds a layer of robustness that purely quantitative systems lack.

### 3 View Integration and Bayesian Updating

Raw moment estimates, however carefully constructed, reflect only backward-looking statistical patterns. Active portfolio management requires incorporating forward-looking views: beliefs about expected returns, volatilities, correlations, or tail behavior that deviate from what history alone implies. This chapter presents three complementary frameworks for view integration and establishes LLM-driven view generation as a systematic approach to populating these frameworks.

#### 3.1 The Black-Litterman Framework

##### 3.1.1 The Foundational Problem

Traditional **mean-variance optimization** suffers from extreme sensitivity to expected return estimates. Small changes in forecasted returns produce dramatically different portfolios; the optimizer aggressively exploits estimation error by taking extreme positions in assets with overstated returns. The resulting allocations are unstable, unintuitive, and dominated by noise rather than signal.

**Black and Litterman (1992)** resolved this instability through Bayesian inference: rather than treating expected returns as known inputs to be estimated directly, the framework begins with **equilibrium returns** as a neutral prior distribution and then incorporates **active views** with explicit confidence levels. The posterior distribution over expected returns blends market-implied information with the investor's subjective or model-driven beliefs, producing portfolios that are stable, diversified, and responsive to conviction in proportion to its precision.

##### 3.1.2 The Equilibrium Prior

Market equilibrium returns  $\Pi$  represent the expected returns implied by current market capitalization weights under the assumption that markets are in equilibrium:

$$\Pi = \delta \Sigma \mathbf{w}_{\text{mkt}}$$

where:

- $\delta$  is the **risk aversion coefficient**, reflecting the aggregate investor's trade-off between expected return and variance
- $\Sigma$  is the **covariance matrix** of asset returns
- $\mathbf{w}_{\text{mkt}}$  is the vector of **market capitalization weights**, serving as the prior portfolio

The interpretation is straightforward: if all investors held the market portfolio with risk aversion  $\delta$ , equilibrium expected returns must satisfy the reverse optimization equation above. This provides a **neutral starting point** reflecting collective market wisdom rather than any individual forecast. The equilibrium prior anchors the model, ensuring that in the absence of active views, the resulting portfolio coincides with the market portfolio.

**Risk aversion estimation** derives from observable market quantities:

$$\delta = \frac{E[R_{\text{mkt}}] - R_f}{\sigma_{\text{mkt}}^2}$$

where  $E[R_{\text{mkt}}] - R_f$  is the equity risk premium and  $\sigma_{\text{mkt}}^2$  is market variance.

### 3.1.3 Active View Specification

Investors express **active views** through the linear constraint system:

$$\mathbf{P}\boldsymbol{\mu} = \mathbf{Q} + \boldsymbol{\epsilon}, \quad \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Omega})$$

where:

- $\mathbf{P}$  is a  $K \times N$  **pick matrix** defining which assets each view concerns
- $\mathbf{Q}$  is a  $K \times 1$  vector of **view returns** (expected outcomes)
- $\boldsymbol{\Omega}$  is a  $K \times K$  **view uncertainty matrix** (diagonal for independent views)

Three canonical view types arise from different structures of  $\mathbf{P}$ :

**Absolute views** specify an expected return for a single asset. The corresponding row of  $\mathbf{P}$  contains a single entry of one in the column of the target asset and zeros elsewhere, with  $Q$  equal to the expected return:

$$\mathbf{P} = [0 \ \cdots \ 0 \ 1 \ 0 \ \cdots \ 0], \quad Q = q$$

For example, an absolute view that a particular equity will return ten percent annualized places one in its column and sets  $Q = 0.10$ .

**Relative views** express an expected outperformance of one asset over another. The pick matrix contains +1 for the outperforming asset and -1 for the underperforming asset:

$$\mathbf{P} = [0 \ \cdots \ 1_{(i)} \ \cdots \ -1_{(j)} \ \cdots \ 0], \quad Q = q$$

where  $q$  represents the expected return differential.

**Basket views** extend relative views to groups of assets. The pick matrix assigns fractional weights to group members, enabling views on sector-level or theme-level performance:

$$\mathbf{P} = \left[ \frac{1}{n_1} \ \cdots \ \frac{1}{n_1} \ 0 \ \cdots \ 0 \right], \quad Q = q$$

Views can be specified using **string expressions** that are automatically parsed into the  $\mathbf{P}$  and  $\mathbf{Q}$  matrices. This design permits programmatic view generation, an essential capability for LLM-driven pipelines that produce views as structured text.

### 3.1.4 Bayesian Posterior Returns

The **Black-Litterman posterior expected returns** combine equilibrium and views via Bayesian updating:

$$\mathbb{E}[\boldsymbol{\mu}] = \bar{\boldsymbol{\mu}} = \left[ (\tau \boldsymbol{\Sigma})^{-1} + \mathbf{P}^\top \boldsymbol{\Omega}^{-1} \mathbf{P} \right]^{-1} \left[ (\tau \boldsymbol{\Sigma})^{-1} \boldsymbol{\Pi} + \mathbf{P}^\top \boldsymbol{\Omega}^{-1} \mathbf{Q} \right]$$

with posterior uncertainty:

$$\mathbf{M} = \left[ (\tau \boldsymbol{\Sigma})^{-1} + \mathbf{P}^\top \boldsymbol{\Omega}^{-1} \mathbf{P} \right]^{-1}$$

The parameter  $\tau$  is the **uncertainty scaling parameter**, typically set between 0.025 and 0.05, representing the relative uncertainty in equilibrium return estimates compared to the covariance matrix. A smaller  $\tau$  implies greater confidence in the equilibrium prior, anchoring the posterior more firmly to market-implied returns.

The **intuition** behind the posterior formula is that of a precision-weighted average. The precision of the equilibrium prior is  $(\tau \boldsymbol{\Sigma})^{-1}$  and the precision of the views is  $\mathbf{P}^\top \boldsymbol{\Omega}^{-1} \mathbf{P}$ . The posterior combines these two sources of information in proportion to their respective precisions:

- **High view confidence** (small  $\boldsymbol{\Omega}$ ): posterior returns tilt strongly toward the active views
- **Low view confidence** (large  $\boldsymbol{\Omega}$ ): posterior returns remain close to the equilibrium prior
- **No views**: posterior returns reduce exactly to the equilibrium  $\boldsymbol{\Pi}$

This graceful degradation ensures that the framework never produces worse results than the equilibrium baseline, regardless of view quality.

### 3.1.5 View Uncertainty Calibration

The specification of the view uncertainty matrix  $\boldsymbol{\Omega}$  is critical to the practical success of the Black-Litterman framework. Three approaches, each with distinct advantages, have gained prominence.

The **He-Litterman proportional approach** sets view uncertainty proportional to the variance of the view portfolio:

$$\omega_k = \tau \cdot \mathbf{P}_k^\top \boldsymbol{\Sigma} \mathbf{P}_k$$

This ensures that views on volatile assets or asset combinations carry proportionally more uncertainty, preventing the optimizer from over-reacting to views on inherently noisy assets. The approach requires no subjective inputs beyond  $\tau$  itself.

The **Idzorek (2005) extension** introduces an explicit view strength parameter  $\alpha_k$ :

$$\omega_k = \frac{1}{\alpha_k} \cdot \mathbf{P}_k^\top \boldsymbol{\Sigma} \mathbf{P}_k$$

where  $\alpha_k \in [0, \infty)$  calibrates the degree to which each view influences the posterior:

- $\alpha_k = 1.0$ : full confidence, the view is treated as effectively certain
- $\alpha_k = 0.5$ : moderate confidence, the view produces a meaningful but tempered tilt
- $\alpha_k = 0.1$ : minimal confidence, the view produces only a marginal adjustment

This parameterization is particularly amenable to LLM-driven confidence calibration, as the strength parameter maps naturally to a continuous confidence score.

**Direct specification from forecast error** derives view uncertainty from the historical accuracy of the view-generating process:

$$\omega_k = \text{Var}(Q_k - \text{Realized Return}_k)$$

estimated from a track record of past views and their outcomes. This approach is the most principled when sufficient historical data exists, as it directly captures the empirical reliability of the forecasting methodology.

### 3.1.6 Black-Litterman Factor Model

Rather than specifying views on individual assets, the factor model variant targets **factor returns**:

$$\mathbf{P}_f \boldsymbol{\mu}_f = \mathbf{Q}_f + \boldsymbol{\epsilon}_f$$

where  $\boldsymbol{\mu}_f$  is the vector of expected factor returns. A view that “the momentum factor will return twelve percent annually” or “quality will outperform value by four percent” operates in the factor space. Asset-level expectations are then derived through the factor loading matrix:

$$\bar{\boldsymbol{\mu}} = \mathbf{B} \bar{\boldsymbol{\mu}}_f$$

where  $\mathbf{B}$  is the  $N \times K$  matrix of asset-to-factor loadings. This approach offers three advantages. First, it is more **parsimonious**: a small number of factor views implicitly generates views on all assets through their factor exposures. Second, it is more **stable**: factors are better-estimated than individual asset returns, producing smoother posterior distributions. Third, it is more **interpretable**: views on well-defined investment themes (value, momentum, quality, size) are easier to formulate, justify, and communicate than views on hundreds of individual securities.

This architecture is realized by combining a factor model with a Black-Litterman prior applied at the factor level, so that views on factors propagate to asset-level expectations through the loading matrix.

## 3.2 Entropy Pooling: Non-Linear View Integration

### 3.2.1 Limitations of Classical Black-Litterman

The classical Black-Litterman framework, for all its elegance, restricts views to **linear functions of expected returns** ( $\mathbf{P}\boldsymbol{\mu} = \mathbf{Q}$ ) and assumes Gaussian uncertainty on those views. This formulation excludes a wide range of practically important beliefs:

- Views on **variance**: “volatility of asset  $i$  will double over the next quarter”
- Views on **correlation**: “the correlation between assets  $i$  and  $j$  will increase to 0.80”
- Views on **skewness**: “asset  $i$  will exhibit negative skew due to downside tail risk”
- Views on **tail risk**: “the 95% CVaR of asset  $i$  will reach eight percent”

These limitations motivate a more general framework capable of incorporating views on arbitrary distributional properties.

### 3.2.2 Minimum Kullback-Leibler Divergence Framework

**Entropy Pooling**, introduced by Meucci (2008), overcomes these limitations by adjusting **scenario probabilities** rather than modifying return parameters directly. Given  $S$  historical or simulated scenarios with baseline (equal) probabilities  $\mathbf{p}_0 = (1/S, \dots, 1/S)$ , the framework seeks new probabilities  $\mathbf{p}^*$  that satisfy the investor’s view constraints while remaining as close as possible to the prior distribution.

The optimization problem minimizes the **Kullback-Leibler divergence** between the new and prior probability vectors:

$$\mathbf{p}^* = \arg \min_{\mathbf{p}} \sum_{s=1}^S p_s \ln \left( \frac{p_s}{p_{0,s}} \right)$$

subject to the normalization constraint  $\sum_s p_s = 1$ , non-negativity  $p_s \geq 0$ , and arbitrary view constraints expressed as moment conditions on the scenario-weighted distribution.

The Kullback-Leibler divergence provides the natural measure of **information loss**: among all probability distributions satisfying the view constraints, the solution  $\mathbf{p}^*$  is the one that introduces the least additional structure beyond what the views require. This principle of minimum relative entropy ensures that the posterior distribution reflects the views and nothing more; no spurious correlations or distributional artifacts are introduced.

### 3.2.3 View Types Supported

Entropy Pooling accommodates views on any moment or distributional property that can be expressed as a constraint on scenario probabilities. The following catalogue covers the principal view types, each with its mathematical constraint formulation.

**Mean views** on expected returns take the form of equality or inequality constraints:

$$\sum_{s=1}^S p_s \cdot r_{i,s} = \mu_i^{\text{view}} \quad (\text{equality}), \quad \sum_{s=1}^S p_s \cdot r_{i,s} \geq \mu_i^{\text{lower}} \quad (\text{inequality})$$

Equality views fix the expected return precisely, while inequality views set a floor or ceiling, allowing the optimizer to determine the exact level that minimizes information loss.

**Variance views** constrain the second central moment:

$$\sum_{s=1}^S p_s \cdot (r_{i,s} - \bar{r}_i)^2 = \sigma_i^{2,\text{view}}$$

enabling beliefs about future volatility levels to be incorporated without altering expected returns or correlations beyond what the variance constraint requires.

**Correlation views** constrain the normalized cross-moment between two assets:

$$\frac{\sum_s p_s (r_{i,s} - \bar{r}_i)(r_{j,s} - \bar{r}_j)}{\sigma_i^{\text{view}} \cdot \sigma_j^{\text{view}}} = \rho_{ij}^{\text{view}}$$

This is particularly valuable for expressing beliefs about regime changes in co-movement structure; for instance, that correlations will converge toward one during a stress period.

**Skewness views** constrain the standardized third moment:

$$\sum_{s=1}^S p_s \cdot \left( \frac{r_{i,s} - \bar{r}_i}{\sigma_i} \right)^3 = \gamma_i^{\text{view}}$$

allowing the investor to express beliefs about distributional asymmetry: negative skew for assets facing downside tail risk, positive skew for assets with optionality or convex payoff structures.

**CVaR views** constrain the expected loss in the tail:

$$-\frac{1}{\alpha} \sum_{s \in \text{worst } \alpha S} p_s \cdot r_{i,s} = \text{CVaR}_i^{\text{view}}$$

where  $\alpha$  defines the tail probability (typically five percent). This enables direct expression of tail risk beliefs without parametric distributional assumptions.

**Group views** aggregate across asset classifications:

$$\sum_{i \in G} w_i \cdot \bar{r}_i^{\text{view}} = \mu_G^{\text{view}}$$

where groups  $G$  can represent sectors, themes, factor exposures, or any other asset classification.

Group views are natural for top-down allocation, where the investor has convictions about broad categories rather than individual securities.

### 3.2.4 Views Relative to Prior

A particularly powerful feature of Entropy Pooling allows expressing views **relative to the prior distribution**, using multiplicative adjustments without specifying absolute values:

- “Variance of asset  $i$  will be four times its prior level”
- “Correlation between assets  $i$  and  $j$  will be half the prior level”
- “Expected return of asset  $i$  will be twenty percent above its prior level”

Relative views eliminate the need to compute absolute moment values manually, substantially reducing calibration error. The investor need only specify the **direction and magnitude of deviation** from the current distributional estimate, which is often a more natural and robust form of belief expression than absolute forecasts.

## 3.3 Opinion Pooling: Combining Multiple Expert Views

### 3.3.1 The Multi-Expert Problem

When multiple sources generate views (quantitative models, fundamental analysis, macro strategies, LLM-based signals) their opinions must be combined into a single coherent input for the optimization framework. Each source may possess different accuracy characteristics, different domains of expertise, and potentially conflicting conclusions. Naive approaches such as simple averaging ignore these heterogeneities, while ad hoc reconciliation introduces subjective bias and lacks theoretical grounding.

### 3.3.2 Consensus Distribution with Credibility Weights

**Opinion Pooling** provides a principled framework for multi-expert combination. Given  $M$  expert posterior distributions  $\{P_1, P_2, \dots, P_M\}$  and a base prior  $P_0$ , the consensus distribution is formed as:

$$P^* = \left(1 - \sum_{m=1}^M \pi_m\right) P_0 + \sum_{m=1}^M \pi_m P_m$$

where  $\pi_m \in [0, 1]$  represents the **credibility weight** assigned to expert  $m$ , and the constraint  $\sum_m \pi_m \leq 1$  ensures the base prior retains residual influence. The base prior acts as a shrinkage target, anchoring the consensus toward the equilibrium or empirical distribution when expert opinions are uncertain or contradictory.

The credibility weights  $\pi_m$  govern the influence of each expert on the final distribution. When  $\sum_m \pi_m = 1$ , the base prior receives zero weight and the consensus is a pure mixture of expert opinions. When  $\sum_m \pi_m < 1$ , the residual weight  $1 - \sum_m \pi_m$  accrues to the base prior, providing a regularization effect that prevents any single expert from dominating the consensus.



### 3.3.3 Expert Weight Calibration

Calibration of expert weights  $\pi_m$  can proceed along several dimensions:

**Historical accuracy** provides the most direct calibration signal. By tracking each expert’s information coefficient (the correlation between forecasted and realized returns) over time, the weight assignment reflects demonstrated predictive ability. Experts with consistently higher information coefficients receive larger weights.

**Domain relevance** adjusts weights according to the scope of each expert’s competence. A macro-focused model may receive higher weight for interest-rate-sensitive assets but lower weight for idiosyncratic stock-level views. Conversely, a bottom-up fundamental model may be weighted more heavily for individual security views than for broad sector or factor calls.

**Source diversification** penalizes redundant signals from correlated information sources. If two experts rely on overlapping data or similar methodologies, their combined weight should reflect their incremental (not total) information content. This prevents double-counting and ensures the consensus benefits from genuinely diverse perspectives.

### 3.3.4 Conflicting View Resolution

The Opinion Pooling framework handles disagreement naturally, without requiring manual reconciliation. If one expert is bullish and another bearish on the same asset, the consensus distribution reflects the **probability-weighted blend** of their views. The base prior acts as an anchor when experts disagree sharply, tempering the consensus toward the neutral equilibrium estimate.

This automatic conflict resolution is a significant advantage over sequential or hierarchical view integration schemes, where the order of incorporation can influence the final result. Under Opinion Pooling, all experts contribute simultaneously, and the consensus depends only on their respective credibility weights and posterior distributions.

## 3.4 LLM-Driven View Generation and Integration

### 3.4.1 Structured View Generation from Multi-Factor Analysis

Large language models introduce a qualitatively new approach to view generation by processing **comprehensive multi-dimensional data** for each asset in the investment universe. The analytical dimensions span:

**Valuation:** price-to-earnings, price-to-book, price-to-sales ratios compared against sector norms and historical distributions. Persistent deviations from fair value generate mean-reversion views; justified deviations (supported by growth or quality differentials) produce continuation views.

**Momentum:** three-month, six-month, and twelve-month returns, relative strength index, and trend strength metrics. The LLM evaluates whether momentum signals reflect genuine information diffusion or speculative excess, distinguishing sustainable trends from fragile overextensions.

**Quality:** return on equity, profit margins, cash flow stability, and financial leverage. High-quality firms with durable competitive advantages generate positive expected return adjustments, while deteriorating quality metrics trigger negative views.

**Growth:** revenue and earnings growth rates, forward guidance trajectories, and estimate revision trends. The LLM assesses whether current growth rates are sustainable, accelerating, or decelerating, translating the assessment into directional views.

**Technical:** support and resistance levels, volume patterns, and market microstructure signals. These shorter-horizon indicators modulate view confidence rather than generating primary directional views.

**Analyst consensus:** ratings distributions, price target dispersions, and estimate revision momentum. The LLM synthesizes consensus information while detecting potential herding or stale estimates.

The output of this multi-factor analysis is a **structured view** for each asset, comprising an expected return estimate, a confidence level, and a natural-language rationale. The confidence calibration follows a systematic scale:

- **High confidence** (0.8 to 1.0): all analytical factors aligned, clear identifiable catalyst, strong data quality
- **Medium confidence** (0.6 to 0.8): four to five factors aligned, some residual uncertainty, mixed technical signals
- **Low confidence** (0.0 to 0.6): mixed or conflicting signals, data gaps, contradictory indicators across dimensions

These structured views map directly into the Black-Litterman  $\mathbf{Q}$  vector (expected returns) and  $\mathbf{\Omega}$  matrix (uncertainty), with the confidence score governing the Idzorek-style view strength parameter  $\alpha_k$ .

### 3.4.2 Macroeconomic Regime Classification and Regime-Adjusted Confidence

LLMs contribute to view integration not only through individual asset analysis but through **macroeconomic regime classification**, which calibrates the global parameters of the optimization framework.

The classification processes a broad set of leading indicators:

- Manufacturing purchasing managers' indices and their diffusion components
- Yield curve slope and term structure dynamics
- Credit spreads and their rate of change
- Unemployment claims and labor market breadth indicators
- Corporate earnings trends and profit margin trajectories
- Central bank communications and forward guidance shifts
- Commodity price dynamics and supply-demand balances

From these inputs, the LLM assigns the current environment to one of four canonical **business cycle phases**:

1. **Early-cycle (recovery)**: leading indicators inflecting upward, credit conditions easing, earnings troughing
2. **Mid-cycle (expansion)**: broad-based growth, moderate inflation, stable credit conditions
3. **Late-cycle (slowdown)**: leading indicators peaking, credit tightening, margin compression beginning
4. **Recession (contraction)**: negative output growth, rising unemployment, credit stress, earnings declining

Each phase implies distinct **regime-adjusted parameters** for the optimization framework:

**Risk aversion**  $\delta$  scales with the business cycle: higher values in late-cycle and recession phases reflect more conservative risk preferences and the empirical observation that risk premia widen during downturns. Lower values in early-cycle phases permit more aggressive positioning when the risk-reward balance favors risk-taking.

**Uncertainty**  $\tau$  scales with market volatility relative to its long-term average. Elevated realized and implied volatility during stress periods widens the uncertainty band around equilibrium returns, causing the posterior to anchor more firmly to the prior and reducing the influence of active views, a prudent response when forecast reliability deteriorates.

**Factor emphasis** shifts across the cycle: quality and value factors receive greater weight in late-cycle and recessionary environments, where capital preservation and fundamental soundness dominate. Momentum and growth factors receive greater emphasis in early-cycle and expansion phases, where trend persistence and earnings acceleration drive returns.

The regime classification feeds directly into Black-Litterman prior calibration, adjusting  $\delta$ ,  $\tau$ , and the relative weighting of factor-level views according to the identified phase.

The LLM-based regime classification discussed here can be complemented by the formal Hidden Markov Model framework developed in the moment estimation chapter. HMM filtered probabilities provide a quantitative, data-driven regime indicator that serves as an input to, or a cross-check on, the LLM’s qualitative assessment. When both approaches agree on the prevailing regime, the resulting parameter adjustments carry greater conviction. When they disagree, the practitioner can increase uncertainty parameters ( $\tau$ ,  $\Omega$ ) to reflect the ambiguity, a conservative response appropriate for periods of conflicting signals. Deep Markov Models extend this further by providing continuous regime characterizations that capture subtler transitions than the four discrete phases above.

### 3.4.3 Multi-LLM Opinion Pooling

Different prompts, models, or analytical frameworks applied to the same data function as **independent experts**, each producing a distinct set of views with associated confidence levels. This multiplicity maps naturally to the Opinion Pooling framework:

- Each LLM-expert generates a posterior distribution over asset returns (views with confidences)

- The opinion pooling estimator combines these posteriors with calibrated credibility weights  $\pi_m$
- The base prior  $P_0$  absorbs residual weight, ensuring stability when LLM-experts disagree

The principal benefit of multi-LLM opinion pooling is **robustness to individual model biases**. Any single language model may exhibit systematic tendencies: recency bias, anchoring to salient narratives, or over-weighting certain data types. By combining multiple models with diverse analytical orientations, these biases partially cancel, and the consensus captures a broader range of analytical perspectives.

**Weight calibration** for LLM-experts follows the same principles as for traditional expert combination. Each model’s historical accuracy is tracked through its information coefficient over successive rebalancing periods. Models demonstrating consistently higher predictive accuracy receive larger credibility weights, while models with poor or deteriorating track records are down-weighted. This adaptive calibration ensures that the opinion pool evolves toward reliance on the most reliable analytical sources.

#### 3.4.4 News Sentiment as View Confidence Modifier

LLMs process news feeds, earnings call transcripts, and analyst reports to extract **sentiment signals** that modulate view confidence rather than serving as primary view sources. The distinction is important: sentiment is treated as second-order information that adjusts the precision of existing views, not as a direct input to expected return estimation.

The sentiment-confidence interaction operates through several channels:

**Reinforcement:** strong positive sentiment aligned with a bullish fundamental view increases the confidence parameter  $\alpha_k$ , tightening the view uncertainty  $\omega_k$  and amplifying the view’s influence on the posterior. The rationale is that when both quantitative signals and qualitative narrative point in the same direction, the probability of the view being correct is higher.

**Contradiction:** negative sentiment opposing a bullish fundamental view decreases confidence, widening  $\omega_k$  and attenuating the view’s effect. Similarly, positive sentiment opposing a bearish view reduces the conviction in the negative outlook. The framework does not discard the view entirely but softens its impact in proportion to the degree of contradiction.

**Ambiguity:** when sentiment signals are mixed or inconclusive, overall confidence is reduced across all views for the affected asset. This conservative response reflects the increased uncertainty that accompanies conflicting information sources.

**Temporal decay** governs the weighting of sentiment signals over time. Recent news and earnings releases receive substantially more weight than older coverage, reflecting the rapid incorporation of information into market prices. A half-life parameter controls the rate of decay, with typical values ranging from one to four weeks depending on the asset’s information environment and liquidity.

## 4 Risk Measures, Diversification, and Hierarchical Methods

Portfolio construction requires explicit choices about how risk is measured, how it is distributed across positions, and whether the optimization should respect the natural clustering structure of assets. These choices are not incidental; they define the very meaning of an “optimal” portfolio. A variance-minimizing allocation differs fundamentally from one that minimizes conditional drawdowns, and an allocation that respects hierarchical asset structure differs from one that treats the investment universe as a flat collection of securities. This chapter presents the full taxonomy of risk measures, the risk budgeting and maximum diversification frameworks that allocate risk systematically, and the hierarchical methods that exploit asset correlation structure to produce stable allocations without matrix inversion. The integration of large language models into each of these layers introduces adaptive, forward-looking intelligence that complements the purely quantitative machinery.

### 4.1 Risk Measures: From Variance to Tail Risk

The choice of risk measure determines what the optimizer considers “risky” and therefore shapes portfolio composition in profound ways. Variance treats upside and downside symmetrically, penalizing favorable returns with the same weight as adverse ones. Tail risk measures concentrate attention on the extreme left tail of the return distribution, where losses are largest and most consequential. Drawdown measures target the investor’s lived experience of peak-to-trough declines, which often drives behavioral responses more powerfully than abstract distributional statistics. Any of these measures can serve interchangeably as the objective or constraint in mean-risk optimization, and this architectural uniformity allows the practitioner to explore the full risk-measure landscape without altering the optimization framework.

#### 4.1.1 Variance and Standard Deviation

Variance remains the foundational risk measure in portfolio theory, originating in the mean-variance framework. Portfolio variance is defined as

$$\sigma_P^2 = \mathbf{w}^\top \Sigma \mathbf{w}$$

where  $\mathbf{w}$  denotes the vector of portfolio weights and  $\Sigma$  the covariance matrix of asset returns. The portfolio standard deviation follows directly:

$$\sigma_P = \sqrt{\mathbf{w}^\top \Sigma \mathbf{w}}$$

Variance is symmetric: it penalizes upside deviations from the mean with precisely the same severity as downside deviations. This symmetry is both its analytical strength and its conceptual limitation. From an optimization standpoint, variance is convex and quadratic in the weight vector, admitting closed-form solutions under linear constraints and efficient numerical solution under more general constraint sets. Its analytical tractability has made it the default risk measure for decades, and its

properties are thoroughly understood.

When variance serves as the objective, the resulting optimization problem is a convex quadratic program, solvable with high reliability by standard solvers.

#### 4.1.2 Semi-Variance and Semi-Deviation

Semi-variance addresses the conceptual limitation of variance by restricting attention to returns that fall below the mean. It is defined as

$$\text{SemiVar}_P = \frac{1}{T} \sum_{t=1}^T [\min(r_{P,t} - \bar{r}_P, 0)]^2$$

where  $r_{P,t}$  denotes the portfolio return in period  $t$ ,  $\bar{r}_P$  the mean portfolio return, and  $T$  the number of observations. Semi-deviation is the square root of semi-variance. By zeroing out positive deviations, semi-variance captures only the downside component of return dispersion, better reflecting investor loss aversion. An asset that frequently produces large positive surprises but rarely declines will appear far less risky under semi-variance than under variance.

The semi-deviation follows as

$$\text{SemiDev}_P = \sqrt{\text{SemiVar}_P}$$

Semi-variance is not as analytically convenient as full variance, since it depends on the realized return path rather than solely on the covariance matrix. However, it remains a convex function of portfolio weights under standard conditions, and its optimization proceeds through sample-based formulations.

#### 4.1.3 Mean Absolute Deviation

The mean absolute deviation offers an alternative dispersion measure that is more robust to outliers than variance:

$$\text{MAD}_P = \frac{1}{T} \sum_{t=1}^T |r_{P,t} - \bar{r}_P|$$

Because MAD uses absolute deviations rather than squared deviations, extreme observations exert less influence on the risk estimate. This makes MAD particularly suitable when the return distribution contains occasional outliers that would disproportionately inflate variance. The MAD-based optimization can be formulated as a linear program through standard auxiliary variable techniques, which can be computationally advantageous for large portfolios.

#### 4.1.4 Value-at-Risk

Value-at-Risk at confidence level  $\alpha$  represents the loss threshold that is exceeded with probability  $\alpha$ :

$$\text{VaR}_\alpha(P) = -F_{r_P}^{-1}(\alpha)$$

where  $F_{r_P}^{-1}$  is the inverse cumulative distribution function (quantile function) of portfolio returns. For instance, at  $\alpha = 0.05$ , VaR is the loss level such that only five percent of return realizations are worse. VaR is intuitive and widely reported, but it possesses a critical theoretical deficiency: it is not a coherent risk measure. Specifically, VaR can violate subadditivity, meaning that the VaR of a combined portfolio can exceed the sum of the individual VaR contributions. This implies that diversification may appear to increase risk under VaR, a paradoxical and undesirable property.

VaR provides no information about the magnitude of losses beyond the threshold. Two portfolios with identical VaR may have vastly different tail behavior. Despite these limitations, VaR remains useful as a reporting metric and regulatory benchmark.

#### 4.1.5 Conditional Value-at-Risk

Conditional Value-at-Risk (also known as Expected Shortfall) remedies the deficiencies of VaR by averaging over the entire tail beyond the VaR threshold:

$$\text{CVaR}_\alpha(\mathbf{w}) = -\frac{1}{\alpha} \int_0^\alpha F_{\mathbf{w}^\top \mathbf{r}}^{-1}(p) dp$$

CVaR represents the expected loss conditional on the loss exceeding VaR. It is a coherent risk measure, satisfying all four axioms of coherence: monotonicity, translation invariance, positive homogeneity, and crucially, subadditivity. This last property ensures that diversification is always recognized as risk-reducing under CVaR. CVaR is convex in portfolio weights, making it amenable to efficient optimization. The sample-based formulation replaces the integral with an average over the worst  $\lfloor \alpha T \rfloor$  return observations, and the Rockafellar-Uryasev reformulation converts CVaR minimization into a linear program.

CVaR is more conservative than VaR because it accounts for the severity of tail losses, not merely their frequency. This makes it the preferred risk measure for portfolios where tail risk is a primary concern.

#### 4.1.6 Entropic Value-at-Risk

The Entropic Value-at-Risk provides an even tighter bound on tail risk than CVaR:

$$\text{EVaR}_\alpha(\mathbf{w}) = \inf_{z>0} \left\{ \frac{1}{z} \ln \left( \frac{M_{\mathbf{w}^\top \mathbf{r}}(z)}{\alpha} \right) \right\}$$

where  $M_{\mathbf{w}^\top \mathbf{r}}(z) = \mathbb{E}[e^{z \cdot \mathbf{w}^\top \mathbf{r}}]$  is the moment-generating function of portfolio returns. EVaR is

derived from the Chernoff bound on tail probabilities and dominates CVaR: for any portfolio,  $\text{EVaR}_\alpha \geq \text{CVaR}_\alpha$ . This makes EVaR a more conservative risk measure that is better suited for distributions with heavy tails, where the moment-generating function captures information about higher-order moments that CVaR does not fully exploit.

EVaR is coherent and convex, and its optimization can be handled through exponential cone programming. It is particularly valuable when the return distribution exhibits significant kurtosis or skewness.

#### 4.1.7 Worst Realization

The worst realization is the most conservative point-in-time risk measure, defined simply as the maximum single-period loss observed in the sample:

$$\text{WR}(\mathbf{w}) = \max_{t \in \{1, \dots, T\}} (-r_{P,t})$$

Minimizing worst realization produces a minimax portfolio that guards against the single worst historical outcome. This measure is extremely conservative and may produce overly defensive allocations in practice, but it serves as a useful bound and stress-test criterion.

#### 4.1.8 Maximum Drawdown

Maximum drawdown captures the largest peak-to-trough decline in portfolio value over the evaluation period:

$$\text{MDD} = \max_{t \in [0, T]} \left( \frac{\max_{s \in [0, t]} V_s - V_t}{\max_{s \in [0, t]} V_s} \right)$$

where  $V_t$  is the portfolio value at time  $t$ . Unlike point-in-time risk measures, maximum drawdown reflects the path-dependent experience of the investor. A portfolio may have low variance yet suffer severe drawdowns if negative returns cluster in time. Maximum drawdown is directly relevant to investor psychology, as prolonged declines often trigger premature liquidation and behavioral errors.

Maximum drawdown is not convex in general, but appropriate reformulations enable its use as an optimization objective.

#### 4.1.9 Average Drawdown

Average drawdown computes the mean of all drawdown episodes over the evaluation period, providing a more representative picture of typical drawdown behavior than the single worst case captured by maximum drawdown. A portfolio that experiences many moderate drawdowns may have a high average drawdown despite a moderate maximum drawdown, signaling persistent capital impairment. This measure is less sensitive to a single extreme event and therefore more stable as an optimization objective.



#### 4.1.10 Conditional Drawdown-at-Risk

Conditional Drawdown-at-Risk applies the CVaR concept to the distribution of drawdowns rather than returns:

$$\text{CDaR}_\alpha(\mathbf{w}) = -\frac{1}{\alpha} \int_0^\alpha F_{\text{DD}(\mathbf{w})}^{-1}(p) dp$$

where  $F_{\text{DD}(\mathbf{w})}$  is the cumulative distribution function of drawdowns for portfolio  $\mathbf{w}$ . CDaR represents the expected drawdown in the worst  $\alpha$  fraction of drawdown episodes. It is the drawdown analogue of CVaR and inherits its coherence properties within the drawdown domain. CDaR is particularly valuable for investors whose risk tolerance is framed in terms of peak-to-trough declines rather than return dispersion.

#### 4.1.11 Entropic Drawdown-at-Risk

Entropic Drawdown-at-Risk applies the EVaR tightening principle to the drawdown distribution, yielding a more conservative drawdown risk measure than CDaR. Just as EVaR dominates CVaR for return distributions, EDaR dominates CDaR for drawdown distributions, providing tighter tail bounds when drawdown episodes exhibit heavy-tailed behavior. EDaR is suitable for portfolios where protection against extreme drawdown scenarios is paramount.

#### 4.1.12 Ulcer Index

The Ulcer Index computes the root mean square of drawdowns, thereby penalizing both the depth and duration of drawdown episodes:

$$\text{UI} = \sqrt{\frac{1}{T} \sum_{t=1}^T D_t^2}$$

where  $D_t$  is the drawdown at time  $t$ . By squaring drawdowns before averaging, the Ulcer Index places greater weight on severe drawdowns while still accounting for persistent moderate drawdowns. It provides a single scalar summary of the entire drawdown experience.

#### 4.1.13 Risk Measure Selection Guidance

The choice of risk measure should reflect the investor's loss preferences, the distributional characteristics of the asset universe, and the computational requirements of the optimization. The following table summarizes the recommended mappings:

Investor Preference	Recommended Risk Measure	Rationale
Traditional mean-variance	Variance	Analytical tractability, well-studied properties

Investor Preference	Recommended Risk Measure	Rationale
Loss aversion (downside focus)	Semi-variance or CVaR	Penalizes only adverse outcomes
Tail risk management	CVaR or EVaR	Captures extreme loss behavior coherently
Drawdown sensitivity	CDaR or maximum drawdown	Targets peak-to-trough investor experience
Conservative tail protection	EVaR or EDaR	Tightest bounds on tail and drawdown risk
Robust to outliers	MAD	Less sensitive to extreme observations

No single risk measure is universally superior. In practice, the most informative approach is to optimize under multiple risk measures and compare the resulting allocations, identifying positions that are robust across measures and those that are highly sensitive to the risk definition.

## 4.2 Risk Budgeting and Equal Risk Contribution

### 4.2.1 Risk Contribution Framework

Risk budgeting decomposes total portfolio risk into contributions attributable to each asset, enabling the portfolio manager to specify how much risk each position should bear. The framework rests on Euler's theorem for homogeneous functions: if the risk measure  $\rho(\mathbf{w})$  is positively homogeneous of degree one in weights (as standard deviation is), then it decomposes exactly into asset-level contributions.

The marginal risk contribution of asset  $i$  measures the sensitivity of portfolio risk to a marginal increase in the weight of asset  $i$ :

$$\text{MRC}_i = \frac{\partial \sigma_P}{\partial w_i} = \frac{(\Sigma \mathbf{w})_i}{\sigma_P}$$

The component risk contribution multiplies the marginal contribution by the asset weight, yielding the absolute risk attributable to position  $i$ :

$$\text{RC}_i = w_i \cdot \text{MRC}_i = w_i \cdot \frac{(\Sigma \mathbf{w})_i}{\sigma_P}$$

The Euler decomposition guarantees that these contributions sum exactly to total portfolio risk:

$$\sum_{i=1}^N \text{RC}_i = \sigma_P$$

This decomposition provides a complete accounting of risk: no risk is “unattributed,” and the manager can assess whether the risk allocation across positions aligns with investment convictions.

#### 4.2.2 Equal Risk Contribution Portfolios

The equal risk contribution portfolio requires each asset to contribute identically to total portfolio risk:

$$w_i \cdot (\Sigma \mathbf{w})_i = \frac{\sigma_P^2}{N} \quad \text{for all } i$$

This system of  $N$  non-linear equations does not admit a closed-form solution and must be solved iteratively. The resulting allocation lies between equal weighting and minimum variance: it avoids the extreme concentration that minimum variance can produce while incorporating correlation structure that equal weighting ignores. Equal risk contribution portfolios are particularly attractive because they require no expected return estimates, relying solely on covariance information. This eliminates the largest source of estimation error in portfolio optimization.

The equal risk contribution approach embodies a principled agnosticism: absent strong views on expected returns, the most defensible allocation is one where no single asset dominates portfolio risk. In the risk budgeting framework, this corresponds to setting the budget vector to equal allocations across all assets.

#### 4.2.3 Custom Risk Budgets

When the investor holds differential convictions about asset attractiveness or risk tolerance, custom risk budgets allow explicit specification of the desired risk allocation. Given a budget vector  $\mathbf{b}$  with  $b_i > 0$  for all  $i$  and  $\sum_{i=1}^N b_i = 1$ , the custom risk budget portfolio satisfies:

$$w_i \cdot (\Sigma \mathbf{w})_i = b_i \cdot \sigma_P^2 \quad \text{for all } i$$

Assets assigned larger budgets receive greater weight (all else equal), while those with smaller budgets are constrained to contribute less risk. The budget vector translates qualitative views about asset desirability into a quantitative risk allocation target.

#### 4.2.4 Extension to Tail Risk Measures

The risk budgeting framework generalizes naturally beyond variance to any convex, positively homogeneous risk measure  $\rho(\mathbf{w})$ . The generalized risk contribution of asset  $i$  is defined via the gradient:

$$\text{RC}_i = w_i \cdot \frac{\partial \rho}{\partial w_i}$$

and the Euler decomposition continues to hold:  $\sum_i \text{RC}_i = \rho(\mathbf{w})$ .

CVaR-based risk budgeting accounts for fat tails and asymmetric return distributions, allocating risk based on each asset’s contribution to expected tail loss rather than to variance. CDaR-based risk budgeting targets drawdown contributions, directly addressing the investor’s experience of capital impairment. The risk budgeting optimizer accepts any of the risk measures discussed above, enabling seamless exploration of risk budgeting under alternative risk definitions.

### 4.3 Maximum Diversification

The maximum diversification portfolio maximizes the diversification ratio, defined as

$$\text{DR}(\mathbf{w}) = \frac{\mathbf{w}^\top \boldsymbol{\sigma}}{\sqrt{\mathbf{w}^\top \boldsymbol{\Sigma} \mathbf{w}}} = \frac{\sum_{i=1}^N w_i \sigma_i}{\sigma_P}$$

where  $\boldsymbol{\sigma}$  is the vector of individual asset volatilities and  $\sigma_P$  is the portfolio volatility. The numerator represents the weighted average of individual volatilities, which is the portfolio volatility that would prevail if all pairwise correlations were unity. The denominator is the actual portfolio volatility, which is lower due to imperfect correlation. The diversification ratio therefore quantifies the risk reduction achieved through diversification.

A diversification ratio of exactly 1.0 indicates zero diversification benefit: the portfolio behaves as though it holds a single asset. Higher ratios indicate greater diversification, with the maximum achievable ratio depending on the correlation structure of the asset universe. The maximum diversification portfolio is equivalent to the minimum variance portfolio when all assets are first standardized to unit volatility, revealing an elegant connection between the two approaches: maximum diversification seeks the allocation that extracts the greatest correlation-driven risk reduction, independent of individual asset volatility levels.

### 4.4 Distance Measures and Codependence

Hierarchical portfolio methods require a distance or codependence measure between assets to define the clustering structure. The choice of distance measure determines what notion of “similarity” drives the hierarchical decomposition, and different measures can produce substantially different dendrograms and therefore different allocations.

#### 4.4.1 Pearson Distance

The most widely used distance measure transforms the Pearson linear correlation coefficient into a metric:

$$d_{ij} = \sqrt{\frac{1}{2}(1 - \rho_{ij}^{\text{Pearson}})}$$

This transformation maps perfectly correlated assets ( $\rho = 1$ ) to zero distance and uncorrelated

assets ( $\rho = 0$ ) to distance  $1/\sqrt{2}$ . Perfectly negatively correlated assets ( $\rho = -1$ ) map to distance 1. Pearson distance captures linear dependence and is the standard choice for most applications.

#### 4.4.2 Kendall Distance

Kendall distance is derived from Kendall's tau rank correlation coefficient. Kendall's tau measures the concordance between two variables: the probability that two randomly selected observations exhibit the same ordering in both variables minus the probability they exhibit opposite orderings. The resulting distance is robust to non-linear monotonic relationships, as it depends only on the ranks of observations rather than their magnitudes. This makes it less sensitive to outliers and distributional assumptions than Pearson distance.

#### 4.4.3 Spearman Distance

Spearman distance is based on the Spearman rank correlation, which is simply the Pearson correlation computed on the ranks of the observations. It captures monotonic but not necessarily linear dependence, occupying a middle ground between Pearson's linearity assumption and the more general dependence captured by non-parametric measures. Spearman distance shares Kendall distance's robustness to outliers while being computationally simpler.

#### 4.4.4 Distance Correlation

Distance correlation, introduced in the statistical literature as a measure of dependence between random vectors, captures both linear and non-linear dependence. Its defining property is that distance correlation equals zero if and only if the two variables are statistically independent. This is a strictly stronger condition than zero Pearson correlation, which only implies independence under joint normality. Distance correlation therefore detects non-linear relationships that Pearson, Kendall, and Spearman measures may miss entirely.

This measure is particularly valuable in asset universes where non-linear dependence structures arise, such as during market stress when correlations exhibit threshold effects.

#### 4.4.5 Mutual Information

Mutual information is an information-theoretic measure that quantifies the total statistical dependence between two variables:

$$I(X; Y) = \int \int p(x, y) \ln \frac{p(x, y)}{p(x)p(y)} dx dy$$

Mutual information is zero if and only if the variables are independent, and it captures all forms of dependence: linear, non-linear, and higher-order. Unlike correlation-based measures, mutual information is not limited to monotonic relationships. Its estimation requires binning or kernel density estimation, introducing a degree of sensitivity to the estimation procedure. It provides the most general characterization of statistical dependence among the available measures.

## 4.5 Hierarchical Clustering

Agglomerative hierarchical clustering constructs a dendrogram, a tree-structured representation of asset similarity, by iteratively merging the closest pairs of assets or clusters. Beginning with  $N$  singleton clusters (one per asset), the algorithm proceeds as follows: at each step, the two clusters with the smallest inter-cluster distance are merged, and the distance matrix is updated. This process continues until all assets belong to a single cluster, producing a complete hierarchical decomposition. The choice of linkage method determines how inter-cluster distance is computed from pairwise asset distances:

- **Single linkage** uses the minimum distance between any pair of assets across two clusters, tending to produce elongated, chain-like clusters.
- **Complete linkage** uses the maximum distance, producing compact, spherical clusters but being sensitive to outliers.
- **Average linkage** uses the mean distance, offering a balance between the two extremes.
- **Ward linkage** minimizes the total within-cluster variance at each merge, producing the most compact and balanced clusters. It is the most commonly used linkage method in portfolio applications.
- **Weighted linkage** assigns equal weight to each cluster regardless of size when computing inter-cluster distances.

The resulting dendrogram reveals the natural grouping structure of the asset universe: assets that cluster together at low merge distances are highly similar (close in the chosen distance metric), while those that merge only at high distances are dissimilar. This hierarchical structure is exploited by hierarchical risk parity, hierarchical equal risk contribution, and nested clusters optimization to produce allocations that respect asset relationships.

## 4.6 Hierarchical Risk Parity

Hierarchical Risk Parity, introduced by Lopez de Prado (2016), addresses the fundamental instability of mean-variance optimization by replacing matrix inversion with a hierarchical allocation procedure. The algorithm proceeds in three steps.

**Step 1: Distance-based clustering.** Compute pairwise distances between all assets using the chosen distance measure (Pearson, Kendall, Spearman, distance correlation, or mutual information). Apply hierarchical clustering with the chosen linkage method to build a dendrogram encoding the asset similarity structure.

**Step 2: Quasi-diagonalization.** Reorder the rows and columns of the covariance matrix according to the leaf ordering of the dendrogram. This quasi-diagonalization places similar assets adjacent to one another, concentrating large covariance entries near the diagonal. The reordered matrix is not truly diagonal, but it is organized so that the hierarchical structure is geometrically apparent.

**Step 3: Recursive bisection.** Allocate risk top-down through the dendrogram. At the root, the full asset universe is split into the two clusters defined by the top-level merge. Capital is divided

between these clusters inversely proportional to their respective variances:

$$\alpha = 1 - \frac{\sigma_1^2}{\sigma_1^2 + \sigma_2^2}$$

The left cluster receives weight  $\alpha$  and the right cluster receives weight  $1 - \alpha$ . This bisection is applied recursively at each level of the dendrogram until individual assets are reached, at which point each asset's weight is the product of all the allocation fractions along its path from the root.

The key advantages of HRP are substantial. First, no matrix inversion is required at any stage, eliminating the numerical instability that plagues mean-variance optimization when the covariance matrix is ill-conditioned or nearly singular. Second, the algorithm handles singular covariance matrices naturally, which arise whenever the number of assets exceeds the number of return observations. Third, by respecting the hierarchical clustering structure, HRP produces allocations that are more stable over time: small perturbations in the covariance matrix cause small changes in the dendrogram and hence small changes in weights, in contrast to the erratic weight swings of unconstrained mean-variance optimization.

HRP accepts parameters for the distance measure, clustering method, and risk measure used in the bisection step. The risk measure need not be variance; any of the measures discussed above can be used, enabling CVaR-based or drawdown-based hierarchical risk parity.

#### 4.7 Hierarchical Equal Risk Contribution

Hierarchical Equal Risk Contribution (HERC) extends HRP by replacing the simple inverse-variance bisection with an equal risk contribution allocation within each cluster. At each node in the dendrogram, rather than dividing capital proportional to inverse variance, the algorithm solves for weights such that each sub-cluster contributes equally to the risk of the parent cluster. This combines the stability advantages of hierarchical methods (no matrix inversion, respect for clustering structure, robustness to covariance estimation error) with the risk-parity properties of equal risk contribution.

The result is a multi-level risk parity: each cluster at every level of the hierarchy, and each asset within its cluster, contributes equally to total risk at its respective level. This produces allocations that are both stable (from the hierarchical structure) and risk-balanced (from the equal contribution principle). The diversification is achieved simultaneously across the hierarchical levels, ensuring that risk concentration does not arise at any scale.

HERC accepts the same parameterization options as HRP: distance measure, clustering method, and risk measure. The choice of risk measure determines what “equal risk contribution” means at each level, whether equal variance contribution, equal CVaR contribution, or equal drawdown contribution, among other possibilities.

## 4.8 Nested Clusters Optimization

Nested Clusters Optimization (NCO) takes a different approach to exploiting hierarchical structure. Rather than allocating through recursive bisection, NCO uses the clustering to decompose a large optimization problem into smaller, better-conditioned sub-problems.

**Stage 1: Inner optimization.** For each cluster identified by the hierarchical clustering, an independent optimization is performed over the assets within that cluster. If a cluster contains  $n_k$  assets, the inner optimization operates on an  $n_k \times n_k$  covariance sub-matrix. Each cluster's optimization produces a set of intra-cluster weights, effectively reducing the cluster to a single composite asset with known return and risk characteristics.

**Stage 2: Outer optimization.** The composite assets from all  $K$  clusters are collected, and a second optimization is performed over these  $K$  composite assets to determine the inter-cluster allocation. The final weight of each individual asset is the product of its intra-cluster weight and its cluster's inter-cluster weight.

The dimensional reduction is the key benefit. Consider a universe of  $N = 500$  assets partitioned into  $K = 20$  clusters of approximately 25 assets each. The inner optimizations each handle  $25 \times 25$  covariance matrices, and the outer optimization handles a  $20 \times 20$  matrix. Both are well-conditioned even with limited return history, whereas direct optimization of the full  $500 \times 500$  matrix may be severely ill-conditioned. This decomposition dramatically improves the numerical stability and statistical reliability of the optimization.

A further advantage of NCO is modularity: the inner and outer stages can use entirely different optimization strategies. For instance, the inner stage might use minimum variance to produce concentrated intra-cluster allocations, while the outer stage uses risk parity across clusters. This flexibility enables the practitioner to tailor the optimization approach to the characteristics of each hierarchical level.

## 4.9 Regime-Driven Risk Adaptation

### 4.9.1 Markov-Driven Risk Measure Selection

The risk measures discussed above need not be applied statically across all market environments. The Hidden Markov Model framework formalized in the moment estimation chapter provides filtered regime probabilities that can drive dynamic selection of the risk measure most appropriate to current conditions.

Given  $S$  regimes with filtered probabilities  $p(z_t = s \mid \mathbf{r}_{1:t})$  and a regime-specific risk measure mapping  $\rho_s$ , the effective portfolio risk at time  $t$  is the probability-weighted combination:

$$\rho_t(\mathbf{w}) = \sum_{s=1}^S p(z_t = s \mid \mathbf{r}_{1:t}) \rho_s(\mathbf{w})$$

In a two-state model, the mapping might assign variance as the risk measure in the low-volatility



state and CVaR at the 95% confidence level in the high-volatility state:

$$\rho_t(\mathbf{w}) = p(z_t = 1 \mid \mathbf{r}_{1:t}) \cdot \sigma^2(\mathbf{w}) + p(z_t = 2 \mid \mathbf{r}_{1:t}) \cdot \text{CVaR}_{0.05}(\mathbf{w})$$

As the filtered probability of the crisis state increases, the effective risk measure transitions smoothly from variance toward CVaR, automatically increasing the portfolio's sensitivity to tail risk. The continuous nature of the filtered probabilities ensures that the transition is gradual rather than abrupt, preventing the portfolio instability that hard switching would produce.

When the Deep Markov Model framework is used instead of a discrete HMM, the continuous latent state  $\mathbf{z}_t$  can parameterize the risk measure continuously. A mapping  $\rho(\mathbf{w}; \mathbf{z}_t)$  that depends on the latent state enables risk sensitivity to vary along a continuum rather than switching between a finite set of predefined measures. In practice, this is achieved by defining a convex combination weight  $\lambda(\mathbf{z}_t) \in [0, 1]$  (output by a neural network) that interpolates between two anchor risk measures.

#### 4.9.2 Regime-Conditional Risk Budgets

Regime probabilities also inform the calibration of risk budgets. Denoting by  $\mathbf{b}_s$  the risk budget vector appropriate for regime  $s$ , the blended budget at time  $t$  is:

$$\mathbf{b}_t = \sum_{s=1}^S p(z_t = s \mid \mathbf{r}_{1:t}) \mathbf{b}_s$$

In expansionary regimes, the budget might allocate greater risk to cyclical sectors with strong earnings momentum. In contractionary regimes, risk shifts toward defensive sectors with stable cash flows and lower economic sensitivity. The smooth blending ensures that portfolio restructuring at regime boundaries is gradual, limiting unnecessary turnover.

For a two-state model with sector groups  $\{G_1, \dots, G_K\}$ , the regime-conditional budgets take the form:

$$b_{k,1} = \frac{w_k^{\text{expansion}}}{\sum_j w_j^{\text{expansion}}}, \quad b_{k,2} = \frac{w_k^{\text{contraction}}}{\sum_j w_j^{\text{contraction}}}$$

where  $w_k^{\text{expansion}}$  and  $w_k^{\text{contraction}}$  reflect the desired risk allocation to sector group  $k$  under each regime. The blended budget  $\mathbf{b}_t$  then governs the risk budgeting optimization at each rebalancing date, producing allocations that rotate sector exposures in response to changing macroeconomic conditions without requiring manual intervention.

#### 4.9.3 Integration with Hierarchical Methods

Regime probabilities integrate naturally with the hierarchical methods discussed above. At each rebalancing date, the regime-conditional covariance matrix  $\hat{\Sigma}_t$  enters the distance computation and

clustering step, producing a regime-adaptive dendrogram. Clusters that form during stress periods differ from those in calm periods, as correlations shift and factor structures change. By re-clustering at each rebalancing date using the current regime-conditional covariance, the hierarchical methods adapt their asset groupings to the prevailing market structure.

In HRP and HERC, the risk measure used in the bisection or equal-risk-contribution step can itself be regime-dependent: variance during calm regimes and CVaR during stress regimes. This produces allocations that are simultaneously hierarchically structured, risk-balanced, and regime-adaptive.

## 4.10 LLM-Augmented Risk and Diversification

### 4.10.1 Regime-Dependent Risk Measure Selection

The choice of risk measure need not be static. Different market regimes favor different risk measures, and large language models can provide the regime classification that drives dynamic risk measure selection. In calm, low-volatility environments, variance-based optimization may suffice, as return distributions are approximately symmetric and tail events are rare. When conditions deteriorate and tail risk becomes elevated (as indicated by widening credit spreads, rising implied volatility, or deteriorating economic indicators) switching to CVaR or EVaR captures the increased importance of extreme losses. In environments where drawdowns are the primary concern, perhaps due to leverage constraints or client sensitivity to capital impairment, CDaR or maximum drawdown becomes the appropriate objective.

An LLM processing current economic data, central bank communications, and market commentary can classify the prevailing regime along these dimensions and recommend the risk measure most appropriate for current conditions. This recommendation feeds into a dynamic risk measure selection policy: the optimization framework remains fixed, but the risk measure parameter rotates based on the LLM's regime assessment. The result is an adaptive portfolio construction process that selects the risk lens most relevant to the current environment.

### 4.10.2 Stress Scenario Design

Traditional stress testing relies on replaying historical crisis episodes, an approach that is inherently backward-looking and may miss novel risk configurations. Large language models can complement historical analysis by designing forward-looking stress scenarios derived from current conditions. The process involves analyzing the current constellation of risk factors (leverage levels across sectors, valuation extremes, geopolitical tensions, policy uncertainty, supply chain vulnerabilities) and constructing conditional scenarios that specify how these factors might interact adversely.

For example, an LLM might identify that the combination of elevated corporate leverage and tightening monetary policy creates a specific vulnerability in credit-sensitive sectors, and design a stress scenario in which credit spreads widen sharply while equity markets decline and interest rates remain elevated. This scenario may not correspond to any single historical episode but is nonetheless plausible and relevant. The resulting stress test complements purely historical scenario analysis by introducing forward-looking, context-specific risk assessment.

The stress scenarios can be translated into return shocks applied to the portfolio, enabling the optimizer to evaluate how different allocations perform under adverse conditions that the LLM identifies as currently relevant. This integration produces stress tests that evolve with the risk environment rather than remaining anchored to past crises.

#### 4.10.3 Risk Budget Calibration from Sector Outlook

Custom risk budgets require the practitioner to specify how much risk each asset or sector should bear. Large language models can inform this specification by processing sector-level research, earnings trends, and macroeconomic sensitivity analysis. Sectors with favorable outlooks (supported by earnings momentum, favorable policy conditions, or structural demand growth) receive larger risk budgets, allowing the optimizer to allocate greater risk to these areas. Sectors facing headwinds (margin compression, regulatory pressure, or cyclical vulnerability) receive smaller risk budgets, constraining their contribution to total portfolio risk.

This process translates qualitative sector analysis into quantitative risk budget vectors that the risk budgeting optimizer can consume directly. The LLM serves as the bridge between unstructured research and structured optimization inputs, converting narrative assessments into the numerical budget vector  $\mathbf{b}$  that parameterizes the risk budgeting problem. By updating these budgets as conditions evolve, the portfolio maintains an alignment between qualitative investment views and quantitative risk allocation that would be difficult to achieve through purely manual processes.

The combination of LLM-driven risk budgets with the full flexibility of the risk measure taxonomy creates a powerful adaptive framework: not only does the risk budget shift with the investment outlook, but the underlying risk measure can also rotate with the market regime, producing a portfolio that is responsive to both the level and the nature of risk in the current environment.

## 5 Portfolio Optimization and Robust Methods

With prior distributions constructed and risk measures selected, the optimization stage determines portfolio weights that best satisfy the investor's objectives subject to practical constraints. This chapter presents the optimization formulations, from classical mean-risk to distributionally robust methods, along with the constraint specifications, ensemble approaches, and naive baselines that form the complete optimization toolkit.

### 5.1 Objective Functions

The choice of objective function determines how the optimizer trades off return against risk. Each formulation identifies a distinct point or region on the efficient frontier, and the investor's mandate dictates which is appropriate.

#### 5.1.1 Minimize Risk

The simplest formulation seeks the portfolio with the smallest possible risk exposure:

$$\min_{\mathbf{w}} \rho(\mathbf{w})$$

where  $\rho$  denotes any convex risk measure (variance, CVaR, maximum drawdown, or any of the coherent measures discussed in prior chapters). This objective produces the global minimum risk point on the efficient frontier. Because no expected return estimate enters the formulation, the optimizer avoids the well-documented estimation error in  $\boldsymbol{\mu}$  entirely.

#### 5.1.2 Maximize Return

When the investor has a fixed risk budget, the optimizer seeks the highest expected return that remains within that budget:

$$\max_{\mathbf{w}} \mathbf{w}^\top \boldsymbol{\mu} \quad \text{subject to} \quad \rho(\mathbf{w}) \leq \rho_{\max}$$

The constraint  $\rho_{\max}$  represents the maximum tolerable risk level. This formulation traces the upper boundary of the feasible set for a given risk threshold.

#### 5.1.3 Maximize Utility

The utility-based formulation balances expected return against risk through a risk aversion parameter  $\lambda$ :

$$\max_{\mathbf{w}} \mathbf{w}^\top \boldsymbol{\mu} - \frac{\lambda}{2} \rho(\mathbf{w})$$

The parameter  $\lambda > 0$  governs the tradeoff: larger values penalize risk more heavily, producing more conservative allocations. When  $\rho(\mathbf{w}) = \mathbf{w}^\top \boldsymbol{\Sigma} \mathbf{w}$ , this reduces to the classical Markowitz quadratic utility.

#### 5.1.4 Maximize Ratio

The ratio objective identifies the tangency portfolio, the allocation with the highest reward per unit of risk:

$$\max_{\mathbf{w}} \frac{\mathbf{w}^\top \boldsymbol{\mu}}{\rho(\mathbf{w})}$$

When  $\rho$  is the portfolio standard deviation, this yields the maximum Sharpe ratio portfolio. When  $\rho$  is the conditional value-at-risk, it yields the maximum CVaR ratio portfolio, and similarly for any other risk measure. The tangency portfolio is the point at which a ray from the origin is tangent to the efficient frontier.

### 5.2 Constraints

Optimization without constraints produces mathematically elegant but practically infeasible portfolios. The constraint apparatus translates real-world investment mandates into mathematical restrictions on the weight vector  $\mathbf{w}$ .

#### 5.2.1 Position Limits

Box constraints bound each individual weight:

$$w_i^{\min} \leq w_i \leq w_i^{\max} \quad \forall i = 1, \dots, N$$

The long-only constraint sets  $w_i^{\min} = 0$ , prohibiting short positions. Maximum position limits  $w_i^{\max}$  (typically in the range 0.05 to 0.10 for diversified strategies) prevent excessive concentration in any single asset.

#### 5.2.2 Budget Constraint

The budget constraint governs the total capital deployment:

$$\sum_{i=1}^N w_i = 1$$

This enforces full investment. A relaxed variant permits partial investment:

$$\sum_{i=1}^N w_i \leq B$$

where  $B \leq 1$  allows a cash position.

### 5.2.3 Group and Sector Constraints

Group constraints impose bounds on the aggregate weight allocated to predefined asset groups (sectors, geographies, asset classes):

$$g_s^{\min} \leq \sum_{i \in G_s} w_i \leq g_s^{\max} \quad \text{for each group } s$$

where  $G_s$  denotes the set of assets belonging to group  $s$ . Active weight bounds restrict deviations from a benchmark allocation:

$$\Delta_s^{\min} \leq \sum_{i \in G_s} (w_i - w_{B,i}) \leq \Delta_s^{\max}$$

where  $w_{B,i}$  is the benchmark weight of asset  $i$ .

### 5.2.4 Transaction Costs and Fees

Practical portfolio management incurs costs that must be internalized within the optimization. Transaction costs penalize turnover relative to the current portfolio  $\mathbf{w}_0$ :

$$\sum_{i=1}^N c_i |w_i - w_{i,0}|$$

where  $c_i$  is the per-asset transaction cost (bid-ask spread plus commissions). Management fees impose a continuous drag proportional to position size:

$$\sum_{i=1}^N f_i |w_i|$$

where  $f_i$  captures the expense ratio for ETFs or funds held within the portfolio. Both terms are integrated directly into the objective function, ensuring that the optimizer accounts for implementation costs when determining weights.

### 5.2.5 Regularization

Regularization penalties improve out-of-sample stability by penalizing extreme or fragmented allocations.

The  $L_1$  penalty (lasso regularization) encourages sparsity by driving small positions to exactly zero:

$$\kappa_1 \|\mathbf{w}\|_1 = \kappa_1 \sum_{i=1}^N |w_i|$$

The  $L_2$  penalty (ridge regularization) shrinks extreme weights toward uniformity:

$$\kappa_2 \|\mathbf{w}\|_2^2 = \kappa_2 \sum_{i=1}^N w_i^2$$

The  $L_2$  penalty is mathematically equivalent to adding  $\kappa_2 \mathbf{I}$  to the covariance matrix, which improves its conditioning and reduces sensitivity to estimation noise. Typical values lie in the range  $\kappa_2 \in [0.001, 0.1]$ .

### 5.2.6 Custom Linear Constraints

For investment mandates that transcend the standard constraint types, the framework supports general linear inequality constraints:

$$\mathbf{A}\mathbf{w} \leq \mathbf{b}$$

where  $\mathbf{A} \in \mathbb{R}^{m \times N}$  and  $\mathbf{b} \in \mathbb{R}^m$  encode  $m$  linear restrictions. This general form subsumes turnover limits, tracking error bounds, portfolio beta constraints, and ESG score requirements as special cases.

## 5.3 Robust Optimization Under Parameter Uncertainty

Classical mean-variance optimization treats  $\hat{\boldsymbol{\mu}}$  and  $\hat{\boldsymbol{\Sigma}}$  as known with certainty, yet these are merely estimates subject to substantial sampling error. Robust optimization acknowledges this uncertainty by optimizing against the worst case within a plausible set of parameter values.

### 5.3.1 Mean Uncertainty Sets

The robust counterpart replaces the point estimate  $\hat{\boldsymbol{\mu}}$  with an uncertainty set  $\mathcal{U}_\mu$ :

$$\max_{\mathbf{w}} \min_{\boldsymbol{\mu} \in \mathcal{U}_\mu} \mathbf{w}^\top \boldsymbol{\mu} - \frac{\lambda}{2} \mathbf{w}^\top \boldsymbol{\Sigma} \mathbf{w}$$

The inner minimization finds the least favorable expected return vector within the uncertainty set, while the outer maximization selects weights that perform best under this adversarial scenario.

The ellipsoidal uncertainty set takes the form:

$$\mathcal{U}_\mu = \left\{ \boldsymbol{\mu} : (\boldsymbol{\mu} - \hat{\boldsymbol{\mu}})^\top \mathbf{S}_\mu^{-1} (\boldsymbol{\mu} - \hat{\boldsymbol{\mu}}) \leq \kappa^2 \right\}$$

where  $\mathbf{S}_\mu$  is the covariance of the estimator  $\hat{\boldsymbol{\mu}}$  and  $\kappa$  is calibrated from a confidence level (typically 90% to 95%). The parameter  $\kappa$  controls the size of the uncertainty region: larger values produce more conservative portfolios that hedge against greater estimation error.

Bootstrap-based approaches construct empirical confidence regions by resampling with replacement from the historical return series.

### 5.3.2 Covariance Uncertainty Sets

Analogously, the covariance matrix is subject to estimation error, particularly in the off-diagonal elements that govern asset co-movements. Robust covariance uncertainty sets protect against misestimation of the risk structure by considering a neighborhood of plausible covariance matrices around  $\hat{\boldsymbol{\Sigma}}$ .

Bootstrap-based sampling constructs confidence regions for  $\boldsymbol{\Sigma}$  by repeatedly resampling historical returns and computing the sample covariance for each resample. The resulting distribution of covariance estimates delineates the uncertainty region.

## 5.4 Distributionally Robust CVaR

Distributionally robust optimization extends the uncertainty set concept from parameters to the entire return distribution. Rather than assuming a fixed distribution  $P_0$  (typically the empirical distribution of historical returns), the optimizer considers all distributions within a Wasserstein ball centered on  $P_0$ :

$$\min_{\mathbf{w}} \sup_{P \in \mathcal{B}_\epsilon(P_0)} \text{CVaR}_\alpha^P(\mathbf{w})$$

where  $\mathcal{B}_\epsilon(P_0)$  denotes the set of all probability distributions whose Wasserstein distance from the empirical distribution  $P_0$  does not exceed  $\epsilon$ :

$$\mathcal{B}_\epsilon(P_0) = \{P : W(P, P_0) \leq \epsilon\}$$

The radius  $\epsilon$  (typically in the range 0.01 to 0.05) controls the degree of conservatism. A larger  $\epsilon$  enlarges the ambiguity set, producing more defensive portfolios that are robust to greater distributional perturbations. This formulation provides protection against regime changes, fat tails, and model misspecification without requiring parametric distributional assumptions.

## 5.5 Synthetic Data and Stress Testing

### 5.5.1 Vine Copula Framework

Historical return samples are inherently limited in size and may fail to capture extreme events, future regime shifts, or tail dependencies that are plausible but unobserved. Vine copula models



address this limitation by fitting flexible multivariate distributions that faithfully reproduce the statistical structure of asset returns.

The vine copula decomposition separates the multivariate distribution into three components:

1. **Marginal distributions** for each asset, fitted individually using flexible parametric families such as Student- $t$ , Johnson SU, or Normal Inverse Gaussian distributions.
2. **Bivariate copulas** that capture pairwise dependence structures, selected from families including Gaussian, Student- $t$ , Clayton, Gumbel, and Joe copulas. Each pair may use a different copula family, accommodating asymmetric and tail-dependent relationships.
3. **Vine structure** that decomposes the full  $N$ -dimensional dependence into a sequence of conditional bivariate relationships, organized as a tree structure where each level conditions on variables from previous levels.

This decomposition is both parsimonious and expressive: it requires only  $\binom{N}{2}$  bivariate copula selections rather than direct specification of the full  $N$ -dimensional distribution.

### 5.5.2 Scenario Generation

Given a fitted vine copula model, one can generate thousands of synthetic return scenarios that preserve the statistical properties of the historical data (including marginal shapes, pairwise dependence, tail dependence, and asymmetric co-movements) while extending well beyond the historical sample. These synthetic scenarios serve as inputs to scenario-based optimizers, enriching the information set available for portfolio construction.

### 5.5.3 Conditional Stress Testing

Conditional simulation fixes one or more variables at specified stress levels and generates consistent scenarios for the remaining assets, conditioned on the stressed values:

$$\mathbf{r}_{\text{synthetic}} \sim F(\mathbf{r} \mid r_{\text{market}} = -0.30)$$

This produces scenarios where, for example, the broad market has declined by 30%, and all other asset returns are drawn from their conditional distribution given this market shock. The vine copula structure makes such conditional sampling tractable, as the conditional distributions decompose naturally along the vine tree.

Applications include factor crash scenarios (equity market drawdown, interest rate spikes), sector-specific stress tests, and correlation regime shifts where historical co-movements break down.

## 5.6 Benchmark Tracking

Enhanced index strategies and constrained active mandates require portfolios that remain close to a benchmark while seeking modest outperformance. The benchmark tracking formulation minimizes tracking error:

$$\min_{\mathbf{w}} \text{TE}(\mathbf{w}, \mathbf{w}_B) \quad \text{subject to} \quad \text{TE} \leq \text{TE}_{\text{target}}$$

where the tracking error is defined as:

$$\text{TE} = \sqrt{\text{Var}(r_P - r_B)} = \sqrt{(\mathbf{w} - \mathbf{w}_B)^\top \boldsymbol{\Sigma} (\mathbf{w} - \mathbf{w}_B)}$$

and  $r_P$ ,  $r_B$  denote the portfolio and benchmark returns respectively. The target tracking error  $\text{TE}_{\text{target}}$  constrains how far the portfolio may deviate from the benchmark in risk space, balancing the desire for active returns against mandate compliance.

## 5.7 Naive Allocation Methods

Naive allocation methods require minimal or no parameter estimation and serve as essential baselines against which more sophisticated optimizers must demonstrate improvement.

### 5.7.1 Equal Weighted

The equal-weighted portfolio assigns identical weight to each asset:

$$w_i = \frac{1}{N} \quad \forall i = 1, \dots, N$$

This allocation requires no estimation whatsoever (neither expected returns nor covariances) and is therefore immune to estimation error. Despite its simplicity, the equal-weighted portfolio has proven to be a surprisingly competitive baseline, often outperforming mean-variance optimized portfolios on an out-of-sample basis when the number of assets is moderate and estimation windows are short.

### 5.7.2 Inverse Volatility

The inverse volatility portfolio scales each position inversely to its estimated volatility:

$$w_i = \frac{1/\sigma_i}{\sum_{j=1}^N 1/\sigma_j}$$

This allocation assigns larger weights to less volatile assets, implementing a rudimentary form of risk-based allocation. It is most appropriate when assets have similar Sharpe ratios, in which case equalizing risk contributions approximates the tangency portfolio. Since only marginal volatilities (not the full covariance matrix) are required, estimation error is substantially reduced relative to mean-variance optimization.

## 5.8 Ensemble Optimization

Model uncertainty is pervasive in portfolio optimization: different formulations, risk measures, and estimation methods yield different weight vectors, and no single approach dominates across all market conditions. Ensemble methods address this by combining multiple optimization strategies to reduce model risk.

Stacking, the primary ensemble approach, proceeds in three stages:

1. **Independent optimization.** A collection of diverse optimizers is run independently, each embodying different modeling assumptions. For example, one may use mean-risk optimization with Sharpe ratio maximization, another hierarchical risk parity, and a third risk budgeting with equal risk contributions.
2. **Sub-portfolio construction.** Each optimizer's output weight vector defines a sub-portfolio. These sub-portfolios span a range of risk-return characteristics reflecting the diversity of underlying assumptions.
3. **Meta-optimization.** A final meta-optimizer allocates capital across the sub-portfolios, treating each as a single composite asset. This meta-level allocation diversifies across model assumptions: if mean-variance overestimates expected returns while hierarchical risk parity is overly conservative, the stacking framework blends both, reducing the impact of any single model's errors.

## 6 Validation, Model Selection, and Production Pipeline

**A portfolio optimization pipeline is only as credible as its validation methodology.** Out-of-sample testing, cross-validation adapted to financial time series, and systematic hyperparameter tuning determine whether an optimized portfolio captures genuine risk-return structure or merely overfits historical noise. This chapter presents the validation and model selection tools, establishes the pipeline architecture that composes all preceding stages into a single estimation chain, and addresses the rebalancing frameworks that govern production deployment.

### 6.1 Walk-Forward Backtesting

**Walk-forward backtesting provides the most realistic out-of-sample evaluation** by strictly separating training and test periods in a manner that respects the causal arrow of time. No future information contaminates any training window, and the resulting multi-period portfolio faithfully simulates the experience of an investor who retrains and redeploys the strategy at regular intervals.

#### 6.1.1 Procedure

The full sample of  $T$  return observations is partitioned into successive non-overlapping segments. At each step  $k$ :

1. **Train** the optimization pipeline on the window  $[t_k, t_k + T_{\text{train}} - 1]$ .
2. **Predict** portfolio weights for the subsequent test window  $[t_k + T_{\text{train}}, t_k + T_{\text{train}} + T_{\text{test}} - 1]$ .
3. **Advance** the origin by  $T_{\text{test}}$  observations and repeat.

Typical calibrations set  $T_{\text{train}} = 252$  trading days (one calendar year) and  $T_{\text{test}} = 21$ –63 trading days (one to three months). The concatenation of all out-of-sample test segments yields a single backtest path in which every observation is genuinely out-of-sample.

#### 6.1.2 Rolling Versus Expanding Windows

Two windowing conventions exist, each with distinct statistical properties:

**Rolling windows** fix  $T_{\text{train}}$  so that the oldest observations are discarded as new data arrives:

$$\mathcal{W}_k^{\text{roll}} = \{t : t_k \leq t < t_k + T_{\text{train}}\}, \quad |\mathcal{W}_k^{\text{roll}}| = T_{\text{train}} \quad \forall k$$

Rolling windows adapt to regime changes because stale data from prior regimes is forgotten. The cost is discarding potentially useful history during stable periods, and the fixed sample size limits the precision of covariance estimates.

**Expanding windows** grow  $T_{\text{train}}$  as the walk-forward progresses:

$$\mathcal{W}_k^{\text{exp}} = \{t : t_0 \leq t < t_k + T_{\text{train}}\}, \quad |\mathcal{W}_k^{\text{exp}}| = T_{\text{train}} + k \cdot T_{\text{test}}$$

Expanding windows yield more stable parameter estimates (particularly for covariance matrices in high-dimensional settings) but adapt more slowly to structural breaks. The choice between the two conventions depends on the assumed stationarity of the return-generating process.

### 6.1.3 Implementation

Walk-forward backtesting is implemented as a temporal cross-validator with explicit test and training window size parameters. Applying this splitter to the full pipeline returns the concatenated out-of-sample portfolio for performance evaluation.

## 6.2 Combinatorial Purged Cross-Validation

**Walk-forward backtesting produces a single backtest path**, and a single path may reflect fortunate or unfortunate timing rather than genuine strategy quality. Combinatorial Purged Cross-Validation (CPCV) addresses this limitation by generating a population of backtest paths from the same historical sample, enabling statistical significance testing of portfolio performance.

### 6.2.1 Construction

The procedure operates as follows:

1. **Divide** the full sample into  $N_{\text{folds}}$  non-overlapping temporal blocks of approximately equal length.
2. **Select**  $N_{\text{test}}$  blocks for testing; the remaining  $N_{\text{folds}} - N_{\text{test}}$  blocks serve as the training set.
3. **Purge** observations near test-set boundaries to prevent information leakage. If the estimation window for any feature (e.g., a trailing volatility calculation) extends across the train-test boundary, the contaminated training observations are removed. The purge threshold controls the number of observations excised on each side of the boundary.
4. **Embargo** observations immediately following each test block to avoid autocorrelation contamination. If returns exhibit serial dependence over  $h$  lags, the embargo period should span at least  $h$  observations.
5. **Enumerate** all  $\binom{N_{\text{folds}}}{N_{\text{test}}}$  combinations, training and testing the pipeline on each.

### 6.2.2 Statistical Output

Each combination produces an out-of-sample portfolio segment. Reassembling the test segments across combinations generates multiple complete backtest paths. The resulting population of performance metrics (Sharpe ratios, maximum drawdowns, cumulative returns) permits distributional analysis:

$$\hat{p}(\text{SR} > 0) = \frac{1}{C} \sum_{c=1}^C \mathbf{1}\{\text{SR}_c > 0\}, \quad C = \binom{N_{\text{folds}}}{N_{\text{test}}}$$

where  $\text{SR}_c$  is the Sharpe ratio of the  $c$ -th backtest path. A strategy with  $\hat{p}(\text{SR} > 0) > 0.95$  provides

substantially stronger evidence of genuine skill than a single walk-forward backtest with a positive Sharpe ratio.

Summary statistics of interest include the mean Sharpe ratio across paths, the probability of positive excess returns, the distribution of maximum drawdowns, and the dispersion of terminal wealth outcomes.

### 6.2.3 Implementation

CPCV is implemented as a temporal cross-validator with parameters for the number of folds, the number of test folds, and the purge and embargo thresholds. Larger fold counts increase the number of combinations (and thus statistical power) but reduce the size of each training set. Typical configurations use  $N_{\text{folds}} \in [6, 10]$  and  $N_{\text{test}} = 2$ .

## 6.3 Multiple Randomized Cross-Validation

**Multiple Randomized Cross-Validation extends CPCV by introducing asset subsampling.** Each trial randomly selects both a temporal window and a subset of assets from the full universe. This dual randomization tests whether the strategy’s performance is robust to both temporal variation and asset composition.

The rationale is straightforward: a strategy that performs well only on a specific subset of assets or a specific historical window is more likely to be the product of data mining than one that performs consistently across many randomly drawn subsets. Multiple Randomized Cross-Validation provides the strongest available evidence that a portfolio optimization pipeline captures genuine predictive structure rather than spurious patterns.

This method combines temporal splitting with random asset selection across a configurable number of trials.

## 6.4 Performance Scoring

### 6.4.1 Built-In Ratio Measures

**Ratio measures quantify the reward earned per unit of risk,** differing in their definition of risk. The principal measures available as built-in scoring functions are:

**Sharpe Ratio,** return per unit total risk:

$$\text{SR} = \frac{\bar{r}_p - r_f}{\sigma_p}$$

where  $\bar{r}_p$  is the mean portfolio return,  $r_f$  is the risk-free rate, and  $\sigma_p$  is the standard deviation of portfolio returns. The Sharpe ratio treats upside and downside volatility symmetrically, penalizing strategies with large positive outliers.

**Sortino Ratio,** return per unit downside risk:

$$\text{Sortino} = \frac{\bar{r}_p - r_f}{\sigma_{\text{downside}}}, \quad \sigma_{\text{downside}} = \sqrt{\frac{1}{T} \sum_{t=1}^T \min(r_{p,t} - r_f, 0)^2}$$

The Sortino ratio penalizes only negative deviations, making it more appropriate for strategies with asymmetric return distributions.

**Calmar Ratio**, return per unit drawdown:

$$\text{Calmar} = \frac{\text{Annualized Return}}{\text{Max Drawdown}}$$

The Calmar ratio captures the worst-case capital erosion experience, relevant for strategies where investors are particularly sensitive to peak-to-trough losses.

**CVaR Ratio**, return per unit tail risk:

$$\text{CVaR Ratio} = \frac{\bar{r}_p - r_f}{\text{CVaR}_{95\%}}$$

where  $\text{CVaR}_{95\%} = \mathbb{E}[-r_p \mid -r_p \geq \text{VaR}_{95\%}]$  is the expected loss conditional on exceeding the 95th percentile loss. The CVaR ratio focuses on tail behavior, penalizing strategies that achieve attractive average returns at the cost of catastrophic downside events.

**Information Ratio**, active return per unit active risk:

$$\text{IR} = \frac{\bar{r}_p - \bar{r}_B}{\text{TE}}, \quad \text{TE} = \sqrt{\text{Var}(r_p - r_B)}$$

where  $\bar{r}_B$  is the mean benchmark return and TE is tracking error. The Information Ratio measures the efficiency of active bets relative to a benchmark, and is the primary performance metric for benchmark-aware strategies.

These ratio measures serve as scoring functions for model selection and hyperparameter tuning.

#### 6.4.2 Custom Scoring Functions

**When no single ratio measure captures the desired objective**, custom scoring functions combine multiple performance dimensions. For example, a score that rewards risk-adjusted returns while penalizing drawdowns:

$$\text{Score}(P) = \text{SR}(P) - 2 \cdot |\text{MaxDD}(P)|$$

The scoring function determines which strategy configuration “wins” during hyperparameter selection. The choice of scoring function therefore encodes the investor’s preferences over the full distribution

of portfolio outcomes, not merely its first two moments. Any callable that accepts a portfolio and returns a scalar can serve as a custom scoring function.

## 6.5 Hyperparameter Tuning

### 6.5.1 Grid Search

**Grid search exhaustively evaluates all parameter combinations** over a discrete parameter grid  $\Theta = \Theta_1 \times \Theta_2 \times \cdots \times \Theta_d$ , selecting the configuration that maximizes average cross-validated performance:

$$\hat{\theta} = \arg \max_{\theta \in \Theta} \frac{1}{K} \sum_{k=1}^K \text{Score}(\text{Model}(\theta), \text{Test}_k)$$

where  $K$  is the number of cross-validation folds and  $\text{Score}(\cdot)$  is the chosen performance measure. The critical requirement for financial applications is that **cross-validation must respect temporal ordering**: shuffling financial time-series data destroys the autocorrelation structure and produces optimistically biased performance estimates. The cross-validation splitter must therefore be a temporal method such as walk-forward or CPCV, never a random-shuffle approach.

### 6.5.2 Randomized Search

**Randomized search samples configurations from specified distributions** rather than exhaustively enumerating a grid. For each parameter, the practitioner specifies a probability distribution (uniform, log-uniform, discrete) from which candidates are drawn. This approach is more efficient when the parameter space is large and some parameters influence performance far more than others, a phenomenon known as low effective dimensionality.

The theoretical justification is that for  $d$  parameters with only  $d_{\text{eff}} \ll d$  parameters materially affecting performance, randomized search with  $n$  samples covers the important dimensions as effectively as a grid search with  $n^{d/d_{\text{eff}}}$  points.

### 6.5.3 Nested Parameter Tuning

**Hierarchical parameter specification enables joint optimization** of the entire pipeline, from covariance estimation method through risk measure selection to constraint calibration, in a single search. Nested parameter paths target parameters at arbitrary depth within the estimation chain.

For example, one can target the shrinkage intensity of the mean estimator nested within the prior estimator nested within the optimizer. This hierarchical addressing enables the search to simultaneously evaluate:

- Covariance shrinkage intensity at the covariance estimation level
- Expected return estimator type and its regularization parameters
- Risk measure selection at the optimizer level



- Constraint parameters (maximum weight, sector bounds) at the optimization level

The result is a principled search over the full configuration space rather than sequential, greedy tuning of individual components.

## 6.6 Pipeline Architecture

### 6.6.1 Pipeline Composition

**All stages of the portfolio construction process compose into a single pipeline object**, inheriting a uniform estimation interface. A typical pipeline chains:

1. **Pre-selection transformers:** retain assets with full history, remove numerically degenerate assets, eliminate redundant assets above a correlation threshold, retain top-ranked assets by a chosen criterion.
2. **Final estimator:** a mean-risk optimizer, hierarchical risk parity, risk budgeting, or any optimization estimator that accepts a return matrix and produces portfolio weights.

The pipeline propagates data through each stage sequentially. Fitting the pipeline fits each transformer in order, then fits the final estimator on the transformed data. Predicting returns the portfolio weights produced by the final estimator after all transformations.

The architectural consequence is profound: **the entire pipeline is treated as a single estimator** for cross-validation and hyperparameter tuning purposes. This means that pre-selection is performed within each cross-validation fold, preventing the subtle but pernicious form of data leakage that arises when pre-selection is applied to the full dataset before splitting.

### 6.6.2 Factor Returns as Metadata

**Factor model priors require factor returns alongside asset returns.** When the prior estimator is a factor model (e.g., a Fama-French three-factor model), the pipeline needs both the  $N$ -dimensional asset return matrix  $\mathbf{X}$  and the  $K$ -dimensional factor return matrix  $\mathbf{F}$ .

Factor returns are passed as auxiliary data alongside the asset return matrix. The pipeline propagates this auxiliary data through all stages that require it, ensuring that factor-based prior estimation receives the correct factor return data at each cross-validation fold.

### 6.6.3 Metadata Routing

**Some estimators require additional data beyond the return matrix  $\mathbf{X}$  and factor returns.** Implied covariance estimation needs implied volatility surfaces; benchmark-tracking optimization needs benchmark weights; transaction-cost-aware rebalancing needs current portfolio positions. These auxiliary data objects do not fit naturally into a two-input paradigm.

Metadata routing solves this problem by enabling arbitrary named data to flow through the pipeline to the specific estimators that consume it. Each estimator declares the metadata it requires, and the pipeline infrastructure ensures that the correct data reaches the correct stage.

The metadata routing mechanism generalizes the pipeline architecture to accommodate the full range of information required by sophisticated portfolio optimization, from market data through implied parameters to benchmark specifications, without breaking the clean estimator interface.

## 6.7 Rebalancing Frameworks

### 6.7.1 Calendar-Based Rebalancing

**Calendar-based rebalancing triggers portfolio reconstruction at fixed intervals**, regardless of how far the portfolio has drifted from target weights. The principal frequencies are:

- **Monthly** ( $T_{\text{rebal}} = 21$  trading days): highest turnover, appropriate for momentum-driven strategies where signal decay is rapid and timely rebalancing captures the bulk of the factor premium.
- **Quarterly** ( $T_{\text{rebal}} = 63$  trading days): balanced between signal freshness and transaction cost mitigation, widely adopted for multi-factor strategies.
- **Semiannual** ( $T_{\text{rebal}} = 126$  trading days): suitable for strategies with slower-decaying signals such as value or quality.
- **Annual** ( $T_{\text{rebal}} = 252$  trading days): lowest turnover, appropriate for strategic asset allocation or buy-and-hold tilts.

The choice of rebalancing frequency reflects the tension between **signal decay** (which favors frequent rebalancing) and **transaction costs** (which penalize it). The optimal frequency minimizes net-of-cost performance degradation, which depends on the specific factor structure, universe liquidity, and cost environment.

### 6.7.2 Threshold-Based Rebalancing

**Threshold-based rebalancing triggers trades only when portfolio drift exceeds specified limits**, avoiding unnecessary turnover during periods of relative stability. Two threshold conventions exist:

**Absolute threshold:** rebalance asset  $i$  when

$$|w_i - w_{\text{target},i}| > \Delta_{\text{abs}}$$

where  $\Delta_{\text{abs}}$  typically ranges from 0.03 to 0.05 (3–5 percentage points of portfolio weight).

**Relative threshold:** rebalance asset  $i$  when

$$\frac{|w_i - w_{\text{target},i}|}{w_{\text{target},i}} > \Delta_{\text{rel}}$$

where  $\Delta_{\text{rel}}$  typically ranges from 0.20 to 0.25 (20–25% deviation from target weight). Relative

thresholds are more appropriate for portfolios with heterogeneous position sizes, as a 3-percentage-point drift is negligible for a 30% position but transformative for a 3% position.

**Hybrid approaches** check thresholds at regular calendar intervals but only execute trades when at least one threshold is breached. This combines the discipline of calendar-based review with the cost efficiency of threshold-based execution.

### 6.7.3 Transaction Cost Integration

**Rebalancing optimization incorporates transaction costs directly into the objective function**, trading off expected portfolio improvement against implementation costs:

$$\max_{\mathbf{w}} \mathbf{w}^\top \boldsymbol{\mu} - \frac{\lambda}{2} \mathbf{w}^\top \boldsymbol{\Sigma} \mathbf{w} - \sum_{i=1}^N c_i |w_i - w_{i,0}|$$

where  $c_i$  is the per-unit transaction cost for asset  $i$  (encompassing commissions, bid-ask spread, and market impact) and  $w_{i,0}$  is the current portfolio weight. The  $\ell_1$  penalty on weight changes naturally induces sparsity in the trade vector: small improvements are not worth the cost of execution, so the optimizer leaves near-optimal positions unchanged.

The net return after rebalancing is:

$$\text{Net Return}_t = \text{Gross Return}_t - \sum_{i=1}^N c_i |w_{i,t} - w_{i,t-1}|$$

Realistic backtesting must account for these costs; a strategy that appears profitable gross of costs may be unprofitable net of costs, particularly for high-turnover strategies in less liquid markets.